

6/4/02
FD



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 648.40: On a convergence Run for QA. I was concerned that it was very weak!
That is: Error related only to past, known, data and we couldn't tell about the future
because the list of Q's was not known.

Well, if \exists a ^{truth} model that was generating the answers to Q's then ALP will have
the Σ error (past+future) that is close to that of the ^{truth} model ("if Any").
I'm not yet sure this is true.

I could perhaps should - for my own use - write a paper on "Convergence Problems
Revisited". Discuss Gacs results. Emphasize final term. Explain how the
k-L distance convergence is a stronger result than Σ error. -
Also using not k cost but Σ (costs of true model) can be significantly better.
Then proof of Bay data non proof for QA data.
Perhaps discuss General "kinds of errors": Σ error v.s. Model error.

Lecture
A "stop loss"
if I can't get
to "Report" don't
continue!

Final Run also can be used to find difference in convergence rates betw. using
 2^{-k} vs ALP. Say 2^{-k} is a factor 2^b worse than ALP (which = 2^{-kb})
Then Σ error for ALP is k vs 2^{-k} is $k-b$ for ALP.

For SM yield, hvn. ALP is 2^b times that of 2^{-k} . This SM yield is very easy to show.
At any time we put bet current fortune x_i on $i=0$ and $i=1$. If PC is prob of
corpus: we started w. #1 fortune we end up w. #PC = 2^N for N bets. Each bet we
double our money unless we lose. So we make $2 \cdot P_i^N$ on the i th bet:
where P_i^N is PC of i th money bet.

If there are k different bets (radix k) we win $P_i^N \cdot k$ on i th bet so total yield
is $\#PC \cdot k^N$ for total corpus. on SM

A very interesting theorem: If you used a "Blind" method to obtain yield over past N yrs,
then this is same as if you proposed strategy N yrs ago & credited for N yrs.

Since ALP is a form of "Blind", it is a v.g. method of Estimating future yield. his discover
Maybe not! Say prob of 1 is $\frac{1}{2} + \frac{1}{4}$ (at N=1) is one ~~more~~ this is 1.5.
Odds are always even. So yield per bet w. future, the error in prob estimate $\rightarrow 0$.

We "know" our yield will \downarrow . Maybe not! Say this is true rule: we use ALP & get
rule that approaches P_i , but we don't know that! We make a complex rule for PC
calculation that is $\approx \frac{1}{2} + \frac{1}{4}$, but we don't know that.

Q: T. random regression rule gives PC of 1 as a cond of past k bits.
It could give $\frac{1}{2} + \frac{1}{4}$ w.o. our knowing it. This sort of bug is likely to occur in SM.

(28) is very interestingly relevant to Action Alpn. yield problems.

But all it really says is that a time series need not be "Stationary"!
Hvn. Note that ALP is by no means restricted to stationary time series.

6/9/02
ID



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 649.40: A poss. way to do Expo (Report):

Describe pure QATM; ~~Give~~ Give format of detail on how it works: Then discuss various diffys
is ~~how~~ how we can deal w. each: Also some further benefits is diffys of the "forms".

So we start w. QATM in which updating occurs before each problem is solved -

We begin w. D-IND probs, indeed as such, so TM knows what kind of Induction Model to use.

At first we use a standard applied on d-functs.

For pure D-IND, the goal is to find a Max PC, D-funct ~~that~~ that generates f. corpus (is maximal Q's correctly)

[Is it poss. to have ^{D-}recognition functions to divide the corpus into parts, is have
d functs (w. possibly diffnt statistics) for diffnt "kinds" of Q's recognized? (see (.26-.31))

I was thinking of Riz's Recog. funct. being part of GPD. Tho the recognis. D-functs, GPD
can have D-components! If fact, for deciding what kind of PST should be used

on a problem, GPD can use D-functs w. impunity (perhaps) since some PST's are
"Universal" - is they can solve any solvable problem (eventually). - But otherwise, it would seem
wrong for GPD to assign a single PST to a problem - given it PC=1.

Hvr, at the begining, we only have a small no. of PST's is it may be first none of them
will solve the assigned problem. Hvr, Riz may be oia. - its just that the "best" PST's take
cc=∞ to solve certain probs! We still get S-mto from first "soln".
to choose PST's

There is a mix-up here betw PC's in GPD is PC's used in "external" induction
problems.

To start, there is only one PST for D-IND probs. - it is AZ141, w. perhaps
some hours for finding common sub. functions, is other things like that.

Then, there is the poss. of (.08-.10)†. I think this gives a new class of functional forms.

I could start by indexing the Q's so that TM could ~~recognize~~ divide
several appropriate D-predictors (Q's). Then, later, we know about the indices

is TM has to invent a d-funct to recognize what the indices "told" it:

essentially a D-INDuction on the indexed info.

O.k., so AZ141, w. a trim of .08-.10; 26-.31 may give some reasonable

~~induction~~ D-induction capability.

Ok.: Now, can I have TM simultaneously work on the S-IND problem of modifying
d. PD that guides TM's search of D-INDuction, D-functions? :: But we already have

AZ141 (.22-.23)† w. some extra hours.

The PD induced by (.22-.23) via AZ141 + hours, it ~~is~~ fixed ~~at~~ 24pt.
414.18 tries to Max it. Perhaps the S-funct involved is cap that ~~is~~ ~~is~~ → 651.00



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 650.90

invariants/norms regularities in the D-induction process. The problem is one of speeding up

to ~~the~~ Maxim of 414.18. This could involve context dependence of pc's of concs, w.r. to D-induction functions. — T. idea of how O^{TM} is a 5-funct of O^n, Q^{TM}, A^{TM}

.03

all previous data of traces of TM Activity.

.04

SN

O_n "Cure Cancer": "Cure cancer" is a BIG problem. If we had some criterion

for Cure, TM could try to satisfy this criterion as soon as poss. Could we formulate this as a regular option problem for TM, w.o. "Look ahead" explicitly spelled out —

So Real TM would automatically, by itself, discover "L.A.?" Actually, this is a sat of thing that I've been worrying about — that a TM with a big problem, is, essentially,

a RTM w. large h. — the BAD kind! It has "Goo": you can't turn it off, after it does

global self improvements, etc.

f2

13

So 2 interesting developments! ① 650.34 ~~650.34 - 651.04~~: Perhaps clarifies meaning of context

dependant pc's of concs, it perhaps generalizes the idea it tells how we evaluate the needed pc's. This is the idea of "internal problems" being "one kind of problem in TM's

corpus! (648.90 - 649.90) ← is on this: I want to clarify it! → 652.40

.19

.20

② (651.04 - .12): That I may not need to teach TM "Look Ahead"!

Re: ② (.20): It may be that improving the GPD (in the narrow sense of improving Guiding PD for Lsearch) may be inadequate for ② to occur. It may be unclear that TM is able to design a "essra" arily good PST's. (which is part of GPD improvement, but it's probly poss. to restrict it so that TM wouldn't get "out of hand"): Just how Brita such a TM could be, is unclear!

Try to work out the details of discovery of ② → see just what features of the system are actually necessary in order to get ~ L.A.

It seems clear that a PST could do L.A. and find solns to problems faster/better.

Probably TM would have to do much "Reasoning" to devise such a PST.

Simply improving Guiding PD would seem maximally steady ∴ could not do L.A.

Would Long during Lsearch be better in this respect? Superficially, No.

For "Cure Cancer" TM would have to do expts. in R.W. Hrr, this is not superficially ident. from doing "internal" Expts

Whoops!

Any OZ problem: A finite time limit, T. ∴ TM will use L.A. to solve such problems! Note that OZ problems are an impt subclass of probs involving PST optimization.

Small T values don't give very interesting results, but we gradually ↑ T in the TSO!

Normal IND probs (i.e. concrete type of problem TM does) are OZ probs of this kind!

6/10/02

11894 205: Jan but v of 202: "Not Explorer activation. See p 673: vaguely to go



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 651.13-19! 652.34-651.04; (646.40-647.40) → (647.12-20) is wrong idea

02: 651.40! Note also that INV probs solved by Lsrch, can also use LA! They try to find a fast Mapping ~~func~~ from prob dom to soln. — T. fastest such mapping functs could use LA, before they generate their final output.

05: So it looks like LA, is (as a special case) Lrng during Lsrch, and commonly occurring things. Hvr, the way I have been doing "Lrng during Lsrch" is different from .02. In .05 we use normal Lsrch w. known bound on CJS for known solr. My "Lrng during Lsrch" converts it into "act-Lsrch" — so we can't (easily) estimate CJS.

02 So .05 is probably much better, less EM, than what I had been proposing for "Lrng during Lsrch".

.02 (is what precedes it) seems to be a radical discovery & perhaps GREAT Simplification of TM!!

— We use straight Lsrch all the time for all problems!

← Apparently Not! see 04!

05: For D-IND probs; for early work, Prere can be solved as INV probl we look for Q → A funct in PC order, in Lsrch, so Pr's automatically gives us (some) short codes.

06: Hvr, INV Lsrch does it (ordinarily) do .05: It looks at the {Q_i, A_i} set & tries to find a function that maps to set into a Q → A_i function for all i. It is not really a "short" function, but it is fast. Would TLU do it better? (remember: for D-IND prob, there is no ambiguity for the Q_i → A_i; for each Q_i, only 1 A_i is poss). — It would seem that TLU would be a valid, fast soln. to the single D-IND problem, viewed as a INV problem —

So it would be of no value: TLU is no extra pollution.

Anyway, the moral is: Regular Lsrch (w.o. ll lrng) can do LA & probly other v.g. things: including Lrng during search. — For INV & IND & 02 & Opt probs

Hvr, while straight Lsrch can be v.g., because of scaling problems, we need Adaptive Lsrch. Whether it's best to have Lrng during Lsrch is unclear. T. advantage ^{of Lrng during Lsrch} is Prere

② it speeds up adaptation ① It seems to deal w. the dirty of many ~~cond~~ cond's being "about the same" — which make normal Lsrch very slow ③ Prere may have been a Real advantage of Lrng during Lsrch. ④ another is that one could dispense w. Lsrch

completely (as Prere was supposed to be v.g. because Lsrch is not optimum under usual into conditions)

— by working on a cond until it was (by // evolution process) clear that that cond was no longer optimum, ^{choice} best switching would be expected time to soln. ← (Prere last in spirit of my "New soln" to Hutter's search problem).

This suggests that "B (ind) Lsrch" could eventually find a use ably clever Heur!

Well, it could, but at enormous cost! w.o. Adaptive Lrng, Lsrch is very Limited!

As I see it, there are 2 extreme versions of TM: ① pure Lsrch w. lrng. betw. problems, only.

② No Lsrch, but ll evaluation as (31-33).

Versions anyway distance betw. these 2 are possi.

6/11/02
ID

653



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

30: 652.40: I want to write up clearly the old version of TM, w. reasons for various parts, before switching to the New version (published).

But first: 651.13-19; (652.00).

①. Key: Consider first INV prob^s (2) One (not exactly L such) way is to do trials for direct soln. in PC order. All trials take about same time so L such is not very relevant. This is usually not a v.g. way to solve INV. ~~prob~~ probs.

② Register-Ind L such; Tries various funcs to map problemn. into soln. Note 652.02: This method has all the power of a general PST, & can do L.H. It can make hundreds (ex. permutational) trials & Lrn. from those trials. — These trials are within each PST used.

I guess that early TM INV. solving routines will first try to find simple, general func^s that maps all problem decns into solns. Next, it tries to categorize problems, so it can use a different function for each problem type (to map problemn. into soln).

then Making Measurements ("Obs") on the problem & deciding what to do on the basis of these Obs

③ "Ob-OP Algebra" can then be used. Further developments may wait for TM to learn how to "reason".

INV

The improvement of the "Guiding" P.D. for ~~the~~ L such; This is its "updating" problem.

The updating problem for INV Al such: This is a QA type induction problem.

The Q for each A is ~~the~~ data points; The first part is the soln to a INV prob. CC (in the sense of Mapping from ~~the~~ prob to soln) & the second part is the info. T had at that time: This will be all previous prob. soln pairs & all traces involved in their solns, & the times taken for these solns.

T. data pair derived is assoc. w. each INV prob. solved by TM. T. complete Q at time to will be all such pairs for problems given at to + to.

What we want as output, is a func that is able to look at the Q & all previous data in TM's TSCQ; For any PST's co-argument, it is able to give a PD for CC needed by that PST to solve that problem. More generally, it will look at the data (including latest prob decn) & output a set of PST's w. assoc. pd's on CC needed to solve the problem.

28

~~How to Select a PST or give a guiding PD to~~ How to Select a PST or give a guiding PD to

29

the list of ~~the~~ PST's, is not clear! One possy: For each CC value T, there

30

will be a PC that PST_i will solve it before that time. Its like a horse race!

So one sets the prob^s that a given PSM_i will solve it quicker than any of the others:

To obtain this PC is easiest via M.C. Carlo. From each PST_i obtain a M.C. Carlo time of soln. (it can be ∞), then compare values for all PST_i's & pick the min.

Do this many times to get D.F. for each PST "winner".

37

T. ~~just~~ just for using the former D.F. as a L such "guiding PD" is not obvious

& may be ~~found~~ it is not a particularly good way to assign PC's to PST's for a given

problem — in fact, w. known T. into we have on CC d.f.'s for each (problem, PST) pair, L such is usually the best way to go!



20: 653.40 : One justification of L search is that it exercises a variety of PST's so that one can better choose betw. them for a given problem.

T. Uncertainty of 653.28 ff is something I decided to "live with" for the present. [WON do seem to offer a better soln.]

Indiv. of t. Q of 653.28 ff, we would like to create a pseudo d.f. in to 653.29-30 as t. output of an inductive inference process — or some other kind of data that could be directly used to guide problem solving.

T. update model types discussed above is one input type. Another kind of update!

20
10 We are trying to Build a problem \rightarrow soln function (i. prob type can be any type: INV, IND, OZ, opt.) We want an induction scheme that looks at all previous trials of functions like that in .09 & induces good trials. On the basis of all such trials in the past, & their degree of success or unsuccess, TM should induce a pd on all possl function, F, for any particular problem (?). But wait, F is a funct that's supposed to work for all problems & all problem types. The "past data" should include the CC of t. function in obtaining a soln. F. funct F could involve search, or recursive operations that are very time-consuming.

20 I want to desc. what I had in mind when I wanted to pc of a conc. to depend on "Context". (i.e. .23-.29) \rightarrow see 653.28-40 (23-29) deal w. Generation of PST's One way to generate them.

21 **2 aspects of Updating** ① wrt PST's : Generation & Modifi of PST's

22 re assignment of pcs to them (conditional on problems).

23 ② Generation of new/modifi. of PST function by constructing them from "concept functions" — assigning conc's pc's reference functions of "Context" of PST; so pc depends on present "state" of function being created, as well as on t. problem being solved. — Thus we obtain both a PST & an assignment of pc to a particular problem (or to actually sup problem!). So ② fits into ① as part of (21)!

23 **VERY GOOD!** .21 & .23 ff make much of the update procedure clearer by integrating

24 2 diffmt aspects of it.

25 Is it clear that .23-.29 is generating a PST? Is it trying to do $O^m = \text{Stnet}(O^0, \text{Gen. Atn}, \dots)$?

It does seem to be "Updating". First, be clear on the legitimacy of what's being done.

Well, O^0 is part of GPD, so it would appear that it is "UPDATING"; so "It is legit — & where are Details?"

We have O^0 w. assoc frequencies of use of various concs... & various "Definitions". \rightarrow 653.31 is "details"

Consider Updating O^0 in view of Q_{m+1} Atn: First attempts: O^0 as O_{m+1} directly,

so how bad $P(Q_{m+1} | Q_m)$ is. ② In view of R_m s.s.z, modify continuous params of O^0 ;

③ make new definitions @ response O^0 (This ~~is~~ last is tricky! I'm not sure about its legitimacy;

I probly works for English text... but may be "special".

6/14/02
ID



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20:654-90: (3) If the pc of Ant1 ~~is~~ still/looks small, try to find a "recognizer" for Ant1 & make a special function for it to \rightarrow Ant1. We can use the default ~~function~~ & pc's in the old O⁹; since it's "all one function" (including the recogn. part!)

I don't know how well (100) fits into my previous work on "Recognizers & their updating" ...
As in "CPM": Nov 8, 01. Section on "Updating". ^{CPM N8:01} ~~Updating~~ (100 may be \equiv bid, pt §2; 2.2)

Note: In CPM N8:01, I was using hard (d-functs) recogn. functs. I may want to use soft (s-funct) recogn. functs. \rightarrow .28

T. UPDATING ^{Algs.} in CPM N8:01 (i.e. others that I may derive) are meant to help TM "Get off the Ground". Hvr, I want to make the Algm. "Simple" so \odot I & other humans can understand & ; pgn, & debug it \odot TM can understand, ~~modify~~ modify & use parts of it for other problems, other parts of the System

The 3 ~~to~~ "Scaling" diffcs:
1) using Non-adaptive search works for small problems only - that are within the vocabulary of the Reference Machine
2) If cps of concs are problem independent, then, as no. of concs grow w. corpus size, the pc's of concs tend to \downarrow . I'm not sure this is Nacy, but it did seem to occur in the TSG
I did ~~it~~ for ANL in Saarbr.....

3) As corpus increases, it becomes very expensive to test entire corpus on trial Phidius of existing PST's. We can use statistical testing, & finding of "diff problems" to reduce stat. of testing. Also, partitioning corpus: so that only a small part of corpus need be tested on a PST - because that PST is applicable to only a small part of the corpus. This has to do w. "Recogn. Funct's" \downarrow

~~As~~ in .07, there is Q of "Hard" v.s. "Soft" recogn. functs. Hard are usually much better, if one can make them. - in v.8, it applies to .19 as well as in .07 \leq CPM N8:01 §2.

In .15, in the Saarbr ANL ~~test~~ TSG is "solution": I felt that assigning just Barn seq. pc's to concs wasn't "specific" enuf. - which is perhaps how my idea ~~of~~

The importance of "Context" evolved

Well, one way to automatically get context is to look at funct's around the funct to be measured. This is like T.S. / proln. using previous ntues as "conditionals". This perhaps ~~equiv.~~ ^{invents} inventing "Macros" or "combination functions". Q.k. But I wanted to invoke context on a broader more "Global" level, as well - so the pc of a conc should depend not only on local context (w/In the function being created), but context like previous history, historical traces, the nature of the problem(s) being solved; previous solns to problems.....



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 655.40 : New Contrast 654.31 w. 655.31 : They are about similar things, 654.23 - 22 is longer
of a contrast : 654.21 v.s. 654.23.

02 "Updating" is the response of the system to a new Qnt, Ant; (problem soln pair).
It tries to make (Qnt, Ant) fit into the system.
It does this by adjusting the functions used (654.23-29; 655.31-40) as well as the assignment of adjusted functions to new Q's.

T. concept of "Updating" is changed and it is done before each problem - specifically for that problem
Sec 657.00-03

Also, it can invent new functions (PST's) as well as the p's that associate new Q's w. new PST's.
Updating is any thing the system does to prepare itself for new problems: Usually this involves
Integrating the latest QA into the existing "system" (or equivalently modifying the existing system to react to new data sets as well as poss. (in the available time)).

10 Sounds Good for Exposition: That TM works by assigning Q to
11 (prob list) to 1 or more PST's. Post solve the problem.
Updating is .02-09

So .10 could be a v.g. introduction (or even abstract) ...
After prob. is solved (or after soln is known) (TM need not solve any problem)
We "Monday Morning Quarterback" it is try to see how TM should have acted (given prob's soln)
by knowing prob's soln is asking how to run the system so it could have solved it "best" (fastest, or best for available CB).

20 T. System itself is .10-.11; .02-09; but I'm still not quite sure of the Details of Operation.
Say we give it an INV problem. The GPD assigns it a PST: Say it
uses LSrch using GPD for "Guiding PD": to solve it.
To do QA probs: GPD assigns to new Qnt to a PST - which is itself:
It then poses a GPD(A|Qnt) - a d.f. as output.

W.B. T. factor of 4 of SSB should be "factor of 2" only. Due to extra number of Direct P's to Updating: In advanced systems, this discrepancy may be very large! However, that the system would not be doing LSrch (because of long) is not relevant.

30 To Update the part of the pd devoted to INV probs. At first, it does it update part first assigns to INV prob to the PST that does LSrch. It does hrs update the part of GPD that is the Guiding PD for LSrch. To define that update, we have to first define the corpus to be used. This will be $\{ \text{Problem } P_{ij}, F_{ij}, T_{ij} \}$
31 templates of the part: P_{ij} is the d.f. of the INV problem; F_{ij} is the function successfully solves it. T_{ij} is the cc of that soln. We also may have P_{ij}, F_{ij}, T_{ij} for problems in which after $cc = T_{ij}$; i.e. problem had not yet been solved.

35 T. pd for updating has to make a function, Z that takes P_{ij} & F_{ij} as input & give d.f. on T as soln for output. We want a Z that maximizes 4.14.18. This will be fine for the existing F_{ij}, P_{ij}, T_{ij} corpus of .30-.31

hrs, we also want to invent new F_{ij} 's.
This can be done by making a Grammar for the F_{ij} 's: hrs, we also have to make a grammar in which (P_{ij}, P_{ij}, T_{ij}) 's are elements! I don't know how to do this!
HVR Note 657.2-1!




JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 656.40: From a "pure induction" pt. of view; Given a new P_{old} , to find/invent a P_{new} \rightarrow its likely to have a small ϵ . (say $\epsilon \uparrow$ is small) — This amounts to updating before working on the problem which I felt/was better because ϵ could update "selectively" for the new problem — a decline in how much time to spend on it. (this note got rid of "factor of 2" in "optimality ratio" \odot).


03: As an "example" (not actually hard INP problem): consider ANNs. We want a "short" function from "3+7" to "10". This is an IND problem; it has a particular "Lurch" that can work: unlike from $3+7$ to "10". This is on IND problem; it has a particular "Lurch" that can work: unlike from $3+7$ to "10". This is on IND problem; it has a particular "Lurch" that can work: unlike from $3+7$ to "10".

i.e. make $f(x) \rightarrow A$, finds in pc order (at least order); We used a Guiding P.D. for the Lurch. The "Goodness" of a $f(x)$ is given by its pc — and perhaps also, acc needed to execute it for ϵ , entire TSO (perhaps?).

0: T. Updating that Turing did for "symbolic regression": He used as TSO, pairs of rules. He had a simple evolve on cond. functions given by a few pc's. After solving many problems the pc's migrated to $\phi \approx 1$. He used a reinforcement mechanism to modify the pc's. My impression was that he should have used Laplace's Rule vs Reinforcement. (he had to try various levels for the reinforcement parameter). T. doubts of his reinforcement mechanism are "inf. Certain kinds of mechanisms can be easily refined to Laplace's rule". I think one uses a 

His method seems close to feed forward Neural nets. "adjusted" "to be chosen by trainer" reinforcement parameter. Anyway: He used a particularly simple space of functions & PD on those functions, so finding a soln. by simple hill climbing worked fine. (33)

20: what I'm trying to know, is get examples of various kinds of Updating.

21:  on Second Post in 658.35, F should be problem independent! In fact, 658.35-40 is praby ALL WRONG! There are different amounts of nonlinearity in TM Models!

The most General: Input a Problem: Output = soln. Next: Input problem, output pd on soln; then use Lsun. Next: Input problem, output's pd on PST's. Each PD could derive for Run category to other non-PST's to problem or work on it directly (or usually) probabilistic (mix of these).

29: Or: this last can occur at all levels! Even to top 3 primary input level. So when we insert problem into TM, it may try to solve it directly & probabilistic & local & show to other PST's. .29 is v.g., but perhaps draw up simple, prelim System, then show how various

32: improvements can be made. J. had a pd on all instructions to reach "position" in the pgm. So, with 7 instructions opcodes: pgms of length 10, he had 7^{10} pc params. Here, most never occurred. 658.10 658.21

34: On the other hand, I had been making pc's of conds indep of where they occurred in the pgm. Invoking "context" as relevant to pc was a way of looking at previous parts of a pgm to induce pds on the "next one". Mutating (in a very concrete — xtreme) old, successful pgms, is a way of inducing a pd on future independent on the past. 658.00

Perhaps list some kinds of problem solns. — then show how system gets to them.

6/16/02

GA(.00-.20)

A Grammar for (Funcs &)
.15-.20 in particular.

Z141(30)
Looks like V.G. Approach

658



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 657.40: In line w. this: Say we have a Sample of k G-params: We can then make a probabilistic study of various features (conc.) in that sample & indicate a pd on forms from this.

Naturally for a BAG of forms, we have a case count of each form type. We extend this to concs in some environments, but the case counts could be various: Say for each instance of a form, we use as "case count", some monotonic funct of its G-params. In studying "pc's" of various concs in such a sample, we can use those word "case counts" to get some thing like probability of forms: They will be biased toward by G-params if we do a assignment of case counts, to concs of concs properly ...

perhaps if we didn't we could get ~~output~~ output Dif. on expected values of G-params of forms. Since we are using a Bernoulli on concs, (w. indep concs: ~~we combine concs for interaction~~)

for each G-para, we have a set of G-values: ~~is~~ ~~is~~: do we want a simple condense them into a scalar "case count" (perhaps + case is var.)

from (11L) — What's the best way to make a dt. on (func, G)? $[F_i, G_i]$
A better approach than (11L): List the various funcs & their assoc. G-values: How can we Postcode this data? Make a lot of data of concs (i.e. piece of Discourse dur is as yet unclear!) \rightarrow some how do it as max for corpus.

For each conc, have a pc def. \rightarrow a func of G. We have a prob. corpus in terms of ~~more~~ concs: the pc of each conc, depending on G. We want to do for corp. w. ~~min~~ min cost.

657.34 A poss. modify/improvement: try Juerson's 657.33-34: A direct pd of concs for each instruction address in the Form: Poss: Postcode "Pushing": Post pts in space will have null data. \rightarrow 669.30

20: Re: .15-20: We'd like a 2 or 3 param way to express pc of a conc. as a func of G. Think of it this way: for each G-value, we have a Grammar. Each Grammar has its own concs, but ~~data~~ costs of concs are shared because they are all part of the ~~same~~ same corpus.

Single Corpus.

As in ~~the~~ General Case, values of real params need not be known exactly: ~~then~~ the pc of the corpus is the integral over all param space of the pc's of each param. in param. space: This integral is usually easy to approximate.

Expressing pc as a ~~conc~~ conc as a func of G (or a "G-fermate")
 $pc = e^{kx+a}$ i.e. ~~ln~~ ln is linear in G. (2 params)
More complex: ~~ln~~ ln is quadratic in G (3 params)

35 We may want to modify G a bit w. a monotonic function before doing these things. \rightarrow 661.12

36 SN on A2141: We look at the corpus & we see ~~conc~~ concs that ~~may~~ may be indep. ~~then~~ then we look for dependencies: If any apparent dependency is strong enuf, we can make definitions ~~of new~~ of new concs that are indep. ~~then~~ then the resulting expression should have fewer dependent concs. If strong enuf dependencies exist, we again try to define ~~the~~ the minimum of existences. \rightarrow 661.20
[Hur, Note ("objection") of 659.32-40 \rightarrow 659.32-40]

SCALING
For just the context of concs: w.o. context, pc's are too low. Context affects pc's: Dif. of pc's concs at pc's of context choice.

The "OMO Scan" Output Missing Way OMO way out of many for WOMO scan.

This is the principal of A2141

\rightarrow 661.20
Spec 659.00

6/17/02

2141: .22-31

The problem of defining symbols of unusually low frequency (But note: 32-40!)
NOTE

659



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 658-40: Could 658.36-40 be used to deal w. unusually low freqs of concs? One trouble w. Pami's
Pam realizing that they occur, doesn't ↑ pc of corpus much! Could a m pc's of complete concs be

02 to only/best way of dealing w. Pami's? (.22)
03 From a practical pt of view, for analysis of failure of systems; unusually low pc's,
is the main problem. — estimating upper buds on very small pc's.

Usual way of doing so: List various failure modes, for each give necessary conds/params for failure;
Coincidence of Events vary for each failure mode. These events are rare & (perhaps) uncorrelated: If we can
empirically show each may cond has pc < 10⁻³ & if rare conc, P_{is} → 10⁻¹².
We can estimate pc of events empirically: That they are uncorrelated cannot be seen empirically;
Coincidence is too rare: So we must use "reasoning" to estimate correlation
(this is very uncertain reasoning!)
Other big difficulty is listing all failure modes.

.00-.02 is ok for not so small pc's but can't get pc < .01: It would seem
that freqs of .01 or .001 could be obscured, but their effect on total pc of corpus would
be quite small! Maybe not so! If SS is 10, known pc of event is .01 divides

pc of corpus by $(\frac{1}{100})^{10}$ error of 1 in 10 gives factor of 10 in pc.
T. probly that pc is smallly .02 & it occurred twice as often & should have very small
likelihood: its ≪ 2¹⁰ (= 1000).

22: In .00-.02 it would seem that .18-.21 would be relevant. If we think a pc = 2
and it occurs x times out of n symbols is $\frac{x}{n} \neq 2$ by an adequate amount, then we would do well
to realize that pc is not ≫ 2, but ≪ 2. This is true if $\frac{x}{n} < 2$ just as much as if

$\frac{x}{n} > 2$!
e.g. A & B symbols: we note that AB has a freq < f_A · f_B. We certainly should be able
to get corpus pc by defining C ≡ AB (!), because: before it does, the pc of AB occurring
would cost pc $(f_A · f_B)^n (1 - f_A · f_B)^m$ (m+n symbols in corpus)

After defn it would cost pc = $(\frac{n}{n+m})^n · (1 - \frac{n}{n+m})^m$; which is the corpus P_{AB} since
 $(\frac{n}{n+m})^n (1 - \frac{n}{n+m})^m$ has peak when $x = \frac{n}{n+m}$.

32: Hm, to forgetting argt (as well as 658.36-40) would seem very Simplistic in view of
my more complex models for 2141, involving alternative parsings — The idea that not
34: all cases of "A.B" should be parsed as "C". While we try to get a max pc parse,
we realize that this may not be the only "best" parse — many others will be as good or almost as good.
T. finally best "parses" are those that get high pc in parsings (b) enable new deduc
that a pc for parse. This suggests that much backtracking may be (acceptable/necessary).
Normally the proper frequency will be considerably less than that obtained if all A.B's were
parsed as "C". There will be usually a very large no. of parses w. same pc whose total pc is > P_{AB}

6/7/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 659.40: Hvr., wrt 659.32 & 40, ~~There does seem to be a paradox!~~ ^{There} If we code a corpus (using 2 processes to get Laplace's rule) by using pc of each symbol \equiv rel freq of that sym bol (in equiv. context) in past corpus, we get a lapit coding — but it doesn't take ll posns into account!

[SN] A possl. explainn. of $\frac{1}{e}$ factor in Z(14) analysis: [f is observed freq. of A.B is α ,

then usually, if A.B is a legit. data, w frequency β , then A.B will usually occur w.

frequency $e \cdot \beta$: A fraction $\frac{1}{e}$ of those (i.e. β of them) are to be passed as AB (\equiv C)

T. rest are legit. meta random juxtapositions of A & B. ^{Actually, it's not so simple,}

because when $C = A.B$ is defined, the frags of A & B \downarrow .

[f b. large w/ins, then, it ~~is~~ $F_{AB} = e f_A \cdot f_B$ (for ~~the~~ $\gg 0$)

then ~~if~~ replacing any of the AB's by C would ~~decrease~~ \downarrow pc of sequence.

Hvr, this ~~seems~~ wrong! AB will occur w. freq. $F_A \cdot F_B$ if A & B are indep.

for ~~the~~ C : say we define C w. fc: $(C \equiv AB) \leq 0$

now pc of corp ~~is~~ $F_A^{f_A} F_B^{f_B}$ $\sqrt{x^4, y^4} = xy$

now pc of corpus is $(F_A - f_c)^{f_A - f_c} (F_B - f_c)^{f_B - f_c} f_c^{f_c}$

under what condn will it be $> .15$? write ~~the~~ ratio (u3 mod operate) ~~gets rid of n~~

$$\frac{(F_A - f_c)^{F_A - f_c}}{F_A^{F_A}} \cdot \frac{(F_B - f_c)^{F_B - f_c}}{F_B^{F_B}} \cdot \frac{f_c^{f_c}}{F_C^{f_c}}$$

[forget rest of symbols!] we actually \downarrow no. of symbols, so ~~their~~ $f_c \leq A$.

This ~~pc~~ enters in both .15 & .16. Could I elegantly represent it by 4 smly symbols? D is it? $F_D + F_A + F_B = 1$

$$F_D^{F_D} \rightarrow (F_D + F_C)^{F_D}$$

$$F_D' + F_A - F_C + F_B + F_C = 1$$

$$F_D' - 2F_C + F_A + F_B = 1 \quad \left| \begin{matrix} \cdot 2 \\ \cdot 1 \end{matrix} \right.$$

$$F_D' - 2F_C - F_D = 0 \quad F_D' = F_D + F_C$$

[SN] Make some kind of INDEX of empty works

How to reference them is unclear: Try to standardize on names of various effects.

Start out w. various questl. synonyms: followed by standard index number. P. 667 is a start

32: (28) Check this $n=00$ stuff w. Algebra I did in finiten — use ~~the~~ formulae. ~~larger~~, & see if I get ~~the~~ results: This could easily point out ~~the~~ disparities (if any)!

Bugs in Reasoning! In particular, check ~~the~~ ~~to~~ ~~make~~ ~~sure~~ I have the no. of each symbol. Remove simple formulae for pc.

Q: could I incorporate the alternatives parity into a kind of code which it didn't insert new symbols, but just that if A occurs, B has a special frequency ~~insertion~~. to follow. I could use a code n to Z(14). This would be a conditional pc for certain symbols.

6/17/02

INDEX

GA vs STACK
611

662



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

- 1) Z141: ^{Z141} 659.32 (A new way to look at + problem)
 639.32
 660.15 - .40 : A simple way to check old Z141 formula. (see if $\frac{1}{2}$ factor is legit)
 661.12 on gen-derived circumstances
 J an earlier ~ @ within 6/17/02 minutes (HAR) on Z141 : re: $\frac{1}{2}$ & parallel passages: Note 659.32 - .40
- 2) GA: 658.00-35 : Idea on a P.P. for [Cand₁, G₂] - a "language" : (missings?)

6/21/02
ID

UPDATING INDUCTION : 100 # (V.G.)

see 663.10 for update ENV



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00 : 663.30 spec : INDUCTION & its Updating: We want to solve problem to take $\Sigma [QA]$ corpus

to find a good O^N (4.14.18). Since it's an OZ problem, the "normal" LSCH soln. is to LSCH over poss. mappings from $[Q, A]$ to O^N .

03 \rightarrow One elementary mapping is to simply try O^N words in PC order like LSCH, if it works. Its just another possible mapping method to go from $\Sigma [QA]$ to O^N .
using a Least Squares criterion. The PC's themselves look for Cond

06 There are certainly other ways to go from corpus to cond. — Say look for regys in corpus. New regys are suspected from known regys. — In general, regy types are probabilistically related.

GPD can contain info of that sort... to suggest a new class of regys to try out. Basis of regys found this far. — Also GPD can suggest Obs to use that inform GPD users to suggest regys to look for

0 One unita devise a "grammar of regys": any Universal mechanism can occasionally "List all poss. regys". I want fast tests for them, & of course the probabilistic relations

14 Partw. them as (.06-.10). .20
15 In .03 updating wants to improve .03: This can be done (one way) whenever a P.D. used in a PST, we can usually improve that PST by improving that P.D. We start w. an uncond. P.D., then we learn to

19 "contextual clues" we may be able to find. standard techniques make it conditional out: Corpus is on O^N values for earlier corpus & any other T. p.d. on O^N space is initially independent of PC's of concs. Also defining new "lines"; Then PC's of concs can depend on context! NOTE 658.36 ON Finding Regys in a Corpus .26

20 One good type of regy: This type is sorted Additio. We find one of them, then find another & indep. of it: first, & find new ones, w.o. having to examine the whole corpus. looks v. good but is relevant here? this is on AZ-141, which would seem to make it Relevant here

22 Many kinds of regys are testable. smaller. We can make various Obs out: corpus that correlate w. various regys. Some sets of OPS are v.g. because as a vector, these OPS

24 "narrow down" & PC of certain regys vary much.

So: .15-.19 talks how to "improve" (\equiv "update") + PST that I'd been running about in .03; .20-.24 talks how to "improve" (\equiv "update") any given a formalized general class of regys (to be used for predn).

30 Note: "Updating" need not take advantage of/ use "latest data": It just does whatever it can to improve the GPD, work with available CB. Ideally, it should be done

31 before a particular problem, so TM knows what part of GPD to "improve". Quite General, but I do know of specific PST's, but maybe not completely general?

(664.00 #) & (663.10-.30) look fine! That I can start w. to techniques I had for updating those PST's Do fit into "T. General plan" — particularly, if I consider that updating can just improve old & create new PST's (among its other sub-goals).

36: (.22 B) N.B. In .03, 15: Improving the Guiding P.D. is an plan. of + more General problem of improving & Update Algor. Presumably, the Non-el. Update might be dealt w. by a later, "More Problems" TM. So at the beginning, TM really is set up to work on non-el. update problem. T. Q's, can TM advance and, (using the el. update) to usually work on Non-el. update?

spec spec 665.20 664



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 664.40 ON T. "50% Soln.": It would seem very el. to concieve of r. ratio of time spent on updating to time spent solving "external" problems, to be a useful param. to adjust. That in general, the system should decide when to update, & how much cc to spend on it. On the other hand, general improvement to system (i.e. updating) can really be done indiply ~~at~~ a limit, of when external problems are worked on. The low range "S.I. activity", just needs a Per cent. of cc — i.e. "proof" of 50% soln. allows for cc relatively indiply of when external probs are worked on, i.e. relative "pressing" of S.I. & external prob. soln. Doing these 2 activities simultly. (Time Share) (in 11) would seem to be fine!

11 Expo on "50% soln": Say we have 2 machines with w. same clock speeds working on Problem: M_1 has ratio of internal to external time share. M_2 " 50% " " " " " " " " M_1 will run $\leq M_2$
 M_2 will be at least as good. — it will spend if one slow down M_1 by factor of 2; at least as much time on both internal & external problems. It is not to say that normal M_1 is not > twice as fast as M_2 in solving problems — M_2 could be way faster than M_2 — but I think it's "unlikely"! — \$73.11

19 23 664.90 Sol. T. System: A problem comes in; It is assigned (probabilistically) to one or more PST's, by the G.P.D. They work on it (possibly by recursively breaking sub-problems to M. Total system) until the probs are solved, or until the cc is exhausted (i.e. the system has done as well as it could w. that (B)). parallel (time share ~~2/a~~ 2/a (H.W.) to P's, the system is working on the problem of self-improvement. (S.I.). S.I. is same as "updating". There are various elements of the system that can be improved — based (partly) on info obtained from problems solved. Some elements: (1) Assignment of total P's to PST's as func. of problem. (2) Improvement of ~~PST's~~ PST's (3) invention of new PST's. (4) Improvement of PD's used ~~by~~ by PST's (these are all "END" type problems) (5) Improvement of $f(QA)$ problem's function. [Better induction: Present S.I. problem can be expressed as 1 or more QA's] The problem of S.I. is "just another problem" for the system that solvers in it w. its other "External" problems. Most of the S.I. prob. & sub-problems can be expressed as part of the [QA] corpus.

39 Soln in addition to .20-.35, what's needed is dev of 2 types of probs solved (b) Various PST's up to solve these probs (c) just how S.I. can be usefully formulated as a prob. & sub-probs that are normally solvable by the system. T. exp. stuff does breakdown & decom of the system into parts that can be worked on/direct

6/22/02
JD

TSQ's. 16ff, 120ff How to write them



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:665.40 separately! ^{GREAT BREAKTHRU} so this looks like a Small Bro & Peru in book exposition is int. design of TM itself!

In exposition: Explaining what GPD is, will probably be quite diff. Comparing it to Google can help a lot. It's more problem oriented than Google. How does it ~~diff~~ differ from Maple or Mathematica? I can go into differences - Perhaps however, "talk" from in the ~~paper~~ paper/report.

Re: Google: it is constantly being updated - as new info appears & as new topics of interest develop; & as interest (output/output) or old topics change.

- 0 List things that have to be done for report/system Dev.
- 11 1) List & explain types of problems: Tell which types will be ~~in~~ in initial system. ^{663.10ff on INV 664.00 on INV}
- 12 2) List some PST's assoc w. various prob types. Spend much time on QA prob.
- 13 3) Look at various aspects of up dating. ^{663.10ff} ^{663.40 ff on up dating INV.} ^{664.00 ff on up dating INV.} ^{1st very very careful about it, how to update.}
- 15 4) T. TSQ: I have at various times, discussed various pieces (measures) of TSQ writing.

Perhaps try to know inst. idoes. (see 20).

Re: TSQ's: some ways (1) Top down (SBS) start w. goal problem; design conc. next to goal, design probs for each conc. (or set of concs).

- (2) Start w. Human ^{AI} Text, see what can be reduced conc next to TSQ to solve it.
- (3) Look at Solns to old AI problems. Try to design conc. next for Procs to MS - ^{in classic AI probs:} Symbolic Integration; Proving Theorems; Solving diff eqs, differential eqs; probs. solved by GPS. Probs solved by Humans in Newell, Simons "Human Prob. Solving".

(4) I do of getting TM to be familiar in say, Algebra, so it could learn Euclid's Prat talked about Algebra.

(5) writing TSQ's ^{for} various levels of TM knowledge: e.g. Pro solve linear quad, cubic equ. say, using one (merely substituting) heuristic.

(6) Try other - TSQ's for learning ^{today} certain symbolic Integrals.

(7) Try to get TM to Learn ^{today} various things that Maple does. (Reverse Engineering Maple)

(8) What about analogy? Learning to reach certain control goals. The beyond learning problem, etc.

(9) Look at various problems Solved by GA. See if I can get TM to do no Oram Mich Factor.

(10) Once I get the details of the (very system) (status/decided upon): Try a not very carefully constructed TSQ on Algebra, say, & see what messages TSQ's are - Recursion

6/23/02
CP

TSG's (cont) How to write Perm



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

60: 666.40 needed concs (via ~~proper~~ problems) (like a virus inserting a Gene or DNA snippet).
One way to ~~beat~~ TM to Inv a lot of Algebra: Just list imp ideas I want to Inv. in an "epistemologic
order" - look at CJ's betw. Perm, & determine needed concs betw. Perm &
design need problems/examples/HINTS or other ways to BRIDGE & massive CJ's.
E.g. I had some trouble w. ANh: so assume TM got over Perm hurdle & continue!

004
The TSG suggestions of 666.20-667.04 seem fine (particularly #10).
It may be that I have the problem "under control": Try to list the weakest parts: See if I have
rats-to-solns of Posc problems. The idea here, is not to forget f. solns!

09
10
666.10-15 is a good start.
T. INV approach of ~~666.10-15~~ ^{658.00-35} + 661.12 Af! seems fine for an advanced TM!
I need ~~an~~ examples of how to work simple problems, Inv. T. Prop. is based on ~~elementary~~ &
grammar of PST's structure for INV probs (PST's used in early AL work on INV probs).
Are there some simpler, more direct ways to solve some Inv. probs using Lsrch?
Perhaps look into those early A.C. probs, Tower of Hanoi, Missionary Conversions, etc. in Post Book on
GPS.

10
Actually, the method ~~666.10-15~~ ⁰⁹ are not bad for ~~an elementary~~ INV probs as well -
we just have a far simpler grammar for PST construction! - A so only very simple INV
probs can be solved. GPS has a kind of Grammar of sub-techniques: Backward
& forward chaining; ~~the vector~~ Construction of a Vector Calc. Deriving OPS &
Obs is an Ob-OP algebra to ~~the~~ ¹¹² of vector ~~calc.~~

10
Hvr, I may want to not discuss "elementary problems" in my early TM
discuss - Mainly work on QA Induction. ¹⁰
I write discuss ~~the~~ advanced INV soln methods (09) as a kind of "side"
remark in f. report (or as a remark for me).

124
T. most imp kind of problem(s), of course: QA & its updating. (666.00-40)
I DO want to (only E. context dependence of concs in construction of O^k conds. Just how to
context (& extended context) words, & just what part of system to PC of f.
definitions of f. context are "changed to"
Also look at Post shift in $CPM(N, B, 0)$ on Raspphinity func & Perm "Updating" → (668.00
SPEC

31
SN INV probs can be solved by QATM in what looks like close to normal INV work.
We present QATM w. is (large?) set of INV prob. derivs & their ~~the~~ correct answers ~~constructions~~
(that is "show answers"), QATM was to find a possible (or exact) relation from these Q to these A's.
This "relation" could be the same universal INV function that INV Lsrch looks for. The problem is
aspect of it can be used to find f. soln. by straight Lsrch!
SN + **PSLQ** ("Partial Sums" of squares vector is LQ Lower-diagonal-overlapped) Matrix technique. **Alstr.**
For Diophantine relation detection: Given $[R_i]$ real or complex, to find \vec{z} : (integers)
 $\vec{x} \cdot \vec{z} = \phi$ ($\phi \neq 0$). This is a polynomial eqn. (look in Google under "PSLQ") I don't think
it used to find show items (= integers) - but this is clearly not unusual. & only time back!

6/23/02
ED



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

30: 667.31

* T. main Discuss. will be IND & its updating: 667.31
667.03-04; .15-.19 are on the "simple" way to do IND. (QA); It (supposedly) describes "contents" of cases.
How does this treatment of "simple" QA, fit in w. a protocol CPM-N.4.01 - involving say "recognition functions" & Metrics?

106

In the forgo. discuss. of IND (QA) probs: I have been mostly (redundantly) discussing D-functs.
rather than S-functs. To what extent can the discuss. be extended to S-functs?

0

S-functs would seem to be a "larger space" than d-functs — But one can somehow get
space just as one searches over d-funct space. In addition, the S-funct space is "loft",
so we can start out w. ~~the~~ "Broad" s-functs, that fit every Q_i to some extent — ~~then~~
narrow to D.f. by "Hillclimbing": It may end up being faster than hunting for d-functs!

11

The shape of the hill (~~is~~ is ~~correct~~ correct) may be smoother for S-functs than D-functs, so sometimes
easier to climb. (Hill shape is ~~is~~)

14

Hvr. A big problem that I haven't worked on much: for each S-IND problem, one
has to choose at least one S-funct representation to use. There are several such
representations. Each PST for a problem may have a diff. 1 or more S-funct representations. At some
pt., TM has to interconvert them, so it can make comparisons!

0

What I want to do, is give more detail on how QA is done.
One method is to find a single d-funct (O_i) for a set of QA's. This is uncond.
feasible.

Trainer can give TM several such sets, identified by indices on O_i 's.

TM has sets w. separate O_i functs: cones can be shared betw O_i 's:

The recognition problem is very simple because of the indices.

Later, TM has to index unindexed O_i 's & into previous set of indices.

25

Later, TM creates new ~~indices~~ index types. T. problems needing new index
types are not solved by ~~old~~ "present" functs, so they are categorized together.

In .25 TM creates "recogn. functs" In the update process, S.I. consists of
a corpus p.c. This can involve changing recogn. functs.

30

When we "graduate" to S-functs, we can start w. d-recognition functions
than do S-recogn., because its "softer" & easier to hill climb. ~~is~~

We will always try to get all functs as close to 0 & 1 as possl. (to a pc of corpus)

So we will want them to evolve toward d-functs:

37: (11)

But note that S-functs can be as "zigged" as d-functs — (even more ragged!) —
In the report, I can just do b. to forgo. stuff on d-functs, S-functs: i (d/s) recogn. functs.

Having much detail on just how to do the updating is not a bad ~~idea~~ idea: I do want
to write a lot of (comments/hints on how) in my copy of the report

6/24/02

669



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

009.668.40: Next, I want to descr. how concs feel PC's & how they may depend on ^{context} context.
For d-functs, in which we have d-recogn. functs, we can show overhead of defns. but without

02 "branches" of r. recogn. functs. (15)

TNB In **CPM N.902** I had a set of R^* functs that ~~partitioned~~ partitioned \mathbb{R} corpus.
This is only 1 way of doing t. partitioning. This partitioning can ~~be done in many ways~~.
have many poss. forms. t. one in (CPM-N) was simple-looking, but not nearly
best, or even very good!

HVR, Partitioning (usually ~~d~~ d-functs, but often s-functs) is very imp.

When TM is first given a problem, t. GPD partitioning corpus (usually s-function)
by assigning 1 or more PST's to it.

14 T. partition question is "What ^{or what} part of (problem) is this?" \rightarrow 670.12
One very BIG, IMPT function of partitioning: To PC of soln: To make Lsrch "Easier" \rightarrow

15: 02 For t. "Context" problem: I want to ~~define~~ define 1. T. goal of t. process: What's to be
Maxd?

(2) T. corpus being vsd. (3) T. existence of CB < DA is ~~imp.~~ imp. Just what t. CB is may not
be so imp. (But I'm not sure of P's last).

19 T. "top Problem": Given corpus $[C, A]$ to find O^* 's of by PC method:

Subclass of methods: Build up O^* from sub-functs. (i.e. concs). This is guided by PC's obtained
by looking at previous problems like this one. T. PC of a conc. will depend on t. S. method.

"S. method" \equiv "Context".

Note .20: This is a subclass of methods to solve .19. { Juanon's method (30) may be a subset of this class
if it is by no means t. only one poss! }

26 .19 can be regarded as one of t. problems in t. $[C_i, A_i]$ corpus.

.15-26 does seem to clarify t. meaning of "Context": I think I have a similar, earlier
writing Postdoc about t. same thing.

Note .15-26 discusses a particular form of function Generation — specifically, AZA1.

Using Juanon's method (657.33, 658.21...) the learning (S.I.) would seem to be quite direct.
His priority params for each instruction are different variables position in t. code string.

His codes are strings; Mine are more like trees.

A seeming Advantage of Tree structure is that it facilitates defns. of common ~~sub~~ sub trees
(corresponding to common concepts (sub-functs). Whether P's is easy to do with "linear"
string format, is unclear: I don't immediately see how to do it.

Maybe "common substrings" (\equiv unions)?

T. 2 methods seem to not write directly. J's method seems to approach a d-funct at t.
end of t. such. AZA1 doesn't seem to do this at all! It remains a very "s"-function $\frac{1}{2}$ by root of t. such.
It would seem like good idea to study just how t. 2 methods differ. We could mix



6/25/02
f. 2 methods in a STEIN ish way (i.e. partial pooling of J's pc's 65821) — Hurs, I'd like to mix t. 2 methods better so J could have $\hat{\alpha}$ strus "or equivts": Benefit of a conc. prod. is used α times. (maybe put Pe's Q "on" stack" \rightarrow 611.) In J: pc at conc. is distinct only of nature of conc; a position in form (linear) in J: " " " " " " of nature of conc., position in form neighboring concs, natural problems, previous probs, etc. \rightarrow 673.3

[5N] Intth k-L form of STB3: f. loss of t. truss & f. approx pc's \rightarrow 0. This is very imp! If mean's part in case of very small pc's that t. 2 will still converge (?) Maybe not so simple!

It's $p \ln(p/p') + (1-p) \ln((1-p)/(1-p'))$ that \rightarrow 0. (p \rightarrow actual, p' approx)

for small p, p' $\ln((1-p)/(1-p')) \approx -p$ so $\approx p \ln \frac{p}{p'} + (1-p)(p-p')$: so $p \ln \frac{p}{p'}$ is 1.
 $p \ln \frac{p}{p'} \rightarrow 0$: well, p is very small, t. Q is how big is $\frac{p'}{p}$? $p \ln \frac{p}{p'}$ is 2 usually $< p$ i.e. $p \ln \frac{p}{p'}$ \rightarrow $p \ln \frac{p}{p'}$
 so again $\left[p \ln \frac{p}{p'} + p - p' \right] \rightarrow 0$ so it's not clear as to what's happening: $\leq p \ln \frac{p}{p'}$ $-p' \rightarrow$ const. ??

- 2 669.14 T. partitioning problem: Changes/improvements ① Add new partition, ② Modify old partition ③ Merge 2 or more partitions ④ Delete ~~the~~ partition (part of ② or maybe results of ① & ③)

In ④ we can sometimes delete a partition if its function is now done by one or more other partition functs. If we have S- (soft, stochastic) partitions, we may delete one if all of its PC's are very small & other partitions do have by PC's in comparison (for t. same softs.)

GPD vs. 01 - Discusses updating of d-Partitioning.

Within a partition, searching for good ~~partitions~~ P (prod) functions is somewhat easier. Contact for concs will include t. e. i.d. of t. partition, perhaps part of t. partitioning function that are closely assoc. w. this partition. [Partitions \equiv "Categories" = "Types" = "Kinds" = "Properties"]

Often, a new problem will have several "properties": Some will be more useful for partitioning than others.

T. entire system as I see it now, is rather easy to desc: A problem comes in; selection function

GPD (PST, Q, I) obtains a PD on PST's: T. index, I means that this is what we want from Q.

(we could, alternatively write GPD, (PST, Q) — where GPD, is a ~~at~~ different element of t. set of PD's that constitute GPD; But then we would not have a heavy of into "betw. t. sections (transfer of lang.)

Improvement: Next, the PST's will work (usually in (T.S.) on t. problem, Q.

They will obtain a PD over past answers, A_i. Is this pd 'part of' GPD? —
It so, what about t. PST's that worked on/solved t. problem?

In case of QA problems, t. answer will be simply obtained by direct applic. of GPD for updating acts ~~on~~ GPD so t. real is over. Updating can also improve items to max. 4/14/18.

In case of IIRV problems:

Clearly a better way to update is for t. PST to be able to call t. update func's decide



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 670.40 What is ^{25 fact} effect of updating it should emphasize. T. PST may decide that for their particular problem is amount of new info since last update, that it would be better to spend all problem time on problem directly. So: updating can be done before, during, ~~after~~ or ad libitum ← choice by PST.

It is perfectly OK for a PST to ~~simply~~ simply go to GPD (say a QA problem).

See 670.34-35: When a PST finishes a problem, it just stops! It can finish due to CP, or for INV probs when it finds a soln. We would have an operation mode for optzn problems (including IND prob), in which TCM periodically out-pub ~~the~~ value of Gorc obtained plus fav. Trainer can stop problem at any Gorc level desired!

well, say I can get it to work problems: Now, how does it (update/salt improve)? (19)

SN A [QA, AI] set can be integrated in any other such sets, or have its own "01" maximum 414.18

If 2 such sub-corp' are "mix" how do they share data costs? (ie, transfer fees).

Well, "mix" a "mix" means that we have a recept function that can validate digits in h before 2 sub-corp' & use different recept for them. Actual Hrs: system as a whole, has a single set of QAs & 414.18 is applied directly. The 2 P's for f, 2 sub-corp' can share data costs, in a common "data distribution module". — Actually, there is only one 01 function for both sub-corp' but this 01 is realized by many recept functions within itself.

19.10 I had Prof Peter Egan PST could be updated in its own particular way: but I think that ~~the~~ The update methods have a lot in common.

10¹² by c
10008 ≈ 1000
times

T. more update problems ~~for~~ for IND problems (INV probs can be formulated as IND

Probs in a way so that f. solns. are actually 2 (identical) L-shaped ^{21K} Prognosis INV problem, simply as a special kind of QA corpus. It ends updating to usual methods ~~for~~ L-shaped soln. for INV.

So, list problem types & assoc. PSM's. Then dev How each PSM is updated & how the system as a whole is

Does this help with problem of updating?
Guiding pd? → (28)

Updated (ie, SJ). → (672.00)

Re. (21-25) **T. IND** treatment of INV probs is not exactly E. Same as T. direct ~~IND~~ INV. method.

In c. IND ~~IND~~ treatment, we want a func that maps from problem data to soln. We want a type (short) func. Whether cc is larger or not irrelevant, but most methods of soln (say using L-shaped) must also give smaller. (The L-shaped gives max (pc/cc))

The soln. by IND may not be "too bad": ~~so~~ which suggests that we may be able to solve optzn problems w.o. 2 stage construction & needed Guiding PD's:

This 2 stage construction could be very time consuming — the 1st stage involves optimizing a 3 param (say, pd on cc for solns. max integration to each D.P. for "best" cond.) — ~~IND~~

The IND approach, while not exactly fine, may be OK. — it may actually be better because some cc is saved

Could we use a for approx. for general optzn. problems? Certainly not obvious! → Note 672.16



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: [SN] In Baran's 1993 paper (on translators - nonlinear approximations) he seemed to get faster

convergence than ALP does. This may be because he imposes constraints on the function being approximated, that are different from constraints I used in ALP analysis. He assumes that by ω components of f_i true hypersurface \downarrow at $z \gg z_0$ certain rate for $\omega > \omega_0$, or something like that. T. constraints I used were that the true system ~~was~~ was derivable by k continuous forms \rightarrow \downarrow if non-linear, was locally linear.

Notes: This is "Curve fitting" not Time Series Regression.

[SN] ^{Solovay & Solovay} ~~showed~~ ^{showed} that ALP will eventually make progress in only large PC for "0" (undefined) ^{In L.V. it does (collecting)} ~~in~~ ⁱⁿ conditional PC of 1 or 0 \downarrow . Can we also show it will eventually pinpoint by large ~~errors~~ errors in conditional PC of 1 or 0?

11:665.15 [SN] On "50% Solution": Say we have 2 machines w same cpus: (A) ^(A) ~~is~~ ^{is} optimum

ratio of internal to external problem time slow, (B) has 50% time share.

If we double cpu clock rate, it will be at least as fast as (A) since it will spend at least as much time on internal \rightarrow external problems.

Yet I had to realize that A might work problems much faster than 2x speed of B.

Thinking about it now, it seems that speeding up (B) by factor of 2 doesn't double its problem solving speed exactly. How did I get it that A could be much faster than B?

Say A spent 90% of time on S.I. ~~and~~ first, then 10% on external probs: which it then solved 200 times as fast

S.I. time is exponential in fraction external probs, so

woops! I assume time share, so A can't wait until S.I. is done before working external problems

Also must A can be faster than B is by $2^2 = 4$ (if S.I. brings exponential speedup).

Maybe not: exponential yield can beat \gg squaring by doubling S.I. time! Is this true? \rightarrow T. system will take

23: 611.27

Backtrack! How can Logical Reasoning (as a way to solve problems) fit in to 67k 27?

I have some recent good discn. of this problem. (I scanned 570.00 - 673.23, didn't find any references [I'm sure there is a ref. in those pp. (> 1 ref!)])

Anyway, I think y. idea was that TM would (via "logical reasoning" as part of its TQA:

It would eventually solve problems \rightarrow to those needed for "Updating" - i.e.

24: 670.02

On "J's method vs. AZAT! In J's method: It is easy for the system to narrow down on particular pm given it's PC \downarrow . In AZAT: If rows have PC's \downarrow we would allow machines T. only way to get by pc for a pm. is for that pm to be \geq Macro. My impression is that it would have to be used > 1 time (which may be impossible!).

Alternatively, AZAT could be allowed to have much context dependence of rules \rightarrow

End of this "Context dependence" can fix it so that only 1 pm can be generated

having PC $> .5$, say (or $> .9$ or $.99$!)

No More Query Failure
 \rightarrow as long as optimum system
to get as good a score as possible
Optimum System

"on J's method"
vs. AZAT
669.30
654.23ff
654.23
contact
dep. of
Kodak

2/27/02

There is a GTS: ID, 6/30/02!

data course, this should be

673

IO / SOY

SOY (Spurious Optimal Yield)

528

673 1/2



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:

: When one picks n points of n ^{individual} ~~independent~~ samples, the peak is obviously by, but it is \approx n \times largest expected individual (??)

Well, I suspected "No": It's a sample case: The population is in 2 parts.

(Set 1) has mean ϕ & $\sigma = 1$; (Set 2) has mean $.5$ & $\sigma = .1$.

One would usually pick a point from Set 1, which is bad!

00

So it may be best to pick the peak of $(x_i - k\sigma_i)$, where k depends on parameters of i ~~circumstances~~.

10

If σ_i 's of individuals are unkn. one can do better. The one can pick individuals who are in subsets of known σ_i , it might be possible to ~~do~~ do so.

Does there exist a function $f(x_i, \sigma_i) = \alpha_i \Rightarrow$ picking α_i w. peak α_i has expected error of zero? (or some other criterion)

For SOY betw. ~ 2 or 3 is 100. SOY doesn't change much, it's, say $\sim .76$ or $.86$

So \int picking i w. max $x_i - .76$, say, would not be bad in ~~many~~ for many ~~of~~ samples.

I want $f(x)$ to be \Rightarrow picking peak from subtracting proper SOY, gives least expected error.

$f(x, \sigma) = x$ is the default.

20

30

7/1/02
FD

SERIOUS DIFFY .11ff

PERHAPS one of "TOP Problems"



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

6/11.10.17

[SN] Out to GPD_1, GPD_2 (Quality Problem): Since the process needed for GPD_2 ($\rightarrow P \rightarrow PST_i$) is not large, we can do integrations rapidly, or do it. Can be merged w. only 2 or 3 "pts" for each PST_i . So probably the process would not take much cc.

.07
.08

In fact, it may be possible to estimate optimum precision to use! It will be when the expected G loss by using less precision is offset by the fact we have more cc to ↑ G in other ways.

Re: Dasirata: ① pc's sum to 1 or ≤ 1 . Eg. Such "CJS" estimate is meaningless.

② if $G \rightarrow f(G)$ (f is unknown ↑ w. only 2 or 3 values ... at least 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000)

Bad features: Inter-^{active} PC's of PST's is not captured. — unless we're rational. If PC's after would tend to Best ^{one} for a while; As it, several PST's may be very similar (say 10 of them). They all get 2 wt. of μ . (i.e. some works on them (1) (efficiency = .1 (!))

It would be better to work on Best PC one, Best for a while & do simultaneous eval's, or sequential alternative eval's; then note that others have had PC's reduced — also first PST has been worked out, & switching a not as quickly (This is WON approach which is certainly better than L such)

So, + Best way is WON! Using Best Prob vs. (cc, μ) curves for each (Problem, PST's) pair.

Note PC's is Won for 02 probs — which is different from Won for INV probs. Won for 02 only uses OR. — But interactive PC's make it very diff?

Bad feature of WON approach: It gives no CJS estimate.

A possible way out is to consider fewer PST's for problems (Humans do this). When we do use several, they are 2 indep or interactive simple way.

→ PST's that use trials that can be also used by a few PST's do pose an implicit and } to trial sharing problem
of problem: It is desirable to take advantage of this kind of interaction, but }
I haven't begun work on problem! There is certainly a vast no. }
of optimal user trials over consecutive trials (may be most such trials!) }
The way we do PC's would depend (partly) on whether this is a cheaper expensive

In .12-.13 times spent on 2 PST's is not entirely wasted, since it gives useful info for industrial corpus to help evaluate the value of PST's.

Consider a System to Solve QA's only; We have 2 PST's of 664.03 & 664.06

[N.B] In general, when "passage grammar structure" data" to update PST's & their PC's!
Grammars also considered do not allow prohibit dependence of PST's. This may be a very Normal characteristic of Grammatical Induction! — Just to PC's out set of PC's & their PC's. Consider "Grammar" of each language...



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

675.40 8) T. idea that cards could be PST's - this could be any ~~thing~~ form - including things better than L-trich, i.e. things to replace it. "Top Pgm".

9) T. scaling problems! 2) L-trich \rightarrow adaptive L-trich. b) % of success depends on context (or other ways to deal w/ pricing) c) \rightarrow low down to time spent ^{testing} ~~abstractly~~ (abstract model of "system function". (we now deal w/ this by breaking system func into parts, & by "structured testing" also ^(finding) ~~(selecting)~~ particular programs that are particularly good for ~~clabbing~~ a "system function")

10) T. new idea of Recog. Algs. helps w/ such a ~~thing~~ 9)

Other Pgm (10) - what ~~was~~ ^{nowadays} ~~was~~ EMCA 4.21.01 (74 marks ~~ago~~), \rightarrow
143.55 \rightarrow 144.00ff has an interesting possibl. formal method of ~~EMCA~~ General Prob. Solv.

The Update Criterion is not obvious hr.

11) \rightarrow Also, t. idea of how to de-use S-functions: T. input time was one ~~thing~~ that I worked on, but actually, TM will have many different methods to represent S-functions - The exact form will depend on the problem: hr. all methods should be interconvertible so I can compare results for updating. \rightarrow Note advance!

12) Application of System to button-type Probs:
A kind of Generalization of H's method: Say H's size of problem; Solver/^{random} problem in small n, using scalable techniques (techniques scalable to larger n), then apply method to larger problem. If small probs are larger & we may be randomly constructed & chances of ~~them~~ v. successful methods working for larger probs, is quite high!
Actually, this method is often close to one used by humans: we study how to solve small probs, but we only use ^{trial} methods that are scalable to large problems.
We may occasionally obtain method that do not scale (we have no proof that they will scale) but this will occur occasionally. To comb. t. method, mention use of hyper-heuristic digress to \approx "Prova" terms in Geometry.

[SN] New Invarson, APL, J programming lang: J is t. result of many yrs. work in looking at use, dialects of APL. I have several papers on this: Perhaps best is Invarson's history of it - 1991. See my file on "J-languages" - folder in "SS2". I'll have to be more familiar w. APL & J better & can get much from Invarson's history of J -
TGo it gives motivation for various features. I may disagree (much) w/ him on his ^{letter} ~~letter~~ reasons: But he may have unaware of!

\rightarrow Concious ~~the~~ reasons: But he may have unaware of!
Anyway, I want to use J/APL in TM, a/o have TM able to invent it -

Since TM is always moving toward "large notation".
Invarson's latest idea of "fact's Pgm" - this is functional pgm, which ^{"fact"} ~~"fact"~~ \rightarrow facts are not explicitly stated.
The top intro to J by H Darras is brill, good. I think he claims he can manipulate formulas as Arguments (like LISP).

7/7/02
ID



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20:676.90: There are ~ 200 "Primitives" in J: These are "wired in" funcs, probably there are only a few "Primitives" in the usual sense. T. Q is: ~~how many of J's "Prims"~~ how many of J's "Prims" would TM be able to ~~discover~~ discover? Also, would it dev some opress? Essentially, all these "Prims" are measures - near well founded but it is useful to learn, because they compress corpus. I want to make sure that TM has adequate means for discovering any/all of the 200 ~~Prims~~ "Prims".

20:06 So this is interesting! - That I have this "Grand Loop" = TM that should be able to acquire any ° of skill, "A few" "Get off to end". A I have the 200 "Prims" that are a kind of measurement of TM's heuristic power. A Q is: "How are these a better 'Measurement' from other skills: (e.g. ability to discover (linear/nl) regression in a real Time Series'?). I guess its because the "200 Prims" are closer to "primaries" - "near to bottom" - very "fundamental". We could, of course, wire in those 200 prims - i.e. give TM "J" to start. Would it then have an adequate set of "Hours"?

T. basic idea of the "TM Loop", is that whenever TM solves a problem, it updates Algum ^{present} # makes a small new into flow into the "next update Algum". IMPLICIT Inversion facts that = J's ability to do "implicit func" - i.e. define funcs w/o using Dummy args - is very imp. This amounts to being able to directly manipulate strings that can later be told to what then represent a function (see in 1.18: excessive quote)

3 IN GENERAL Re the "TM Loop" (206). 3 imp components: 1. Set of PST's, 2. The PST used for "updating" 3. T. TSO. During TM's Development, I Monitor (keep watch on) 1, 2 & 3. If 2 isn't growing properly, I have to put stuff in TSO to guide it a/o hand in improvements to PST's in General, & update PST in particular. The "Monitoring" is perhaps via the DFTY of the prims but TM can cancel work if "cannot" means maybe "can't", but KB is quite too far out.

35 AS IS, I'm starting w. Lsrch (w. various forms) for the PST of the Update Algum. - Also some (w) same PST for induction is prob INV prob (if any). I could add other Inductive methods of 663.06 ft ("regularity"). Also various partitioning methods ("Recapit funcs").

35 Then, there was the idea I had in mind in 589: Simple Adaptive Lsrch (no line during Lsrch). If the CJS was too large, I put stuff in the TSO to give "Bridge Concs" & guess what I was worried about in 35 was that certain things were not accessible this way, i.e. certain hours involved learning from sequenced trials.



30: 678.40:

"Good Scientist" does.

So it may be that there is a natural "logit" division of the updating process between After new data arrives, a "just before each problem is solved".

Hvr, the update after ~~the new data arrives~~, may be an artifact of human memory.

It seems to have to do w. (where / in what form / whether) the data is stored. ("whether" is perhaps of most import! — since humans probably do not store all input; two a machine could.)

Even w. machines, the Q is At what level of ~~accessibility~~ accessibility (cost) of data

(is its partial compression) should be stored — this is common to humans & machines

At first, I was thinking that GPD, would always be "completely" updated & revised

(Problem soln / ~~new data~~ input) — This is clearly incorrect. "completely" is meaningless, since CB is $< \infty$.

For Humans perhaps Much of "post data update" involving "looking for an Explain"

has to do w. expected future utility of such data/explanans. Also, at the post-data time, one has time to do the update, & it seems that one will probably have to do it eventually, so one uses that time for updating. This would seem to be

marginally relevant to machines: A machine should be using all of its time usefully. If the machine has some apparent "free time" it should be updating.

Hvr, what part/aspect of the GPD, (or "GPD₂") should be updated, is unclear.

It would seem that raw data coming in would be ~~the~~ a candidate for

update. Also (perhaps of most import) — one must compress the input data enough to know how much CC should be spent updating it! There is, in general,

this uncertainty about how much CC to spend updating something, because one doesn't know if its going to be useful. — ~~Minimum~~ So risk is involved:

One may find the update gave very little of expected value... so one has to guess or make decision about whether & how much to compress (i.e. update)

data ... that this decision is being continuously made as the data is being (processed / updated).

35: 679.24:

In fact instead of fitting to "crossion", or "clustering" ~~the~~ model to the

complete ("exact") GPD, ~~we~~ we will very probably try to express the GPD

int. first place as a crosscor or clustering or ~~equivalent~~ equivalent classes model.

JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:679.40 + u is (particularly f. ~~is a set of equivalence classes~~ "is a set of equivalence classes" model) could very simply
 (construction/computation/formulation) of f: GPD₁.

.04 { So 678.19, 679.35 ~~680.00~~ 680.00: looks like an adequacy problem }
 { Soln. tot. "correct" problem. } → 684.31

05 : Another problem perhaps much more important is J's method v.s. AZ141. Mainly it has to do w. adequacy of AZ141.

Perhaps a related Q: How adequate is 677.35-40? [conc net, "Bridg" conc ...]

Note: in 677.35-40 (con. net...) OK @ our time used a conc to solve a problem, it's possible will, but it may not be enough to warrant a "definition"... even when it is used 1 more time.

In general, I'm worried about the "Scaling" problem that occurred w. ANL at Saarb.

That ~~is~~ ^{for} ~~u = no. of concs~~ $u = \text{no. of concs}$, what pc of each $\propto \frac{1}{n}$. So solns w. k concs have $pc \propto n^{-k} = \frac{1}{n^k}$. Main ways [Prot of] to ↓ this effect: ① Breaking composition to 2 parts, so $n \rightarrow \frac{n}{2}$ ② have pc of conc dependent constant [① will be considered a "special case" of this. — that "Type" of problem was a kind of "constant".]

Say a (soln) problem involved 3 min (not so by pc) concs & some "connecting" concs — so for human to discover such a soln, heur. would be very fast: A good part of this speed is due to search analogies w. other problem solns & heuristics (Heuristic tricks. I really have to find more heur. search heur. for problem solving. The conc. net itself should tell mainly how to solve the problem, but other heur. are like catalysts to enable solns at "reasonable pcs"

So what I may need to do is first make the conc net, but see what other heur. "Catalytic" heur. are used by humans.

~~was going to say that Lurch uses additional techniques not used by humans in "Conc net"~~ I was going to say that Lurch uses additional techniques not used by humans in "Conc net": Not so! Lurch gets solns w. cc of $\frac{\text{conc soln}}{\text{pc of soln}}$. We prob. can't do much better!

0 Except that Lurch is definitely not optimum: WON can often give much better results! WON takes advantage of heuristics that are in f. GPD₂

[Perhaps heuristic weakness of conc. net analysis: Often, in working on a problem, various unsuccessful ^{finds} trials will suggest much better conc. Lurch does (such) heur.

This amounts to gross changes of pc of concs achieved by learning from conc. trials. Which is like: chosen by pc rand. try (it) from (what heur.) decide on one or more new conc. to try; go to next

7/9/02

681



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 680.40: So: 680.05-40 points up to need of a high level study of just how problems are solved — in conc. nets is "assoc. Heurs" E. act.

582: Conclusions: T. best way of deciding on which PST to use, ~~minimum~~

1) In all cases do GPD, first to get pd of G as a function of $(Prob, PST_i, CC)$

2) Do 678.19, 679.35, 680.00 > 2 soln. to "correct" problem

3) Or ~~if~~ if d.f.'s for the PST's are indep, I can use a W/O/N soln. (PST needs be developed, but old ideas at 680.05, is more recent ways would be useful for review.)

4) if assuming indep. of the PST's d.f.'s is poss. T. classic PST₂ approach, using

MT Carlo to make choices, is not bad. Hvr, see recent note: → .10R

.09-10: If may be possible to find θ (or 2) params of the GPD, d.f. of a PST: But characterize it w/len of so that (it's (param) success (increasing) param; partially If 2 param, partially order is use tables or more fine consuming methods for deciding order

Probably to must diff part of TM design is 680.05 A-T. design of TSCQ's: # Concepts of "Heuristic factorization" (finding set of adequate heurs (which may not be so easily expressed in the Conc. net!)

[The perhaps By Definition, all needed Heurs are int. Conc. net]

So, if a heur is a conc. int. net, then it would seem essentially that these conc. (Heurs) be context-dependent.

SN

I have this idea of "Heurs" (as components of the Conc. net)

is also idea of PST's — which are not part of a conc. net. Certain of the PST's (maybe 1 or 2 only) use L such as conc. nets. — Other PST's can use any derivable methods of solving problems. (Generally a "Heur" could "call" a PST!

It ^{probably is} may be necy to allow Heurs take PST's or parts of PST's, if we want all Heurs to be expressible in the system (since a "Heur" could be practically anything, a Heurs in the Conc. net seem

to be rather restricted in what they can do: ^(approximate) Modifying pd on ~~some~~ order to L such trials (?) — unclear as to whether the Heurs in the Conc. net can do all possible rearrangements of trials) — A even reordering of trials ~~doesn't cover~~ "Quick about" type of heurs.

What has to be done for TSCQ analysis is decision: To take various Disparate kinds of Heurs is see if I can fit them into the models of .21-31. If not, fix the models. What I really need is a better ~~simple~~ way of categorizing ("typing") heurs.



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 681.40: [uSN] answer My first about hours is "like", I have solved many ^{different} problems; & some, have certain features in common - I want to emphasize these features in trials for solution of new problems. "Induction" sound to be more like \approx to same problem - or, essentially only one problem w. no antecedents ("Max #19.18")

01: In fact "Max #18.18" has a "history", (= corpus of $\sum_{i=1}^n [Q_i, A_i]$ corpus $\rightarrow Q_{n+1} \rightarrow$ ~~next step~~)
It's not exactly r. vito ~~way~~ way to look at it. problem, hvr. T. ~~handwritten~~ demand
 \Rightarrow implied is that we use a system to that ~~we~~ used successfully in the past to obtain a new Q_{n+1} : C.O.D is not quite vito! T. corpus of pairs is

02: $[[Q_i, A_i]_{i=1}^k, \dots]_{k=1 \dots n}$ The implied Q_{n+1} def. for input $[Q_i, A_i]_{i=1}^{n+1}$ (i.e. f. different members of Q_i, A_i)
"A" corresponding to $Q_{n+1} = \sum_{i=1}^{n+1} [Q_i, A_i]$ Then action Q_{n+1} to match desired A_{n+1} d.f. 27

T. problem (for max) is: I don't see any "E-Bern Seq." to suggest induction! -
Little repetition of concs. - Nothing like to repeated use of various hours or various
<concs \equiv sub-functions \equiv sub-themes>

-18 } perhaps would split it, for each $[Q, A]$ problem, we had a recombinant function for each Q type & a separate $f_i(Q) \rightarrow A_i$ d.f. for each Q_i . We could then share concs betw f_i 's.

0 } Later, we try to "integrate" the f_i 's & modify to recognize points \rightarrow also \uparrow part corpus
22 } Earlier, I had ~~the~~ the dept of -18 \rightarrow -21 as a "Not so good, but is all edge field, try it" type of soln. - but it may not be so bad as \approx . usual general way is "first try" to \approx type of soln. We may want to get very many $f_i(Q_i) \rightarrow A_i$ functs so we can \odot all these p's \odot have alternative ways to do the "integration" of -21.

27 } Another Big trouble w. (102-112). Its just extrapolating to corpus; it's not (ordinarily) trying to maximize #19.18. If just looks for $(F(Q) \rightarrow A)$'s that are similar to the successful ones of f. past. - \approx is not a bad way of doing it, but it is 8

30 } Certainly not to Best \rightarrow 683.10 (who could TM realize that performing Optimization was the Best, for we could find another highest PC model for the past date!?!?) \rightarrow 683.10

32 } [uSN] res. (305 \rightarrow 31) This way be fine! The "Optimization routine" is a pgm. It will, indeed, get close to the desired corpus w. high pc (if option pms low high pc - even just to option idea - perhaps have option idea, be expressive) - TM doesn't have to understand what option means, in any way - yet it could use it as a way to (write/dcrb) pms i.e. a way to assign p's to pms.

7/10/02

U40 (1000-1) = 54/p.
+8 = 54/10 = 5.4

Line 21: Switched to Pen
7/8/02 #1

(for Stock 1.05) 611,000

582



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00.64 2.40 : Still not clearing my mind as to whether (682.305 - 1.70) is correct w/o the mechanics of how it would work.

TM could "try" to conc. of "optimization" after being given a set of examples of it. Ordinarily, it takes much cc to verify that something works as an example of optm w a certain CB. Theoretically, even tho the "optm conc" will have an initial (copy) long decn - w. env of SSE, it could be as good (explain/wide!) of a corpus. It's not clear, tho, that having TM try, what optm means is much better. Better than having it be a "locked in" conc.

08 Hvr. (E 682.305-40) (1.00 rat) would seem to be not very critical (or) putitive stock.

09 Note that a pure Q.A. TM write lra "in a way" to do INV problems in a way!

10 (682.30 682.31) Hvr. 682.305-40 has a more serious implication: Unless I show it works fine, I have trouble of 682.27-30 - that TM is not long (optm/updates INV) properly!

A way to (somewhat) deal w this: Consider a few (or one) method of finding best "414.18" maxzns: T. findings of "regys" in t. QA corpus.

15 "Corpus" for this problem is an "ordered" sequence of $[Q_i A_i]_{i=0}^k$ pairs $k=1/n$. The idea here is that the "total regys" for $(z=2+1)$ almost always includes those for $(z=2)$; it usually some new ones. In this sense, t. (code/decn) of 0^{k+1} is usually a smaller mod. of t. (code/decn) of 0^k . A Swallowish notion could be regarded as a "Mutation". So, t. Q is: "What's (t. best/av. p.) way to do this? (A extended cc best dependent Mutation)." } --

20 Most Generally: One has various possl. "Update Algs" (\equiv PST's). To find better ones: T. Gore Berg: "How much \uparrow M 414.18 for a given CB." (for comparison of some Algs, CB doesn't have much effect - for others, it is very imp: {for very large CB, one can do "Look ahead" & do "experiments" to get into.}).

25 Hvr, rite now, I want to simply Derb T. current update/inductn algm. - (not really looking to how to improve it!). I guess this is to show how "Contexts of knice entries" from how it differs from t. "Default" (i.e. $Az(14)$). -- 684.03

31 "Back tracking" in this model: To make truly for 0^{n+1} : creation of 0^{n+1} : western successful $0^{n+1-\alpha}$ Models (α is "amount of Back track") & we make trial "mutations"/additions. Trouble w. this: It 0^{n+1} was made by a 0^{n+1} mod & β $0^{n+1-\alpha+1}$ was made by a mod. of $0^{n+1-\alpha-1}$ (which is quite possl) than modifying $0^{n+1-\alpha}$ will not catch the particular soln, but $0^{n+1-\alpha+1}$ has to its corpus.

7/11/02
ID

684
684



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 682.40! So; for each O^m , we will have a backtrack index bt_m , that tells how far back it goes back to paths soln. Normally ~~to~~ to discover ~~the~~. O^m , we will ~~the~~

02 To Be continued.....

03: 683.31 [7/13/02] One way to modify O^m : (Simply Backtracking by 1, would mean I'd have to take at least as long as it took to find O^m plus more, because I have an additional

for Boston for or for
Wasa.
Even when she has to
Cops.
How how and how bay
(that will be a day)
(Milton round my neck
to life)
I now pronounce you Man
and wife.

Constraint: i.e. $O^{m+1}(A_{m+1}, Q_{m+1})$ must be large. One ~~that~~ is well modified!
TM "recognizes" Q_{m+1} as different from all previous Q 's; It uses O^m for all other
 Q 's $\hat{=}$ deriving a special O for Q_{m+1} . This combined algm $\in O^{m+1}$. I'd
need + additional reash
We do not deriving O^m in this "soln". We may then try to improve this O^{m+1}
by ~~the~~ ^{generalizing} (simplifying) (less constraint in) to recogn. algm, so it recognizes some
of the past Q 's — yet gives O^m results — also, perhaps modifying O assoc
 $Q \rightarrow A$ algm.

So 03 H. draw gives one meaning of "small change in O^m ".
Another possi. meaning is, perhaps, obtained by logical analysis of O^m — but
modifying it in certain ways will only change its response to certain Q 's.

The most imp't ^{ununderstood} Q about TM seems to be 682.00 — 684.17: Essentially: "Just how do we do
IND updating"? Perhaps More to it. Pt.: Just how does one or more of the relevant PST's work?
681.00 $\hat{=}$ a review of how TM decides which PST to use.

Updating involves (least el) Finding a good grammar relating Q 's to PST's ~~is~~
of P's of Gora ~~is~~ a function of CB. At first we ~~elaborate~~ this into getting a good
Set of PST's; $\hat{=}$ getting SPD to relate them to incoming Q 's.
.21 is on how to design a good PST for QA problems. — what I was mainly concerned w.
was a Rational for context dependent fc of concs: (also relationship to J's method.)

1: 680.04 [SN] It will probably be poss. to ~~the~~ design a (approx soln) to the "correl" d.f. \rightarrow it will be very
easy to tell when to drop a particular PST $\hat{=}$ which one to jump to. T. approx. form
of the joint d.f. could have more than a few parames, because ~~need~~ Raveris it much data,
if the models I would be using, usually drift up. (— so I couldn't afford many parames)
So this whole business of sequencing the PST's could be "very easy" — "very smooth"!

36 A ~~reasonable~~ reasonable approach to PST for induction: ~~the~~ ~~grammar~~ of ~~the~~ ~~text~~
 \rightarrow functions that have been used as O^m trials: ~~use~~ ~~use~~ G ($\hat{=}$ corpus PC) as part of corpus,
(See stuff I (relatively) recently wrote on GA for how to do this). From "G corpus", I can
get good trials in a Mt Goto way, or, because of the simplicity of the G-corpus,
I think this \rightarrow correct.

7/18/02

TM in "STABLE STATE".06

685
685



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

50: 684.90 It may be able to (appropriately) integrate from quality study: Do go on to try Deer: it looks quite reasonable

51: T. idea of 684.36 is simply my realization that IND problems should be solved some way as any other problem (INV or OZ). This seems to solve one of last remaining outstanding TM. probs (other than TSCQ!).

52: There certainly are other problems (e.g. see 611.00 ff @ "Steel"), but at present, the systems seems "stable" enough to write rather than draft "report". (But do note 2.1)
53: 667.26 is recent beginning of discussion of updates of IND functs. — If it follows thru to 685.05.

54: T. data of soln. of the IND problem is a recursive data. The "boundary cond" is induction using trials of θ in p order w. Least cutoff: Note that this is not "Lsrch". Lsrch always has possibl. of suby good soln. This "Boundary Cond." does not have that possy.

55: It might be best for Exposition, to describe a pure induction (Q&A problem) TM in detail — Giving some reasonable initial "boundary values" for various types of IND problems of interest. It should be possible to write TSCQs for such a machine. IND & OZ prob capabilities could be added later — but these "recursive

56: IND Engine is what + main "driving force" in the system.
57: NB: in 684.36 ff, we realize that a {GPD, GPD₂} approach is necessary for IND probs, & that it is probably relevant to the evaln. of p.c. of Cond w/ context — But the exact details of this "soln." to the "context evaln" problem has yet to be described!

58: (7/21/02) AH! Re: recent work on "correlations" betw. Lsrch trials ("conds")
59: First GHT says if trials are indep (i.e. a constraint on cc's of trials), then $\frac{P_c}{cc}$ order of best. If trials are not indep, then and joint PD is known, the best path may be diff't to compute (NP hard?). If we cannot do better tricks I discussed about "correlns", this could

60: give "long" during Lsrch a \uparrow speed of soln. considerably: maybe dozens w. Lsrch.
61: 2 ways we can have "long" during Lsrch: (1) T. Joint GPD over conds remains constant during Lsrch (Pro New "correlns" data not taken advantage of.
62: (2) T. GPD is actually modified during Lsrch, as a result of what happens during certain of the trials. : It may be possibl. to assume conds p.c.'s are indep, & do (2) & this may simulate (1) to some extent (?).

63: (3) An approx. GHT (22) for correlated conds: pick cond w. max (incremental) $\frac{P_c}{cc}$.

64: The (known) GHT: we have to Joint PD on conds.
65: Woops! In the original GHT, there was only one prize. Is this a necessary cond. for

7/25/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

CO: 685.40 Max $\frac{p_i}{c_i}$ to be "Best choice" soln.? I would guess not. In which case a joint PD on rents is poss. & meaningful. So consider -1.37-38 / As soon as we make a trial & it fails, we have a

now (Marginal) P.D. We can pick cand. ~~max~~ $\frac{p_i}{c_i}$. (we know c_i 's in advance)

This is a "Max Greedy" Approach. A less Greedy approach looks to trials in the future,

This is a Dynamic Prog. Problem.

Consider $k=2$: We chose our first cand: prob of win is p_1 ; $c_1 = c_{c1}$

So prob of win with $c_2 = c_{c2}$ is p_2 : If c_1 loses; its ~~loss~~ following Least cost p_2 , c_2 follow.

So cost Then prob of win is $(1-p_1)p_2$ at total cost of $c_2 + c_{c2}$

I pick we chose cand of max $\frac{p_i}{c_i}$

So p_2 of c_2 win cost $(1-p_1)p_2$ prob of ~~win~~ total win cost $c_2 + c_{c2}$.

I have 2 trials! At the end of Round, I know certain prob's of win; certain prob's of total c_i .

— which trials to do?

We have choices trials: Objects: Each consists of several poss., a sequence of various p_i :

Say $\left\{ \begin{array}{l} p_1, c_{c1}, \text{win} \\ p_2, c_{c2}, \text{win} \\ p_3, c_{c3}, \text{no win} \\ c = (1-p_1-p_2) \end{array} \right\}$

v.s. $\left\{ \begin{array}{l} p_1', c_{c1}', \text{win} \\ p_2', c_{c2}', \text{win} \\ p_3', c_{c3}', \text{no win} \\ c' = (1-p_1'-p_2') \end{array} \right\}$ $(= c_{c1} + c_{c2})$

Which is better?

Maybe $(p_3, c_{c3}, \text{win})$ is a third possible choice (?)

The Ris Problem Seems Very Imp't! I really don't have to solve it just now! Put it out Stack: 611! Also Note: It is not exactly what IM needs to solve; 687.3

An outline of the Paper:

Introdu: Gives decs of types of problems solved. INV; OZ probs.

Derb. END prob (QA): show how it is special case of OZ problem. Say that it is a special case of OZ, but don't explain how, yet! Give ref. to Factor where it is discussed.

A nice way to organize the paper would be to follow: Abstract: Derb. INV prob. in detail; OZ probs: Give lots of examples. Then discuss QA probs as imp't kind of OZ prob: Give 419.18, w. explain. but no more detail. I.e. don't tell how to get approx of OZ. 687.00 Perhaps Just go into material of Abstract, in intro & 1st sec.

Then, Then, in intro, I list contents of following sections (or New Paper Appendix).

If c_i costs are divided up into several "equon classes" of correlated c_{c_i} , then a good that would be to pick the best guess from each c_{c_i} class. Each time it fails, we get rid of a lot of trials of low c_{c_i} . This is a practical approach in Real TDM as well as a poss.

10

16

20

25

30

36



00- (686.36) : Dcrb. GCPD: Henry General cond P.D : Give Google as ~~an~~ example.

Explain how it contains v all the "knowledge" in TM. (if we include PST's as part of GCPD).
Perhaps discuss "index nos." Part (sometimes) tell idk what "kind" of response is wanted.
Perhaps be a bit vague at first about how f. GCPD figures out what we expect of it.
Plan introduce "index nos."

Unclear as to how much to put in intro. It should include intro from Abstract, but give
Some more detail — also Give section no. Where more detail is given.

Intro: Start w. history of project (86, 89, 94): Give general goal of a very useful
~~Statement~~ "Scribble's Apprentice" The convergence theorem of Alg. probly
if there were
suggests that ~~with infinite time~~ no limitations on computational resources, we could ~~obtain~~ easily design
~~an extremely efficient~~ a machine that would be extremely good at prediction —
at least as good as any existent practical algorithm. ~~Natural~~ obtain

One of the main goals of the present project is to ~~find how good~~ we could do
~~with limited computational resources~~ determine the theoretical and practical
effects of limitations in computing resources. It should be noted that this work

is only distantly related to computational learning theory, which is concerned with the system's
~~asymptotic~~ behavior for ~~very~~ large problems. Most of the solutions
we find for problems ~~are~~ ^{have} computational complexity that is exponential in
problem size. ~~These problems are~~ ^{The system works} ~~excluded by~~ ^{modifying its vocabulary} ~~the system's~~ ^{computational} ~~complexity~~ ^{complexity} ~~is~~ ^{is} exponential in

problems to be solved ~~are~~ ^{are} small. It is our impression that ~~people usually solve problems~~
~~rather than~~ ^{is} the way that people usually solve problems — by reformulation,
~~rather than~~ ^{extend} ~~iterative~~ ^{iterative} search (maybe reference to Minsky
"Society of Mind").

68920

1: 686.25 **SN** Toward a ^{criticism for cond} soln: say ~~that~~ "r" won; want to win in k trials, but
want to order all k cond. sequences ~~and~~ and up in "r", \exists expected
 \ll is min. (Q: for original problem: I think > 1 soln. (is ~~it~~ winning cond)
is ~~probably~~ logical. — does this ~~apply~~ in original GHI proof? — does it ~~apply~~ in
present case?)

A poss. soln. to the problem: Say we decide on a sequence of M cond. to try,
w. understanding that we quit as soon as we find a "winner".
Then the M cond: we have $\{c_1, \dots, c_M\}$; $c_i = 1/M$. ~~Let~~ $c_i \leq 1/M$.
 c_i is the "name" of the i th cond in the list ("Names" = $1/M$).

7/28/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

10

Say c_1, c_2, \dots, c_n is a possible list.
 For that list let p_{c_i} & c_{i+1} be the p_c & c assoc. w/ card c_i .
 Then, for each adjacent pair of cards, we could interchange p_{c_i} to obtain a new list: $p_{c_i}, c_{i+1}; p_{c_{i+1}}, c_{i+2}$ are the original p_c, c pairs.
 $p_{c_{i+1}}, c_{i+2}; p_{c_i}, c_{i+1}$ are the ~~interchanged~~ pairs after interchange.

I suspect that interchange would ~~be~~ give a lower "expected" future cc loss, i.e.,

i.e. $\frac{p_{c_i}}{c_{i+1}} \cdot \frac{p_{c_{i+1}}}{c_i} > \frac{p_{c_i}}{c_i} \cdot \frac{p_{c_{i+1}}}{c_{i+1}}$ i.e. $p_{c_i} \cdot p_{c_{i+1}} > p_{c_i} \cdot p_{c_{i+1}}$

10

But first look at the ~~convex~~ convex.

For original c_i vector of 68739 expected cc before coin is:

$$\sum_{j=1}^M \left(p_{c_j} \prod_{i=1}^{j-1} (1-p_{c_i}) \right) \left(\sum_{i=1}^j c_{i+1} \right)$$

$\left. \begin{array}{l} \text{prob of success} \\ \text{on } c_j \end{array} \right\} \left. \begin{array}{l} \text{quantity of future loss for } c_j \end{array} \right\} \text{cc of succession } c_j$

} want to Minimize this.

(Pens not to product of inter p_c : i.e. to product of cond. p_c 's.)

14

I think I used this proof method for the original GHTI, but ~~it~~ was for Goucas complicated looking? There are only a few terms in the \sum of it that change

20

when we interchange ~~of~~ c_i & c_{i+1} . for \sum_j ~~is~~ I think ~~is~~ ~~is~~ $j=i$ & $j=i+1$ that are critical. c_j, c_{j+1}

For the original order, we have

25

$$p_{c_j} \prod_{i=1}^{j-1} (1-p_{c_i}) \left(\sum_{i=1}^j c_{i+1} \right) + p_{c_{j+1}} \prod_{i=1}^j (1-p_{c_i}) \left(\sum_{i=1}^{j+1} c_{i+1} \right)$$

for the permuted order:

28

$$p_{c_j} \prod_{i=1}^{j-1} (1-p_{c_i}) \cdot \left(\frac{\sum_{i=1}^{j-1} c_{i+1}}{A} + c_{j+1} \right) + \left[p_{c_{j+1}} \prod_{i=1}^j (1-p_{c_i}) \right] \cdot (1-p_{c_j}) \cdot \left(\sum_{i=1}^j c_{i+1} \right)$$

$= B + c_j + c_{j+1}$

30

recursion: $p_{c_j} \left(\prod_{i=1}^{j-1} (1-p_{c_i}) \right) \left(\sum_{i=1}^{j-1} c_{i+1} \right) + c_j + p_{c_{j+1}} \left(\prod_{i=1}^j (1-p_{c_i}) \right) (1-p_{c_j}) \cdot \sum_{i=1}^{j+1} c_{i+1}$

↑ quite similar!

.28 $\Rightarrow p_{c_j} \cdot A \cdot (B + c_{j+1}) + p_{c_{j+1}} \cdot (1-p_{c_j}) \cdot A \cdot (B + c_j + c_{j+1})$

.30 $\Rightarrow p_{c_j} \cdot A \cdot (B + c_j) + p_{c_{j+1}} \cdot (1-p_{c_j}) \cdot A \cdot (B + c_j + c_{j+1})$

.30 $\Rightarrow p_{c_j} \cdot A \cdot B + p_{c_j} \cdot c_j \cdot A + p_{c_{j+1}} \cdot A \cdot (B + c_j + c_{j+1}) - p_{c_{j+1}} \cdot p_{c_j} \cdot A \cdot (B + c_j)$

While .14 is pretty nice, I'm not sure about .20 ff! Are p_{c_j} & $p_{c_{j+1}}$ terms the only difference in the way just .25 & .28?

In the original GHT proof $p_{c_j} = p_{c_j}^1$ & $p_{c_{j+1}} = p_{c_{j+1}}^1$: I think case (.28 = .30) - so something's wrong! Also note, that about my comparison of



8:20 P

00 : 688.40 \approx 528.14 fact: \geq cases was wrong! for $J < J_0$ (say $J_0 \in C_{j_0+1}$ where j_0 is the change)
T. sum of 528.14 up to $J = J_0 - 1$ are the same, but for $J > J_0 + 1$ the $\prod_{i=j_0+1}^M (1-p_i)$ products
differ by $\frac{(1-p_i)(1-p_{i+1})}{(1-p_i')(1-p_{i+1}')}$.

04 : Hur, in spite of 688.40 error, -1.14 does seem to be correct for both the order
5HT is for γ new (i.e. mt. new situation in which P_{ij} 's are dependent,
we still have to choose a C_i ($i=1/M$) sequence a priori. -1.14 is the Gouze was made for
Minz. (Pro as written, it is post). But there is no case which wins. $-P_{ij}$'s
may be \geq deard Gouze!) So Minzing -1.14 looks like a desert formal
Soln.

11 : Hur, (NB!): mt C_i sequence (vector), each P_{ij} is, in general, a function of all C_i 's
that have come before it: (i.e. we know that all of them field. So in this general
case, the problem seems much more diff, the minzing -1.14 still seems like
a formal soln. We can assume each P_{ij} depends only on C_{i-1} , but P_{ij} is a very
definitely, simplified approach... \rightarrow 695.30

20 : 687.30 : Re: "introduction", I had considered giving more motivational background —
like saying that the system is meant to be toward a realization of "STRONG A.I."
Hence this would be best put into the Talk rather than the introduction.

25 : This case
no many
flawed
intro.
perhaps General idea of 687.09.30: That ALP is known to be U.g. w. $CB=0$.
What's best strategy for $CB < \infty$? Or can we make U.g. (not nearly optimized)
strategies for $CB < \infty$? Discuss 586, 89 is the "factor of 4" question.
For my self: I may want to discuss in some detail, the "factor of 4" of
is just why it wasn't exactly like (i.e. not all hours could be expressed as static
Modif. of γ GCPD). But many could: Some hours involve changing Pd doing such;
are not (such. (No CJS estimate is easy (poss.)). But any hour can be expressed as PST.
Then explain factor of 2 in "50% solution" ...
 \rightarrow Hur, If ALP is so good w. $CB=0$, there are obvious ways to test finite CB , Rank +.
TD system I've devised!

28 : So perhaps emphasize this system as an attempt to follow the path of Humay (exp.
Or: Just Dec. Sol 86-89-98: Then say this is continuation of that work: Correction of errors,
development of parts not never made completely deducted.



The present paper will be a continuation of that work, ~~bringing the system~~
in certain cases, replacing conjecture by proof, and bringing the system
much closer to practical realization.

In the form of a simplification and formalization.

The system is meant to ~~imitate~~ operate much as humans do. When the input
system is first presented with a problem,

The ~~initial~~ ^{initial} state The information in the system (knowledge) is contained
in two data structures: The General Conditional Probability Distributions (GCPD)
and the set of Problem Solving Techniques (PSTs)

The GCPD is able to take as input a great variety of strings, numbers or ~~other~~
other data types. Its output is a probability distribution on ~~the~~ strings number or
other data types. ~~It is~~ ^{It may be thought of as a very general}
web search engine - like Google. ~~Search engines~~ ^{Search engines} usually have as output
a list of web pages in order of relevance to query. The GCPD
will have ^{output a} list of possible responses to its input; each with an associated
probability.

The set of PST's at the beginning of system life are inserted by the
system frame.

(GPS)

Though the system is meant to be a general problem solver, ~~the~~ ^{the} kinds of problems
it solves and the techniques it uses are well beyond those described in

Newell and Simon's pioneering work (see reference ~~to~~ ^{to} GPS). It also differs from ~~the~~ ^{the}
in that few problem solving techniques are built into the system. It ~~adds~~ ^{adds} ~~new~~ ^{new} methods
~~to~~ ^{to} ~~the~~ ^{the} ~~system~~ ^{system}. The problems ~~fall~~ ^{fall} into two ~~or~~ ^{or} classes: ~~all~~ ^{all} ~~various~~ ^{various} ~~problems~~ ^{problems}
and time limited ~~optimization~~ ^{optimization} problems. ~~While~~ ^{While} ~~all~~ ^{all} ~~problems~~ ^{problems} ~~are~~ ^{are} ~~not~~ ^{not} ~~all~~ ^{all}
problems fall into these categories. We feel that the system is sufficiently general so that it can learn to solve ~~other~~ ^{other} ~~levels~~ ^{levels}
of problems as well.
One approach: ~~Make~~ ^{Make} ~~intro~~ ^{intro} ~~abstract~~ ^{abstract}, but simply ~~add~~ ^{add} ~~material~~ ^{material}. Perhaps
start out w. $(.00)^t$ or $(-1.33)^t$

~~Do~~ ^{Do} explain of "if all individuals p" "F.N. of '89" ~~that~~ ^{that} ~~why~~ ^{why} ~~they~~ ^{they} ~~are~~ ^{are}
"No + Bad", ^{is} ~~but~~ ^{but} ~~wrong~~ ^{wrong}.

on to ~~the~~ ^{the} ~~time~~ ^{time} ~~varying~~ ^{varying} ~~problem~~ ^{problem} $G(x) = G(x), f(t)$: One has (usually) a finite
set of PST's to ~~max~~ ^{max} $G(x)$. GCPD, gives us a ~~pd~~ ^{pd} on G values for ~~the~~ ^{the} ~~problem~~ ^{problem}
~~is~~ ^{is} ~~of~~ ^{of} $CB = t$, for each PST. We ~~root~~ ^{root} ~~just~~ ^{just} ~~select~~ ^{select} ~~the~~ ^{the} ~~PST~~ ^{PST} w. ~~best~~ ^{best} ~~expected~~ ^{expected}
value of $G(x) + f(t)$. Using the ~~cost~~ ^{cost} ~~function~~ ^{function} (criterion) we choose a PST,
and work on it for time T_0 . We ~~then~~ ^{then} ~~compute~~ ^{compute} ~~as~~ ^{as} ~~to~~ ^{to} ~~whether~~ ^{whether} ~~we~~ ^{we} ~~abstract~~ ^{abstract} $G(x) + f(t)$ to ~~the~~ ^{the} ~~next~~ ^{next}
further work on this PST. - If not, we stop. We may produce our best x

7/30/02

Bibliography

691



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

ACM Aug 2002: 2 ^{papers} art. of int.:

1) Data Mining: pp 28-59 7 papers: Of much interest. Explain various
(compression), pattern schemes.

2) "Swarm Intell." p 62: Explains ~~with~~ how it works! Various forms!

"Ant Colony" act.

Both have lots of References.

7/2/82

805

692



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:640.90 Thus far" or go to another PST if it's an expected f in E present $G(x) f(t)$.

(This all assumes that these decreasing calculations take zero time)

I suspect that this may also solve the "Chess" problem. In fact, it may be a subset

\Rightarrow my most recent ideas about "Chess".

It is likely that any problem solvable by a human can be solved by this system —

If we can get past human's incoherence into a "PST" in the system. I assume that this is always possible — but the requirements TSQ may be difficult to derive.

Mainly, I feel that after TM has done a lot of problem solving of the type of Trener

improving the "ECPA, [PST]" in view of past experience. It should be able to improve PST's that I don't understand, & PST's that humans don't yet know about.

[5N] about the "Convergence Rules" for Bayes & QA:

First do the KL form of sum for Reun: From this, we easily get ① the error form on bits.

② For KL form on (elements of Bay) or on "A's" From Reun (last) \rightarrow corr² form Bay elements or on A's: This last corr² form is not as good as bit format ① but it is equivalent to save.

③ It is actually of interest. $E(S/P_{true}) \cdot \ln(\frac{P_{estimate}}{P_{true}})) < \text{constant}$.

[5N] I was concerned w. deciding which of the non-optimum codes or models to store — expecting future "Backtracking" needs.

Remember ① would save a max of X2 in cc?

② A possibly good way to do this: to "cluster" the codes into classes that are "similar" — (There will perhaps be a "distance" function betw. classes.)

For pruden, we use the best class usually (maybe add some others if they have env. info) for Backtracking, try classes different from the formerly best class.

We want a class that has family by PC, yet is different from the formerly best class that has been rejected by x. new data.

These "Classes" can be obtained by "Clustering" Algs.



20:692

The ^{research} ~~system~~ ^{was} with described is a continuation of earlier work on making (dynamic) (S 86, 89, 94). While the present ~~system is~~ ^{system is} in general ~~similar~~ ^{similar} to the earlier ones, we ^{have} worked out ~~many~~ ^{many} critical ~~aspects~~ ^{aspects} of implementation ~~resulting~~ ^{resulting} in a simpler, more general ~~description~~ ^{description}.

We will describe in ^{some} ~~much~~ ^{more} detail, as parts of a paper that were not ~~clear~~ ^{clear} in the earlier papers.

The general operation of the system has ~~changed~~ ^{changed} much, but our understanding of its operation has. What we ~~will~~ ^{have} ~~presented~~ ^{presented} is a ~~simpler, more general~~ ^{simpler, more general} system ~~description~~ ^{description} clearer presentation of a system that is at once simpler ~~and~~ ^{and} more general — with ~~much~~ ^{much} better understanding of the ~~system~~ ^{system} to which the system is ~~applied~~ ^{applied}.

One thing that ~~has~~ ^{has} ~~changed~~ ^{changed} ~~much~~ ^{much} ~~is~~ ^{is} ~~the~~ ^{the} ~~order~~ ^{order} ~~of~~ ^{of} ~~the~~ ^{the} ~~description~~ ^{description} of the system is ~~optimization~~ ^{optimization}.

Q: How is present system different from S 86, 89, 94? ~~Some~~ ^{Some} aspects of the earlier work was not clearly ~~defined~~ ^{defined}, ~~and~~ ^{and} much of the recent work has ~~been~~ ^{been} ~~done~~ ^{done} ~~to~~ ^{to} ~~work~~ ^{work} ~~out~~ ^{out} ~~the~~ ^{the} ~~details~~ ^{details} of operation in a ~~clearer~~ ^{clearer} ~~simpler~~ ^{simpler} way ~~been~~ ^{been} ~~forwarded~~ ^{forwarded} ~~an~~ ^{an} ~~understanding~~ ^{understanding} of its limitations.

Much of the recent work has been in

Maybe not better in this: few people will be even mildly acquainted with S 86-94.

So: I have indicated 13 P nos. in abstract (L.A.B.)

Modular, for introduction: P 2: Give roots (S 86, 89, 94). ~~remove last line, insert after P 4!~~ P 4.5

Then say: "The generality of the system in that it likely ~~that it can be~~ ^{that it can be} ~~applied~~ ^{applied} ~~to~~ ^{to} ~~bring~~ ^{bring} ~~it~~ ^{it} ~~beyond~~ ^{beyond} ~~its~~ ^{its} ~~limitations~~ ^{limitations}"
 beyond these limitations

P 3 ENV: Give formal Model

P 4 Q2: Give formal Model. In § — we will show an important kind of induction ~~problem~~ ^{problem} ~~can~~ ^{can} ~~be~~ ^{be} ~~transformed~~ ^{transformed} ~~into~~ ^{into} ~~a~~ ^a ~~time~~ ^{time} ~~limited~~ ^{limited} ~~optimization~~ ^{optimization} ~~problem~~ ^{problem}.

resources of the present system are devoted to solving induction problems of this kind.

P 5.5 All of the "knowledge" in the ~~present~~ ^{present} system is contained in two data structures:

1. A set of Problem Solving Techniques (PST's) and a General Conditional Probability Description (GCPD). A PST is a subroutine which, when given a problem description, will attempt to solve that problem. The GCPD ~~is~~ ^{is} ~~a~~ ^a ~~matrix~~ ^{matrix} ~~system~~ ^{system}, may be

thought of as a very large ^{very general} ^{probabilistic} information retrieval system ^{such as} Google, ~~and~~ ^{and} ~~most~~ ^{most} ~~knows~~ ^{knows} ~~what~~ ^{what} ~~is~~ ^{is} ~~going~~ ^{going} ~~on~~ ^{on} ~~in~~ ⁱⁿ ~~the~~ ^{the} ~~world~~ ^{world} ~~at~~ ^{at} ~~any~~ ^{any} ~~one~~ ^{one} ~~time~~ ^{time}. ~~it~~ ^{it} ~~doesn't~~ ^{doesn't} ~~know~~ ^{know} ~~the~~ ^{the} ~~probability~~ ^{probability} ~~values~~ ^{values}.

8/1/02

694



P9 say Lsearch will be discussed in section —

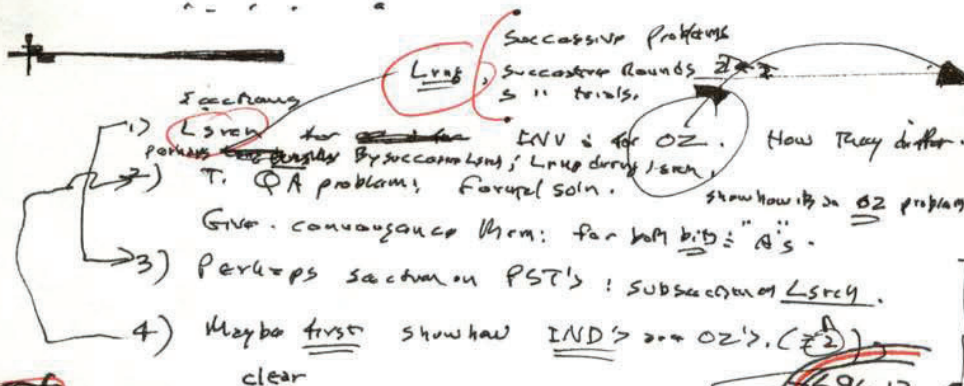
L search uses a probability distribution to guide search for a solution to a problem. Initially, this probability distribution will be the GCPD. However, simple Lsearch ...

As the system evolves, it ~~will~~ ^{may} modify Lsearch and eventually replace it with the better PST's. ~~This is discussed in Section~~ ^{Makes} Section — (on "How Good is the System") ~~discusses this~~

P10 add last line: perhaps in abstract as well as method.

P11 follow w. discuss. of Skuneman Lsearch (small cjs). — Use as large cjs
→ system can afford: gives better search capabilities.

P12 section ... discusses "differences" in the 2 approaches?



SN in General Optzn problem, we want to Maximize some function of C & G. We need linear ordering of C, G pairs using first CB is one way of ordering. In Anytime problems, I'm not sure as to what f. criterion is! Another party: Minimize to achieve a certain G-level.

21 **SN** Is there a ~~clear~~ ^{Big} ~~Goal~~ for the General update of (GCPD, [PST's]) — ? what's Goal for a PST?
AH! Not seeing to be a ~~Big~~ ^{Big}; Only GCPD, it directly updated. This part of the system doesn't care if PST's are wrong. all it is interested in is predicting how good they will be w/ a particular problem. **See 69500-07 for likely soln.**
What would motivate to updating process to divide v.p.

- ① QA: how many are OZ probs given cur. terms, act.
- ② PST's: a) Lsearch & variants of it is an induction problem. b) other methods. Mentions GCPD much
- ③ Updating.

PST's? — Also, what's the Goal of a PST?
We ~~do~~ really need intermediate comparison of PST's:
The idea is that we want to keep a new PST if it is better (for at least one problem) than any other PST (Not exactly, we also have seen how well GCPD does us problems for it).
So a new PST may not be "better" than any other, but it may be a desirable test (when that PST is appropriate for a problem → when new PST's are first introduced/invented, they are usually only known to be applicable to only a few probs — later we expand their utilities. We may ask "why" are they good? — "When" are they good?
So maybe we start out w. a particular problem, & "invent" PST's for it.
Normally, I have no over-all Goal for the system. For **INV probs**, I can always tell if a PST is any good on a particular problem.

69500



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: Look at "CPM" stuff (both of them) to see how parts fit together.

Also, may be use & use of their expos. by SFTS

03: SN on 4/4/88: I was going to use SFTS cor to prove this:

Now I'm not so certain about SFTS cor! It says nothing about the approx of P(x) approximating P(y).

It could be purely ad-hock

On the other hand, if $\prod_{i=1}^n P(x_i | x_{1, \dots, i-1}) > e^{-kn} \prod_{i=1}^n P(x_i | x_{1, \dots, i-1})$
Then $E(\text{error}) < \frac{k}{2}$. This is true, indep of P's & conv. length

10: If guess the badness of an A.H. P.D. is that it is unlikely to remain P(x) small in the future.

If I put in a prior factor, the P.D. amounts to a coding on the reference UMC, and we can compare it directly to the code assoc. w. the "generating P.D. (P.C.)."

If we want to compare 2 cond. p.d.'s w. the true p.d. (P.C.), it may be necessary to put in a prior constant.

Likely that the soln. is what I talked about in "The Discovery of ALP": That SFTS gave a v.f. so that regarding "best" that one could get, using stated approx. — It was "sum of all codes";

That any finite approx. was < that value, but the larger it was, the closer it was to the correct value.

While this is true for the closeness to the "right" semi-measure, there are interesting in the ratio of two SFTS-measures, (this is "familiar territory" — "I've been there before" — what was conclusion (if any)?)

How, taking logs: If A' is closer to A than A'' is, then A'-B' need not be closer to A-B than A''-B'' is.

20: On the other hand: The statement "The Discovery of ALP" did not involve SFTS. It was for semi-measure only. It's bigger better.

If we use the "convergence firm" we get a bound on error.

Actually, the conditions for the corollary are more stringent than I put:

23: $P'(x_i(n)) > e^{-kn} P(x_i(n))$ is $1/2^n$ $x_i(n)$ is a string of length n.

So 23 is true for all strings of length n. Superficially, it would seem that we don't need to be that general! — But 23 only has to be true for a small fraction of strings of length n.

On second thought: the corollary is about expected value so it probably has to cover all " " " "

30: Say $\mu(x)$ is the true generating p.d. $P(x)$ is a universal p.d. : so $P(x) > e^{-k} \mu(x)$ for all x. w. finite dom. wrt. μ

31: Say $P(x)$ is a p.d. that one has found. It has dom $\mathcal{C} = \alpha$ & it assigns $P(x) = \beta$ to a certain corpus.

$P_M(x)$ has dom $\mathcal{C} = 1$, & assigns β_M to that corpus.

33: I expect that $\beta \cdot \alpha < \beta_M < \mu(\mathcal{C} \text{ corpus})$

$\beta \cdot \alpha \leq \beta_M(\text{corpus}) \leq \mu(\text{corpus})$. Actually $\mu(\text{corpus})$ can be any ϵ (i.e. ≥ 0), but for a large corpus, it will have a certain "entropy" — (mean in p paragraph) (usually).

$P^0(x)$ will always have < all of codes of $P_M(x)$.

Any code for a corpus can be regarded as a dom of p.d. + dom of corpus in terms of p.d. (2 part code)

So any p.d. on a corpus (like P^0 on $(-3, 3)^d$), will be a subset of codes of $P_M(\cdot)$

8/7/02

Texas Tokyo
NYC Hiroshima
I/RAQ N. Korea

G-M Toyota

038



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

So it will love to assign a lower β α than β_m (-1.33), i.e. $\beta \alpha < \beta_m$.

If $P^A(x)$ has ~~some~~ $\text{demp } p_c = \alpha$, it assigns β to corpus.

$\alpha, \beta > \sqrt{-1.33}$, then $P^A(x)$ would ~~be~~ better than $P^B(x)$ in predicting x corpus.

The big idea here is that $\alpha \cdot \beta$ is x d.f. of x d.f. by the "pair" $P^A(x)$ to the corpus.

For this corpus, at least, larger $\alpha \cdot \beta$ means that x is closer to that of $P^A(x)$, universal d.f.

But say $\alpha_1 \cdot \beta_1 < \alpha \cdot \beta$ but $\beta_1 > \beta$. So, for the corpus above $P^A(x)$ looks better than $P^B(x)$.

Ah! For $P^A(x)$, we should use α, β for $P^A(x)$ & α_1, β_1 for $P^B(x)$. P^A & P^B combine linearly w. wts w & w_1 resp. - for corpus of length ϕ . As we move thru the corpus, w changes as x and w_1 stay the same like (1.10)

However, no matter how good $\alpha \cdot \beta$ or β is, P^A could assign a

zero to some x that $M(x)$ would assign > 0 to, in which case $\frac{P^A}{M}$ would be > 1 (in e^{-k}) would be > 1 ! So expected \sum return would be $< \infty$.

This trouble occurs whenever P^A is not universal (i.e. that's not exactly it! If P^A were universal, then it assigns ϕ to no poss. seq; but P^A could ~~not~~ assign > 0 to every poss. seq. & yet not be universal. (e.g. say P^A is uniform d.f. - or any particular seq)

So: any d.f. like P^A plus uniform d.f. would never assign ϕ to any corpus.

Uniform d.f. would get very little wt. if P^A were even marginally useful! (Uniform d.f. = "Random")

In general, getting data that would lead up to believe that $P^A(x) \approx e^{-k} M(x)$ for all x seems unlikely! Showing that it is true for just known corpus is also unlikely.

So the main idea here is that P^A is of unknown goodness (i.e. F_H)

And the ~~problem~~ is: larger $\alpha \cdot \beta$ is, the closer P^A is to P^B in the sense of more codes.

So, since P^A is optimally close to M ; P^A is perhaps close to M .

I guess analogous arguments hold for P^B .

Read notes in §78 folder: It may have good ideas on coding! (It may, indeed, but those notes are not "easy to read"!) P^A P^B

Perhaps P^A is as close as we could ever get to M (using M as a priori info).

If we could find out what the code equivalents of P^A & P^B overlap, we could simply add P^A & P^B and cut out the overlap! - the P^A & P^B are sets of non overlapping codes, they still both consist of a set of codes. Undoubtedly, this expression of P^A & P^B as sets of codes

can probably be done in ≥ 1 way. - But if we use Arith. coding, there will be unique code for P^A & for P^B . - In which case, it would seem that we could do a simple wt. mean of ≥ 2 perms Hrr., see +1.09!



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

Well, if we use a wtd sum of p^0 & p^1 for prodn, wts are α, β ; α, β , resp.
Then clearly one wants pairs of high poss. α, β .

And the goal is to get find a set of pairs $\Rightarrow \sum_i \alpha_i \beta_i$ is as large as poss:

best approx to P_M . In general, we don't ask what p^1 is "better" than p^0 :

we just use wtd. mean w. wts α, β .

This sort of "solves" the problem of justifying looking for pairs of large α, β -
They contribute most to the prodn. (Most wts.)

T. fogg. Idea is not new: I just forgot it!

There is some uncertainty about this $\Rightarrow \alpha, \beta$, here: I'm not sure that there can't be
"overlap" in the codes for p^0 & p^1 . The av. code length of -1.37 is amusing, but I
don't have complete confidence in it!

If "overlap" betw. pairs ~~is imp.~~, then finding a single pair of Max α, β
becomes the Main Goal:

500.32

They depend on "No overlap" of codes betw. p^0 & p^1 .

Another Q that I probably would not log: for the next prodn, I can use p^0 & p^1 w.
wts α, β ; α, β , resp. I can devise a new pair that predicts every
bit of the corpus, using a wtd mean of p^0 & p^1 , but the wts change with each bit, depending
on success or failure of each individual pair.

Would A & B end up identical? - it's a "simple" algebraic Q.

On second thought A & B seem to be the same: A just considers the prodn of the $n+1$ th bit
(for a n bit corpus), but B predicts for each of the n bits.

Perhaps write up section of Report that this is about NOW! So I don't forget it!

Give more discuss. in some depth! This is mainly for ME: If necc, I can

trim it from the "report".

Outline of section on the QA problem: (The problem has already been defined in
introduction. ~~Mo!~~ j-st ~~Maund~~ has a kind of OZ problem).

In the introduction, perhaps mention "Ganz OZ problem" of (700.00 = 16) or ~~section~~
in the introduction of Report but have F.N. in my version of the report. F.N. can refer
to notes, or have more content of Notes.

Do, deph Ganz soln w. $CB=0$: This does justify the approx. I'm using

on overlap of p^0, p^1 : say they overlap: so $p^0 = \alpha_0 + C$; $p^1 = \alpha_1 + C$

C = common overlap. α_0, C is α , over 3 nonoverlapping sets of codes,

The best way to predict would be $(\alpha_0 \cdot \text{prodn of } \alpha_0 + \alpha_1 \cdot \text{prodn of } \alpha_1 + C \cdot \text{prodn of } C) / (\alpha_0 + \alpha_1 + C)$.

If we just mix p^0 & p^1 without considering overlap, we use just

prodns of p^0, p^1 & C w. wts p_0, p_1 & z .

of p^1

In assuming that my ~~pair~~ has large α, β is a larger α or β . P_M predicts
"closer" to M , the generator. Larger α, β means closer to P_M : Whether it has to be closer to M , or unclear.

701.00 spec



00: 699.40: **[SN]** On QZ problems: In general, we have $\text{Re maxz } G \leq \text{Minz } C$.
 There are many ways to resolve this (often) conflict. All solns are partially ordered at least.
 1) Best ways to order: fix C ~~at~~ C : Get Max G : This is usual form of QZ problem
 2) fix G : minz C : This is a INV problem form. The idea is to "satisfy & content" in min C.
 3) Any combination of partial to linear order: Step function of G & C to be
 Max G . \rightarrow see 205.00! This becomes a time-varying $[G(x,t)]$ problem!
 4) "Anytime problems" I don't know just how they define it. soln. One way:
 Do regular QZ w. fixed C but periodically $C \leftarrow kC$ k can be 2 (or close to).
 Using large k makes it poss. to change character of soln.
 Another way is $C \leftarrow C + d$ whenever we reach C . The
 degree of Greedy ness depends on size of d or k .
 2 2004
 1.5 100k.
 When a user is faced with a form of QZ problem, we can use the 2-step Update
 of GCPD to solve it: T. second step depends on t / form of QZ problem definition.
 16 705.00

16 **[SN]** After Expo introd: readers should have clear idea of the "parts" of the system
 1) GCPD 2) [PST] 3) Just what the QZ problem is 4) How Update of GCPD, [PST] is done.
 5) How TSC's are written.
 How a QZ problem is defined: Appl, corpus.

20 **[SN]** On sequencial predn is a QZ problem. Normally a code for an "A"
 is a s.d. code (discrete universal pred): If prefix codes are allowed - so we want codes for
 any string of what-if-known A is a prefix, then we can do sequencial prediction.
 In general, the "Q" of a problem has to tell what kind of problem it is - whether it
 is discrete or continuous prob. is to be used. This makes the s.d. code & prefix codes give purchase to some
 P.D. for predictions. To s.d. codes gives cover's "Extension Complexity" for sequencial predn.

30 ask the algn On text for any Q_0 , gives a P_0 on A_0 & its extensions, finite or infinite.
 Also explain how to use 11 time series as Q_i, A_i 's (did I mean 2 or more T.S.'s in direct series w. poss. interaction? - Probly)
 Also explain how to use Encyclopedia as meta
 " " about how to use a large corpus (say SM or Encyc) as data - Root cost of
 codes include cost of referring to data in the corpus
 This expands report a lot! Perhaps leave it out or put "footnotes" in my version, but refer
 to my notes.



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20. (59.40): It may be possible to show closeness to P_M implies (in this case) closeness to M .

The fact that $\alpha \cdot \beta$ is large makes it likely to be close to M 's parameter.

Consider: $E_M \ln(\text{corpus}) - (\ln M(\text{corpus}) - \ln P_M(\text{corpus})) \geq 0$.

Consider: $E_M \ln(\text{corpus}) (\ln M(\text{corpus}) - \ln P^0(\text{corpus}))$

Troubles: Does not need to consider \ln approx of P_M or of P^0 .

Perhaps they have to because we are only

considering P^0 & P^1 as approxs of the component codes of P_M , so the approxs α_0, α , have to be included.

20.09 Given a corpus: M is $M(\text{corpus}) = \text{approx}(M)$: finding a $M \Rightarrow P_M$ is Max, is what we were trying to do. So in this sense maximizing it gives us our best approx to M .

Given $CB = \infty$, we can find Max of .09: it is P_M .

20.11 So in finding a funct. $f(\text{corpus}) = \text{approx}(f)$ is max, we can be trying to approximate M rather than P_M .

We know that the best f that we could find w. $CB = \infty$ is P_M & it is v.g. for prodn. (small error).

But, as I noted before, no matter how good our f (of (1)) is we can't be sure that it is always writing a constant factor of M .

Picking a single $f \Rightarrow$ all is max many sort of all kinds, and has picked f most likely with M .

So it's a Max likelihood choice w. careful consideration of f & prior of f .

20.19 This approx. of $f \Rightarrow$ v.g. for $CB = \infty$.

20.20 Introdn: Explains INV probs, OZ probs: Also 700.00-16 on how OZ & INV is a

How find CB is not the only way to define OZ probs.

20.23 Section in INP: Give name of QA prob. Explain that QA's can be strings of no's. Give some good examples (Also English).

Also explain how Bay induction & Time series can be covered by this formalism also. Give 4.14.18 & explain how it is an OZ prob. & we can have exp. relation.

Explain Primary Ind prob. is defined by corpus, & approx. & in predictability, by FST to be used in CB or $(G, \text{relation})$

Discuss $CB = \infty$ for 4.14.18 & some of its properties. Explain about continuation of the Q sequence:

that $t \leq \text{err}$ or KL disc is bounded, n — what this effect makes — if Q 's available under Q 's v.s. if they are quite different from any part of corpus.



20: : Expo Computing Resources ($\approx cc$):

When I speak of "Time to perform some solve a problem using a particular PST", I usually mean a generalization of time: i.e. "computation cost". In addition to time, this can include memory size, cost of CPUs(s), weight, electrical energy (watts times time), cost of cooling (watts times time), Real estate (both on chip and ~~data center~~ area of over ~~space~~ ^{in human community}).... In general, ~~think~~ ^{relevant} ~~on~~ ^{of} ~~costs~~ ^{of costs} will be different in ~~each~~ ^{each} computing environment. What ever does cost, we can ~~save~~ ^{minimize} ~~it~~ ^{it} \Rightarrow a scalar "money" - one of the costs we want to minimize.

Section on induction:

In this particular paper, we will mainly discuss a certain / limited form of induction, since the updating of ~~Pro~~ / ~~system~~ ^{problem} can be put into this form. More general forms of induction can be implemented with minor modifications of ~~Pro~~ / ~~limited~~ form.

We are presented with a sequence of Q_i , A_i pairs ... $i=1 \dots n$ (the "corpus") The Q 's and A 's can be strings and/or numbers. More generally, they can be any data structure that is uniquely encodable as a ~~finite~~ ^{finite} string of bits.

We are given P_{old} and P_{new} and the problem is to obtain a probability distribution over all possible strings, A_{new} . The ordering of ~~Pro~~

The Q 's and A 's may ~~correspond~~ ^{correspond} to questions and answers in some formal or natural language; they may be the inputs and corresponding outputs of a ~~known~~ ^{known} unknown (or partially known) mechanism - they can be real, integers, or strings. Many induction problems can be put into this form. With very little modification of our formalism, we can deal with the extrapolation of an ordered sequence of strings (language induction) or reals, (curve fitting). We can deal with sequential prediction - in which our corpus is an ordered list of strings and/or numbers.

We can ~~extrapolate~~ ^{extrapolate} ~~very~~ ^{very} data from a large digital text ~~such as an encyclopedia~~ ^{such as an encyclopedia} or dictionary.

~~The solution to our~~ A formal solution to our problem is to find (in the variable ~~time~~) \rightarrow 703.30

30: ~~the~~ a ~~probabilistic~~ ^{conditional} operator O^n such that

31:
$$O_n \triangleq P(O^n) \cdot \prod_{i=1}^n O^n(A_i | Q_i)$$

is as logical as possible.

$O^n(A|Q)$ is the probability of A in view of Q . ~~the probability distribution over all~~

~~possible strings A, can take many forms.~~

$P(O^n)$ is the ~~pro~~ ^{pro} ~~prob~~ ^{prob} probability of the ~~pro~~ ^{pro} ~~conditional~~ ^{conditional} probability ~~of~~ ^{of} O^n . ~~if~~ ^{if} ~~can be~~ ^{can be} ~~represented~~ ^{represented} by ~~the~~ ^{the} number of bits in a short ~~program~~ ^{program} that ~~implements~~ ^{implements} O^n , ~~then~~ ^{then} $P(O^n) \approx 2^{-n}$. ~~Other, more exact, forms can be used.~~ ^{Other, more exact, forms can be used.}



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

Spec 700.16 in particular

SW

HA

J. Generalized \Rightarrow OR problem in which f is ~~linear~~ trying to be Max'd

\Rightarrow a function of $G \in T$, \Rightarrow seems to be identical w/ \Rightarrow General ordering function on $G \in T$! We have a partial ordering of G, T pairs v.z $G' > G \Rightarrow T < T'$.
for a linear ordering we need a function of $G \in T$ — which is \Rightarrow same as $G = G(x, T)$.

How to solve \Rightarrow is still unclear, but normal fixed CB OR probs \Rightarrow INV probs

give "special cases".

\Rightarrow One approach is via GCPD, \Rightarrow GCPD₂; In GCPD, we have to DP's for

each ~~state~~ ~~function~~ \Rightarrow (PST, G, CC) from known ~~function~~

form prob orders G, CC , we can select \Rightarrow PST ~~prob~~ for given problem

\Rightarrow is most likely to have Max G(x, T)

We can obtain this Max Carto wise since we rarely don't need an exact soln.

it may.

So This seems to solve the classic time varying OR problem

Max G(x) * F(t) in f . More Gen / form G(x, T).

20: 704.40 we will then describe improvements in the system that will overcome these limitations \rightarrow

21 These improvements may be introduced into the system in ~~an~~ either or both of two ways: They can be programmed into the system by the operator, ~~and/or~~ ~~the~~ operator acquired by the system v.z its training sequences, while this is much better for the system to acquire knowledge v.z learning, some of the needed improvements take too long and they

25 study reports essential computation time for the system to learn ~~the~~ the operator the normal training sequences in which case learning may be programmed in by the operator.

27 Maybe say its w/ seqs 21-23 — much shorter. \rightarrow while almost all of these improvements can be acquired by the system thru its normal training sequences, some of the needed improvements in \rightarrow (25-27)

Do describe different kinds of OR problems: Until recently improvements the operator variables

OR time varying optimization was not to be conceivable by L such. We have, however, been covered out our updating technique can be readily adapted to this class problem



So try this: After Introduction:

- 1) Soln to Induction Problem. } "Prelimmys"
- 2) How Lsrch works
- 3) The evolution of the system, 704.35 → 705.30 gives initial expo.

996.13 gives a hint of outline of the paper.

$$\frac{1}{1 + \frac{1}{k}} = \frac{k}{k+1}$$

$$\frac{2k}{k^2} - \frac{1}{1-k} = 0$$

$$\frac{2}{k} = \frac{1}{1-k}$$

$$k = 2 - 2k$$

$$3k = 2$$

$$k = \frac{2}{3}$$

$$\frac{1}{1 - \frac{1}{k}} = \frac{1}{\frac{k-1}{k}} = \frac{k}{k-1}$$

$$= -\frac{2}{3}!$$

A "not bad" approach. First do "Prelimmys" IND & Lsrch.

Then we give "bare bones" (prelim) system: What INW posts is modify & P. in accord w. "What we would" in post (ind of cc's of "Murray Cards").

Can this be done for INV/OZ probs still? I guess it can: we would simply modify f. p. d. so it would give by PC's ~~to~~ PST's ~~the~~ Post here in k, post, been successful (ind of cc's involved).

For simple Lsrch, for both INV & OZ problems ~~there~~ is only one d.f. on PST's for all problems (possibly 1 for INV, 1 for OZ), so it's easy to pool data.

The GCPD is not being used to select PST's (There ~~is~~ only 1 PST being used - i.e. Lsrch.)

5X I've forgotten details of the "Conv. thm for Bess". I must write this up! i.e. I don't remember details of Thm. (Get alone details of proof!) use ALP: Sort of 2 part codes: sum from n ||:

$$\alpha_n = \sum_{j=1}^n \alpha_0^j \prod_{i=1}^n O^j(A_i)$$

A_i is the i th object in bag. probably use $(O^j(A_i))^{m_i}$ if A_i occurs m_i times.

O^j is jth part
 O^0 is sprig of jth part.

If $\mu(A)$ is the "true" part generating $\{A\}$

$$\alpha_n > \alpha_0^M \prod_{i=1}^n \mu(A_i) \quad | \quad \alpha_0^M \Rightarrow \text{sprig of } \mu(\cdot).$$

α_n is indep of order of $\{A\}$

(μ) α_n assigns to corpus, α_0^M times not assigned by μ to corpus. - INDICATOR CORPUS.

μ assigns pc's to each of A_i : Does μ ? It can assign a "partial" pc to A_n

- "Given to other $n-1$ A_i 's"

→ 718-22

98TM 106-108; 110, 111, 115

bazmehon 113

(11.01-17) is good.

108.01-02

110.01-10

111.01, -17

Gen "good": 110.08-34; 111.03-17

More recent stuff

ID 378.00 ff!

380.23 -40

381.00 -08

8/19/02

or (Relative entropy)

Kullback-Liebler

DISTANCE

KL Perm. 30 to 708.10

708.10 is a kind of

707 BUG



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:706.40: I was concerned that the "easy" approach to proof, using the TS inequality, assumed
the probs for P' (based on P_A) were sequentially indep. It seems that

they are not all indep
funt. "Bay" / "Poisson".

Rec P_C predn. of the k -th element in the sequence is:

0.06

$$\prod_j \alpha_0^j \prod_{z=1}^k \alpha_0^j \alpha_0^j(A_i) \quad \left/ \quad \prod_j \alpha_0^j \prod_{z=1}^k \alpha_0^j(A_i) \right.$$

normal factor.

If this is so, then there is clear cancellation of all (or all but one?) normal factors.
Maybe not! T. normal constant here is wrong. ~~Not exactly!~~ Not exactly! It

is ~~correct~~ correct for $P_A = \text{some measure}$. The $\frac{P'}{P} > \alpha$ inequality
is true for some measure, but a different one, which is always \geq some measure.

If guess, in a similar way we have an eq. like 06 for individual bits of the $\{A_i\}$ sequence

This way seq. of bits would enable us to apply TS corr. directly, in one step,

to get desired result. So $\frac{P'}{P} > \alpha$ is order to each bit/bits of the corpus —

so TS corr. holds (is applicable).



Rec. to Application of (TS corr.) Does this mean that I must be able to associate

2 PC w. any finite string of bits? (Some of these strings may have $P_C \geq \alpha$ or α .)

Actually I may need only 2 PC, (P') assoc. w. a bit string w. unknown before the A_i 's

and 2 other PC's for that same kind of string.

But it is concerned w. "Expected values" is ~~some~~ we sum over all poss.

sequences of A_i 's. Phrase will be general, have ~~different~~ bit lengths.

A further simplification to prove the K-L part of STS73:

$$KL(\vec{P}_1, \vec{P}_2) = \sum P_i \ln \left(\frac{P_i}{P_2} \right)$$

Theorem: $KL(\vec{P}_1, \vec{P}_2) + KL(\vec{P}_3, \vec{P}_4) = KL(\vec{P}_1 \times \vec{P}_3, \vec{P}_2 \times \vec{P}_4)$

$\vec{P}_1 \times \vec{P}_3$ is the cartesian product of \vec{P}_1 & \vec{P}_3 : if \vec{P}_1 has n elements & \vec{P}_3 has m elements, then $\vec{P}_1 \times \vec{P}_3$ has $n \times m$ elements. We can think of each being product of components PC's.

linear ≥ 5 has $n \times m$ terms. Any way this must be (most likely) $n \times m$ terms combined? — Yes! → 708.00

$$\sum_{i=1}^n \sum_{j=1}^m P_i^j P_3^j (\ln P_i^j + \ln P_3^j - \ln P_2^j - \ln P_4^j) \quad \text{4 mn terms}$$

$$\sum_{i=1}^n P_i (\ln P_i - \ln P_2) + \sum_{j=1}^m P_3^j (\ln P_3^j - \ln P_4^j) \quad \text{2n + 2m terms}$$



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

! work, $\sum_{k=1}^r p_k^k = 1$ for all $k \in r$. So, looking at the various terms of -1.35!

$$\sum_i \sum_j p_i^i p_j^j \ln p_i^i = \sum_i p_i^i \ln p_i^i$$

more generally $\sum_i \sum_j p_i^i p_j^j (\ln p_i^i - \ln p_j^j) = \sum_i p_i^i (\ln p_i^i - \ln p_j^j)$

$$\sum_i \sum_j p_i^i p_j^j (\ln p_j^j - \ln p_i^i) = \sum_j p_j^j (\ln p_j^j - \ln p_i^i)$$

So (-1.36) = (-1.35) + term is proved! $\sum p_i = 1$ \Rightarrow recycles, so we

would have to have the "0" symbol as a proxy.

So the KL's add or incrementally \Rightarrow it is tedious to prove S78 T3 corr. (Per K-L proof)

I looked at Cournot's Day: Several problems on KL dist, but none like this one
Hus says give 2 refs $\left[\begin{matrix} \text{Csiszar} \\ [78] 1967 \end{matrix} \right]$ and also Amari Book Springer 1987 [10] ref

The application of the Perm depends on at least 2 things

- 1) $\sum p_i = 1$; so all probs are normalized (possibly by fixing p_0 to "0").
- 2) \Rightarrow 2 d.p.'s to be combined must be indep so prob of combination = product of probs.

If the p's are not indep (-1.35) becomes $\sum_{i=1}^n \sum_{j=1}^m p_{i,j}^{i,j} (\ln p_{i,j}^{i,j} - \ln p_{2,4}^{i,j})$

In the case of S78 T3; the p's are not indep. so it would seem that the Perm doesn't apply!

On the other hand, if this proof is correct, it would seem very likely that the theorem would apply!

The theorem would apply to successive A_i 's, but not to bits within an A_i , nor to S78 T3

On the other hand maybe $\sum_{i=1}^n p_{i,j}^{i,j} = p_j^j$

In (-1.36): KL (p_3, p_4) is not really meaningful if p_1 & p_2 are not indep & $(p_3$ & $p_4)$ are not indep!

There may be a way to define this so that Perm is true!

At present it's not clear in my mind as to the meaning of KL distance if the second pair of p's is dependent on the first. But the $\left(\frac{KL \text{ dist}}{KL} \right) \Rightarrow$ defined for the entire 2 step system. $\rightarrow 35$

Def

Another approach way This method of getting KL Dist for individual A_i 's (which are indep)

Or more crudely, regarding each A_i as an indep binary sep. to which S78 T3 cor. applies, Perhaps the Perm applies to successive A_i 's, since they are "indep". It is indep, but the wts of the PM components are not.

For the first part of .29: either the first step of system is in a variable p.d. over a set of states. KL D: measure dist betw. 2 p.d.'s over same set of objects.

Peter Gacs was this out in his proof of the conv. Perm!

Levit first edition 1993 pp 235-238

Levit second edition 1997; pp 328-331

The theorem(s) he proves about KLD are a simpler & more general way of ~~doing~~ doing

These are the same



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 209:40 : $P_m \in \mathcal{M}$ for identical seqs of A_i 's

01 For set of all poss \mathcal{M} seqs of n A_i 's: $T: E \text{ of } \left(\sum_{i=1}^n \text{KLD between } P_m \text{ and } P_{A_i} \right) < \frac{K}{n}$.

"E" assigns P_m to each bit m in A_i 's. To get KLD, P_m assigns P_m 's to each bit depending on previous bits in A_i .

04 All previous A_i 's, as well as previous bits in A_i containing bit m (Question).

Under cond. 00-04, Gacs's theorems about KLD [in Livit 1997 pp 329-331] apply

And we get: $\sum_{i=1}^n D(P_m || P_{A_i}) \leq K$

How to get meaningful "E"; in-01 = We have to take E over the bits of the ensemble.

Some members of the ensemble may not have i bits: but that's o.k., the "E" then doesn't include them.

So we know, for n "QA's" an ensemble of strings — not all the same length. Both A_i & P_m assign P_m 's to each bit of each string, & we also have a "run" product for both A_i & P_m .

Perhaps try to explain it as in "a paper":

Int. sequence of bits; we do need "end of A_i " markers: The positions of the markers profoundly influence P_m 's of bits that follow them.

One poss. way to deal w. this way to use SYNTAX on each of the A_i 's as a separate corpus. P_m is diff. for each A_i , here. We may need a "stop" symbol. [Hopefully, Hutter has published proof that it works for 2 symbols]

$$\text{Then } \sum (P_m(s) - P_m(s))^{2^k} + (P_m(0) - P_m(0))^{2^k} + (P_m(1) - P_m(1))^{2^k} \leq K$$

More generally $\sum D(P_m(s) || P_m(s)) \leq K$

22 $\sum D(P_m(s) || P_m(s)) \leq K$; $D(P_m(s) || P_m(s)) = \sum_{y \in \{0,1\}} (P_m(y) \ln \frac{P_m(y)}{P_m(y)})$

So, the ensemble is the set of all binary seqs w. n "s" in them

27 or, the ordered set of all finite binary sequences (w. possible repetitions).

28 For each A_i : $\sum_{k=1}^{\infty} D(P_m, P_{A_i}) \leq E \ln \frac{P_m(A_i)}{P_m(A_i)} = E \ln \frac{P_m}{P_m} = E \leq \frac{K}{n} < K$

But A_i 's will vary in length for each i .

If we have to do 28 for all poss. finite strings, A_i 0, 1, 00, 01, 10, 11, ... - and so on.

34 Or, use a "Blank" symbol, w. understanding that once a blank occurs then all following must be blank.

What this does is make $M = P_m$ for all blanks & for the first m (which includes A_i of length $\leq m$)

28 holds w. M 22 summed over $y = 0, 1, \text{Blank}$.

38 For 28 consider $A_1, A_2 \Rightarrow E \ln \frac{P_m}{P_m} |_{A_1} + E \ln \frac{P_m}{P_m} |_{A_2} = E \ln \frac{P_m}{P_m} |_{A_1+A_2} > \frac{K}{n}$

I used this to be true, if 28 is to be useful

Yes 711.00-10 Previous!



01: 710.90

02

03

05

0

12

21

22

24

0

0

0

$$\sum M \ln \frac{M}{PM} \left(A_1 + \sum M \ln \frac{M}{PM} \right) A_2$$

$$= \sum M \ln \frac{M}{PM} \Big|_{A_1, A_2}$$

$$\sum M A_1 \ln \frac{M A_1}{PM A_1} + \sum M A_2 \ln \frac{M A_2}{PM A_2} \propto$$

$$\left(\sum M(A_1 + A_2) \right) \ln \frac{M(A_1 + A_2)}{PM(A_1 + A_2)} \propto$$

$$\sum_{A_1, A_2} M(A_1) \cdot M(A_2) \left(\ln \frac{M A_1}{PM A_1} + \ln \frac{M A_2}{PM A_2} \right) \propto$$

$$\sum_{A_1, A_2} M(A_1) M(A_2) \ln \frac{M A_1}{PM A_1} = \sum_{A_1} M(A_1) \ln \frac{M A_1}{PM A_1}$$

if $\sum_{A_2} M(A_2) = 1$ so it looks like $\alpha(1.00) = \beta(1.03) = \gamma(1.02)$

Which means $(.00 L) = (.02 L)$ as desired. So -1.38 is true.

Which shows -1.28 can be somewhat ~~smaller~~ $(\leq "A_2")$

So in -1.28 : Summing over all A_i :

$$\sum_{i=1}^N M(A_i) D(M(A_i) | PM(A_i)) \leq \sum M \ln \frac{M}{PM} \Big|_{A_i}$$

So total bit divergence

was -1.28 is correct? do I need E here? Yes!

K_0 - K_0 is some large no.

$$\sum_{i=1}^N M(A_i) D(M(A_i) | PM(A_i)) \leq \sum M \ln \frac{M}{PM} \Big|_{A_i} \leq K$$

I think I can do it w. only 2 symbols: Just consider all possible pairs of A_i .

T. expected $E \leq \text{error}$ is still $\leq E \ln \frac{M}{PM} \Big|_{A_i}$

Not entirely clear how to do this, hvr. \rightarrow Note 712.00, hvr (objection!)

Another Q: Re: ~~of~~ proof of $.00 - .03$ Did I assume $PM(A_1)$ & $PM(A_2)$ were indep?

If maybe that they are not - in which case I'd have to look at new facts about w. M situation.

Look at $.05 \sum_{A_1, A_2} M(A_1, A_2) \ln \frac{M A_1}{PM A_1} : : .05$ may be ok.

$$M(A_1, A_2) = M(A_1) \cdot M(A_2) : \text{But, } M(A_2) \text{ may be a function of } A_1! \quad \text{PM}(A_1), \text{ PM}(A_2) \text{ may depend on } A_1, \text{ but } M(A_2) \text{ is indep of } A_1$$

lets look at $\ln \frac{M(A_1, A_2)}{PM(A_1, A_2)} : \begin{matrix} \uparrow \text{numerator factors out.} \\ \uparrow \text{Denom doesn't.} \end{matrix}$ $PM(A_1), PM(A_2)$ are not indep

I'm not sure this gives trouble, hvr. But I do have to look into it!

$$\sum_{A_1, A_2} M(A_1) M(A_2) \left(\ln \frac{M A_1}{PM A_1} + \ln \frac{M A_2}{PM(A_2)} \right) \geq \sum_{A_1, A_2} \ln \frac{M(A_2)}{PM(A_2)} M(A_1) \cdot M(A_2)$$

\rightarrow 712.02 also



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

01 (711.24) spec (SN) I think i. "B's" of 710.34 (or equiv) may be necy. We do need info about when
to A_2 end, start.

02 (711.40) spec Look at $\sum_{A_1, A_2} \mu(A_1, A_2) \ln \frac{\mu(A_1, A_2)}{P_{M_1}(A_1) P_{M_2}(A_2)}$ more carefully

$$JF = \sum_{A_1, A_2} \mu(A_1) \mu(A_2) \left(\ln \mu(A_1) + \ln \mu(A_2) - \ln P_{M_1}(A_1) - \ln P_{M_2}(A_2 | A_1) \right)$$

$$= \sum_{A_1, A_2} \mu(A_2) \cdot \mu(A_1) \left(\ln \mu(A_1) - \ln P_{M_1}(A_1) \right) + \sum_{A_1, A_2} \mu(A_1) \mu(A_2) \left(\ln \mu(A_2) - \ln P_{M_2}(A_2 | A_1) \right)$$

$$\rightarrow = \sum_{A_1} \mu(A_1) \ln \left(\frac{\mu(A_1)}{P_{M_1}(A_1)} \right) + \sum_{A_1, A_2} \mu(A_1) \mu(A_2) \ln \frac{\mu(A_2)}{P_{M_2}(A_2 | A_1)}$$

$\sum_{A_1} \mu(A_1) (\ln \mu(A_1) - \ln P_{M_1}(A_1))$ is ok. for
to first term, but second term

should be $\sum_{A_1, A_2} \frac{\mu(A_1, A_2)}{\mu(A_1)} \cdot \ln \left(\frac{\mu(A_1, A_2)}{\mu(A_1) P_{M_2}(A_2 | A_1)} \right)$ so T. terms do add up to give desired equality & inequality

So I may have proved it! The A_2 term in both $P_{M_2} \leq D(\cdot || \cdot)$ and $\sum \mu \ln \frac{\mu}{P_{M_2}}$ parts
are a bit complicated, but they "fit together" properly!

If no securing necy, The B (blue) symbol is needed, then we can make it correct by
a simpler ternary sec radix, which is dealt w. in S 28+3. (we have additional
info that $\mu(A_i) = \mu(A_j)$ are indep if $i \neq j$) ($P_{M_2}(A_2) = P_{M_1}(A_1)$ are usually not always)
not independent.

It is my improv'd Gacs' proof that $\sum_{i=1}^n \mu D(\mu || P_{M_i}) < K$
works for any radix. That $D(\mu || P_{M_i}) < 2(M - P_{M_i})^2$: he only proves for radix 2 - which is easy!

To show it for radix ≥ 2 is not so easy; (ES Almost Certainly true): Hutter says he has proved it.

One poss. way! Start w. radix 2, then 3: for radix 2: $D(x || y) = 0$ for $x=y$; $(x-y)^2 = 0$ for $x=y$.
Look at line $x=y$: $\frac{\partial^2}{\partial x^2} > \frac{\partial^2}{\partial y^2} > 1$ for x, y on unit sphere.

Can't be shown for radix 3 & higher:

for radix 2: $f = \mu(x) (\ln \mu(x) - \ln P_{M_2}(x)) + (1-\mu(x)) (\ln(1-\mu(x)) - \ln P_{M_2}(1-\mu(x)))$

$$\frac{\partial f}{\partial x} = \mu(x) \ln \mu(x) + 1 - \mu(x) \ln P_{M_2}(x) - \frac{\mu(x)}{2x} + \dots$$

for four bits ≥ 2 Thy Pits! for radix $r \geq 2$: Use Lag. Mults. star

We constraint $\sum_{i=1}^n x_i = 1$. sep. w/out $f = \sum_{i=1}^n (x_i (\ln x_i - \ln y_i) - (x_i - y_i)^2) = \min$ (p. 104)
is. f is always > 0 .

I'd guess that $\sum y_i \leq 1$ is unnecessary. 2 r equs; 2 r + 1 unknowns! some endup.

2 1 parameter is sets of minim f - time! Its $x_i = y_i$ is link



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

so: 7/24/02: No! $x_i = y_i$ $i=1, 2$ is not a line: It's \mathbb{R}^2 $2+1$ dim space and has r constraints, so it's a $r+1$ dim manifold!
 $x_i = y_i$ ($i=1, 2$) is only part of the solution. Remember r other constraints. - so $\mathbb{R}^2 \supseteq \geq 1$ dim. object

03 -- $F = -1.35!$... $\frac{\partial F}{\partial x_i} = \ln x_i + 1 - \ln y_i - 2(x_i - y_i) + \lambda = 0$
... $\frac{\partial F}{\partial y_i} = -\frac{x_i}{y_i} + 2(x_i - y_i) = \phi$...

$\sum y_i = 1$ is not
vs $\sum x_i = 1$
only $\sum x_i = 1$

If $x_i = y_i$ then $\lambda = 1$ (from 03)

" $x_i = y_i$ then no solution for 05!

If we add constraint $\sum y_i = 1$, we get a line

$\frac{\partial F}{\partial y_i} + \lambda_2 = -\frac{x_i}{y_i} + 2(x_i - y_i) + \lambda_2 = 0$. so $\lambda_2 =$

From $x_i = y_i$ is a soln. and $\lambda = \lambda_2 = 1$. But remember only $r+2$ constraints.

using $(0.5, 0.5)$ as is; $\ln \frac{x_i}{y_i} + 1 + \lambda \frac{x_i}{y_i} = 0$ solving for $\frac{x_i}{y_i}$ gives $\frac{x_i}{y_i} = 1$ if $\lambda = 1$

but if $\frac{x_i}{y_i} = z$ $\ln z - z = -1 - \lambda$ so z can have any value! λ is associated λ . for $z < 0$, λ is complex. - but $\frac{x_i}{y_i}$ has to be a same for all i .

hvr, say $\frac{x_i}{y_i} = \alpha$ $x_i = \alpha y_i$ $y_i \cdot 0.5 - \alpha + 2(\alpha y_i - y_i) = \phi$ $-\alpha + 2y_i(\alpha - 1) = 0$

$-\alpha + 2(1-\alpha)y_i = 0$ $y_i = \frac{\alpha}{2(1-\alpha)}$ $x_i = \alpha y_i = \frac{\alpha^2}{2(1-\alpha)}$

so this is a 1 param manifold: if we choose $\alpha = z$, then $y_i = \frac{\alpha}{2(1-\alpha)}$ and $x_i = \alpha y_i$

$\beta = \frac{\alpha}{2(1-\alpha)}$ $\frac{1}{2\beta} = \frac{\alpha(1-\alpha)}{\alpha} = \frac{1}{\alpha} - 1$ $\frac{1}{\alpha} = \frac{1}{2\beta} + 1$
 $\alpha = \frac{1}{\frac{1}{2\beta} + 1} = \frac{2\beta}{1+2\beta}$

y_i ; $x_i = \frac{2y_i^2}{1+2y_i}$ hvr, if $x_i = \alpha y_i$ $\alpha \neq 1$ then y_i can't be normalized, consequently, y_i never is normalized.

But in my 578 T3, y_i was normalized! and the relation between (x, y) was OK.

Also $x_i = y_i$ is a min point, yet $\frac{\partial F}{\partial y_i} = -1$ at $x_i = y_i$

try $r=2$: $x_1(\ln x_1 - \ln y_1) + x_2(\ln x_2 - \ln y_2) - (x_1 - y_1)^2 + (x_2 - y_2)^2 + \lambda = 0$

$\frac{\partial}{\partial x_1} = 1 + \ln x_1 - \ln y_1 - 2(x_1 - y_1) + \lambda = 0$

$\frac{\partial}{\partial y_1} = -\frac{x_1}{y_1} + 2(x_1 - y_1)$



: Gaci's discussion of top of p 329, Livit 1997

~~$$f(Q_0) = P(0)(\ln P(0) - \ln Q_0) + P(1)(\ln P(1) - \ln Q_1) - 2(P(0) - Q_0)$$

$$f(P(0)) = \dots$$~~

$$\frac{df}{dQ(0)} = \frac{P(0)}{Q(0)} = \frac{P(1)}{Q(1)} + 4(P(0) - Q(0))$$

~~...~~ = 0 if $P(1) = Q(1)$

look at old 78 Notes, early versions.

$$f(Q_0) = \sum_{i=1}^n \left(P_i \ln \left(\frac{P_i}{Q_i} \right) - (P_i - Q_i)^2 \right) \quad \frac{df(Q_0)}{dQ_0} = \frac{P_0}{Q_0} - 2(P_0 - Q_0)$$

unnormalized
In the original STS paper, I had the proof of $KLD \leq \frac{1}{2} \sum (P_i - Q_i)^2$.

Now, I only use $\frac{\partial f}{\partial \mu}$ & $\frac{\partial^2 f}{\partial \mu^2}$! Hvr. b.c. μ & P_i were ~~...~~
Hard Measures.

So, from (-1.05) $\frac{\partial f}{\partial y_2} = -1$ for unnormalized \vec{y} at $\vec{x} = \vec{y}$.

If we normalize at $x(-1.09)$, we can get $\frac{\partial(f + \lambda \sum y_i)}{\partial y_2} = 0$ at $\vec{x} = \vec{y}$.

$$\frac{\partial^2(f + \lambda \sum y_i)}{\partial y_2^2} = \frac{-x_2}{y_2^2} = -2 = \frac{1}{y_2^2} - 2 \quad \text{at } \vec{x} = \vec{y}$$

$$\text{at } y_2 = 1, \text{ it is } -1! \quad \text{at } y_2 > \sqrt{2} = 1.414 < 0!!$$

$$f, \text{ in 712.35 was } f \equiv \sum_{i=1}^n x_i (\ln x_i - \ln y_i) - \frac{1}{2} (x_i - y_i)^2 \left(\lambda \sum x_i + \lambda_2 \sum y_i \right)$$

It works out well for normalized case.

write a Basic form to examine path examples! My impression was that 712.35 as it was, worked ok for $r=3, 4$!

$$x = .5 \quad y = .5 \quad \left[KLD - SQ2. Bas \right] + \lambda \sum x_i + \lambda_2 \sum y_i$$

$$\text{Print } x * \log(x/y) - .5(x-y)^2$$

x, y

I got sign of KLD wrong in 712.35! should be $x_i (\ln x_i - \ln y_i) \dots$

No! sign was o.k.

$$-1.05 \rightarrow \frac{\partial f}{\partial x_1} = -\ln x_1 - 1 + \ln y_1 - 2(x_1 - y_1) + \lambda = 0$$

for $\vec{x} = \vec{y}$, $\lambda = +1$ (was for log mult.)

$$-1.03 \rightarrow \frac{\partial f}{\partial y_1} = \frac{x_1}{y_1} + 2(x_1 - y_1) \quad (+\lambda_2)?$$

... $x_2 = -1$ was ~~...~~ for Log. Mult.

$$\frac{\partial^2 f}{\partial x_1^2} = -\frac{1}{x_1} - 2 \quad ; \text{ always } < 0$$

$$\frac{\partial^2 f}{\partial y_1^2} = -\frac{x_1}{y_1^2} - 2 \quad ; \text{ always } < 0$$

So it works o.k. if \vec{y} is normalized, for \vec{x} by table.

$$KLD - SQ2. Bas \quad (\ln(1-y) + \ln(y_2))$$

$$\geq \text{same, but for } x_1, y_1, y_2 \quad (x-y)^2 + (x-y_2)^2$$

So $y_1 + y_2 < 1$.



50.714.40:

So using normalized \vec{y} (i.e. $\lambda=2$):

~~2.713.05~~
 $2 \frac{713.05}{2x_i} = \frac{1}{x_i} - 2 \implies > 0 \text{ if } x_i < \frac{1}{2} \quad x_i = y_i = .8$

$$2 \frac{713.05}{2y_i} = \frac{x_i}{y_i^2} - 2 \implies > 0 \text{ if } \frac{x_i}{y_i^2} < \frac{1}{2} \quad x_i^2 < \frac{1}{2} y_i^2$$

~~Empirically~~ Empirically for $r=2$; both \vec{x} & \vec{y} normalized

$\frac{\partial R + L}{\partial y_i} > 0$ at $x_i = .8, y_i = .8$; at $x_i = .5, y_i = .5$ $\frac{\partial f}{\partial y}$ seem to be slightly < 0

($f(.5, .5 \pm .05) \rightarrow -2.0478 \text{ E-}08$
I tried various values very close to $y = .503$ $f(x=.5, y=.503) < 0$ but as y varies around that, the behavior is erratic - SPS is only 3 decimal places, Basic should have no problem! - But differences between small quantities! - try double precision!)

Seems to be o.k. in double precision: (I guess it does do $\ln()$ in double, too.)

try $f = x \ln \frac{x}{y} + (1-x) \ln \frac{1-x}{1-y} + 2(x-y)^2$ $x \ln x \rightarrow -x$

$$\frac{df}{dx} = \ln x + 1 + \frac{1}{1-x} + \frac{x}{1-x} = \ln(1-x) + \ln(1-y) + 2(x-y)$$

$$\ln x + 1 + \frac{1}{1-x} + \frac{x}{1-x} = \ln(1-x) + \ln(1-y) + 2(x-y)$$

$$\ln x + \ln\left(\frac{1-x}{1-x}\right) + 2(x-y) = \ln x + \ln(1-y) + \ln(1-x) + 2(x-y)$$

$$x \ln \frac{x}{y} = x \ln x - x \ln y \rightarrow 1 + \ln x - \ln y$$

$$(1-x) \left(\ln(1-x) - \ln(1-y) \right) = \ln(1-x) - \ln(1-y) - x \ln(1-x) + x \ln(1-y) \rightarrow \frac{1}{1-x} + \frac{x}{1-x} - \ln(1-x) + \ln(1-y)$$

$$2(x-y)^2 \rightarrow 2(x-y)$$
$$1 + \ln x \left(\frac{1}{1-x} - \frac{1}{1-x} \right) + \ln\left(\frac{1-x}{1-x}\right) + 2(x-y) = \ln\left(\frac{x(1-x)}{1-x}\right) + 2(x-y)$$
$$= \ln x + \ln(1-x) + \ln y + 2(x-y) \quad \left[\text{if } x=y; \text{ this is } \ln x \right]$$

Maybe not $\ln \frac{x}{y} - \ln\left(\frac{1-x}{1-y}\right) + 2x - 2y$ which is correct, I guess (I guess)

$$\frac{d.30}{dx} = \frac{1}{x} + \frac{1}{1-x} + 2 \quad \text{which is } > 0 \text{ for } 0 < x < 1$$

$$\frac{d(2)}{dy} = -\frac{x}{y^2} + \frac{1-x}{(1-y)^2} - 2x + 2y = 0 \text{ if } x=y$$

$$\frac{d}{dy} = \frac{x}{y^2} + \frac{1-x}{(1-y)^2} + 2 > 0 \text{ for } 0 < y < 1, \text{ } 0 < x < 1$$

$$x_1 \ln(x_1/y_1) + x_2 \ln(x_2/y_2) + (1-x_1-x_2) \ln((1-x_1-x_2)/(1-y_1-y_2))$$

ID



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00

using Maple I did $x_1, x_2, (1-x_1-x_2) : y_1, y_2, (1-y_1-y_2)$
 $f(x) = D(x|y)$; $\text{solve} \frac{df}{dx_i} = 0$ if $\vec{x} = \vec{y}$
presumably same for x_2

try y_1 $\frac{\partial f}{\partial y_1} = 0$ for $\vec{x} = \vec{y}$.
second dir wrt y_1 is $\frac{x_1}{y_1^2} - \frac{(1-x_1-x_2)}{(1-y_1-y_2)^2} > 0$ on unit sq,

So, for any (legal) \vec{x}, \vec{y} pts, show that \exists a \vec{x}_0, \vec{y}_0 pt. with $\vec{x}_0 = \vec{y}_0 \Rightarrow$
 $\vec{x} = \vec{y}$

On second part: if all second derivatives are positive or zero,

The rem is in Maple file: **KWDIST.MS** for (100) can easily get any first or
second derivative

$x^2 + y^2 + z^2 = 3$ (not color)
diff(" ", y);

put cursor at, then "Enter", then
cursor past "after"; then enter —
will get derivative of the top ("") expression.
doing "again" will dit wrt of that
result.

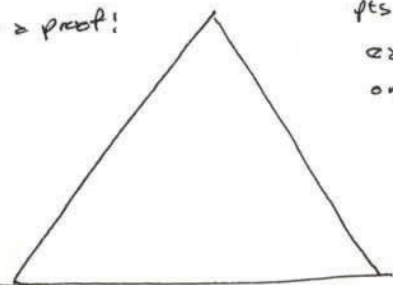
Anyway, for $n=3$,

all derivatives are 0 evolved on $\vec{x} = \vec{y}$ surface

All second derivatives are > 0 inside
to unit cube for all variables

from this info, it may be possible to prove that at all pts, the one must be "uphill" form
in $\vec{x} = \vec{y}$ point.

Maybe a proof!



pts. in $2n = 1$ (at Δ) represent vectors $\vec{x} : \sum x_i = 1$
each vector component = ϵ distant to
one of the sides of Δ . $\vec{y} : \sum y_i = 1$.

02

New tack: we want to show $D(x|y) \geq 0$ for all \vec{x}, \vec{y} ;
Take a certain \vec{x}_0, \vec{y}_0 : we know $D(\vec{x}_0|\vec{y}_0) = 0$. Consider a path $\vec{x}_0, \vec{x}_0 + t(\vec{y}_0 - \vec{x}_0)$ to
for $t=0, D=0$, each component of y is changing monotonically along the path; since
 $\frac{\partial^2 D}{\partial x_i^2} \geq 0$ along the path, D must \uparrow (unless $\vec{y}_0 = \vec{x}_0$)

03

D is a nice, convex function, w. no singularities in the region — except, for edges, where
components of \vec{x} or $\vec{y} = 0$. So inside the region, we have fine, smooth functions,

$\frac{\partial D}{\partial x_i} = \sum \frac{\partial D}{\partial x_j} \Delta x_j + \frac{\partial D}{\partial y_j} \Delta y_j$

Q if we have $\frac{\partial^2 D}{\partial x^2} = x^2 + y^2$ $\frac{\partial^2 D}{\partial x^2} \geq \frac{\partial^2 D}{\partial y^2} > 1$ but $\frac{\partial^2 D}{\partial x \partial y} = -1$ what does $g(x)$ look like?

8/26/02
JD

BUG?



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

$g(x, y) = f(x, y) = 1 - x^2 - y^2$



A small table or matrix with numbers 1, 0, 0, 1, 4.

$\frac{\partial^2 g}{\partial x_1 \partial x_2} = \frac{x_1}{y_1^2} + \frac{x_2}{y_2} + 4$
 $g = x^2 + y^2 - 10xy$ would be more complex

$\frac{\partial^2 g}{\partial x_1^2} = \frac{1}{x_1} + \frac{1}{x_3} + 4$
 $\frac{\partial^2 g}{\partial x_2^2} = \frac{1}{x_2} + \frac{1}{x_3} + 4$

So, P_{10} is second deriv of g_{10} wrt x, y and $10A > 0$.
T. function can still \downarrow much.
So look at $\frac{\partial^2 g}{\partial x_1 \partial x_2} = \frac{\partial^2}{\partial x_1 \partial x_2}$ $\frac{\partial^2 g}{\partial x_1 \partial x_2} = \frac{1}{x_1} + 2 > 0$

50

07

$\frac{\partial^2 g}{\partial x_1 \partial x_2} = -\frac{1}{y^3} - 2$ is always < 0 , but

$\frac{\partial^2 g}{\partial y_1 \partial y_2} = \frac{x_2}{y_2^3} + 2$: also always > 0 .

This is not relevant to the proof.

Whoops! it $\frac{\partial^2 g}{\partial x_1 \partial x_2} > 0$ whenever $x > 0$ and $y < 0$.
So it's cross deriv is small, it's OK.

$x^2 + y^2 + 10xy$: $f(+1, -1) = 1 + 1 - 10 = -8$.

$ax^2 + by^2 + cxy$ what is relation of a, b, c
so f function is > 0 ? well if
say a, b known. when $\rightarrow ax^2 + by^2 > cxy$
or $\frac{ax^2}{cxy} + \frac{by^2}{cxy} > 1$ $\frac{ax}{c} + \frac{by}{c} > 1$
min $\frac{ax}{c} + \frac{by}{c} = 0$ $\frac{a}{c} = -\frac{b}{c}$
 $a = -\frac{b}{2}$ $2a = -\frac{b}{2} = -\sqrt{\frac{b}{a}}$
So c must be $(\sqrt{\frac{a}{b}})^2 \cdot 2 + (\sqrt{\frac{b}{a}})^2 \cdot b = \text{min}$
 $= a + b$ so $|c| \text{ must be } < a + b$

from P_{10} it is clear that the cross derivatives have to be considered - that they have to be smaller by a certain amount than the second partial derivatives. We necessarily concern ourselves with $\frac{\partial^2}{\partial x^2}$, $\frac{\partial^2}{\partial y^2}$ and $\frac{\partial^2}{\partial x \partial y}$.
t. x & y 2nd deriv from g or h from f ! either would work ok.
The g is because of the argument of $-1, 30, 33$!
We can make it $\frac{\partial^2}{\partial x^2}$ or $\frac{\partial^2}{\partial y^2}$ or $\frac{\partial^2}{\partial x \partial y}$
from (y, say) . In $(0, 0 - 0.7)R$, t. x_1, x_2 are from
it is y_1, y_2 from P_{10} .

(But I have to go over this "proof" more carefully
no $2 \cdot 2 + \frac{1}{2} > c$ is what we need
 $2\sqrt{\frac{1}{2}} + 2\sqrt{\frac{1}{2}} > c \Rightarrow 2\sqrt{2} > c$ $4 > c$ is necessary
i. $\frac{\partial^2}{\partial x^2}$ to be > 0 in all directions.

So: conclusion of 712.27-717.21 : Proof to conv. form. for Radix < 2 probably by

shown by $\textcircled{1}$ showing Galt's proof for the DC. (1.1) \neq ! I think this is for any radix.
712.27-717.21 tries to show that " " implies \leq (arrange) : i.e. that

~~$(x-y)^2 < (x+y)^2$~~ for radix < 3 : The idea would involve showing that cross derivatives in $(.00-.21)R$ are smaller than the two "uncrossed" second derivatives (see $(.10-.20)R$). I haven't shown this yet. If true, it would be easy, using Maple to get derivatives. If false, then ~~show~~ this proof as is, it's no good.

Her. Nutter says he proved the conv. form. for any radix: so use his claim: its almost certainly true, indep of whether his proof is correct (Note that both eqns. compared must be "measured")

What I want to do now is write up the conv. form discussion for "J. report/paper(s)".

I'm starting by reading the discussion in "2 kinds" (599). The model of induction I'm using is a mixture of 4.2 models discussed in Sections 3.1, 3.2 & 3.3

A difference may be this: (if this argument may apply to the problem of 599 as well).

8/27/02

Maple: (SD)

Maple file **KL DIST. MS**

ID

> This is "converted" version!
> $x \ln(x/y) + (1-x) \ln((1-x)/(1-y)) + (x-y)^2$: ← not used. Maple Printout

D ≡ (> $x_1 \ln(x_1/y_1) + x_2 \ln(x_2/y_2) + (1-x_1-x_2) \ln((1-x_1-x_2)/(1-y_1-y_2)) - (x_1-y_1)^2 - (x_2-y_2)^2 - (x_1+x_2-y_1-y_2)^2$:
> diff(" , x1);

$$\ln\left(\frac{x_1}{y_1}\right) - \ln\left(\frac{1-x_1-x_2}{1-y_1-y_2}\right) - [4x_1 + 4y_1 - 2x_2 + 2y_2 - (x_1-y_1) + (x_3-y_3)]$$

This should be symmetric in (x_2, y_2) and (x_3, y_3) !!

$$\left(\ln x_1 - x_1 \right) \cdot \left(-(\ln y_1 + 4y_1) \right) \left| \begin{matrix} 1 & & \\ & 3 & \\ & & 2 \end{matrix} \right. - \ln x_3 + \ln y_3 \left| \begin{matrix} & & \\ & & \\ & & 2 \end{matrix} \right. - 2x_2 + 2y_2$$

$$-4(x_1-y_1) - 2(x_2-y_2) \neq \text{something} - (1-y_1-y_2)$$

$$-2(x_1-y_1) - 2(x_3-y_3)$$

$$\frac{\partial^2 D}{\partial x_2^2} = \frac{1}{x_2} + \frac{1}{x_3} - 4 \quad \text{its always } \leq -2 \quad \text{but it should be}$$

Sym in "1" & "2" indices (here, $x_1 \leftrightarrow x_3$, not $x_2 \leftrightarrow x_3$)

try $E = (x_1-y_1)^2 + (x_2-y_2)^2 + (x_3-y_3)^2 - x_1 - x_2 + y_1 + y_2$

$$\frac{dE}{dx_1} = 2(x_1-y_1) + 2(x_2-y_2) = -2((x_1-y_1) + (x_3-y_3))$$

$$E = (x_1-y_1)^2 + (x_2-y_2)^2 + (x_1+x_2-y_1-y_2)^2$$

$$\frac{dE}{dx_1} = 2(x_1-y_1) - 2(x_3-y_3) = +2(x_3-y_3)$$

So P_2 seems to ignore mechanism of symmetry! - x_1 is involved

w. $x_1 \leftrightarrow x_3$, but not x_2 !

But even so: $D = \phi$ for $x = \frac{1}{2}$:
but $\frac{\partial^2 D}{\partial x_1^2} = \frac{1}{x_1} + \frac{1}{x_3} - 4$

$$\frac{\partial^2 D}{\partial x_1 \partial x_2} = \frac{1}{x_3} = 2$$

$$\frac{\partial^2 D}{\partial x_2^2} = \frac{1}{x_2} + \frac{1}{x_3} - 4$$

$$x_1 + x_3 \leq 1 \quad \frac{1}{x_1} + \frac{1}{x_3} \text{ is max at } x_1 = x_3 = \frac{1}{2}$$

So $\frac{\partial^2 D}{\partial x_1^2}$ is always > 0 .

So is $4(\frac{1}{x_2} + \frac{1}{x_3} - 4)(\frac{1}{x_1} + \frac{1}{x_3} - 4) > (\frac{1}{x_3} - 2)^2$
Let $z = \frac{1}{x_3} - 4$ is $4(\frac{1}{x_2} + z)(\frac{1}{x_1} + z) > (z+2)^2$

What's min of $(\frac{1}{x_2} + z)(\frac{1}{x_1} + z)$ as a function of x_3 ? (fixed z)
 $x_3 = 1 - x_1 - x_2$

$$z^2 + (\frac{1}{x_1} + \frac{1}{x_2})z + \frac{1}{x_1 x_2} \quad \text{Drop this for a moment!} \quad \rightarrow 17.6.00$$

Get min of this expression w.r.t. x_1 & x_2 by partial setting

8/27/02

717.6



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 717.6 I took first partial wrt x_1 of $(\frac{1}{x_2} + \frac{1}{x_3} - 4)(\frac{1}{x_1} + \frac{1}{x_3} - 4)$ (Maple? KL dist 2.1ms)

$$\rightarrow (-\frac{1}{x_1^2} + \frac{1}{x_3^2})(\frac{1}{x_2} + \frac{1}{x_3} - 4) + \frac{\frac{1}{x_1} + \frac{1}{x_3} - 4}{x_3^2}$$

03: It = 0 if $\frac{1}{x_1} + \frac{1}{x_3} = 4$ and $\frac{1}{x_2} + \frac{1}{x_3} = 4$ ← solving this would also give $\frac{\partial}{\partial x_2} = 0$

$$\therefore x_1 = x_2 \quad \frac{1}{x_1} + \frac{1}{1-2x_1} = 4 \quad \frac{1-x_1}{x_1(1-2x_1)} = 4$$

$$4x_1 - 8x_1^2 = 1 - x_1 \quad 8x_1^2 - 5x_1 + 1 = 0 \quad \Rightarrow 8x^2 - 5x + 1 = 0$$

looks bad! $25 - 32 < 0$ so no real soln! (on the other hand, I know from symmetry considerations that

I might try to find way do det Maple to try to find ~~soln~~ num. soln. for $x_1 = x_2 = x_3$ (symmetry must be used)

$$\frac{\partial}{\partial x_1} = \frac{\partial}{\partial x_2} = 0.$$

or, just plot the function

formal: plot $\geq d(x, y, z) = -0.1 \dots 0.1, y = -0.1 \dots 0.1$

Maybe plot. 717.6.00!

Looks like \uparrow as at $x_3 = 0$ and $x_1 = x_2 = 0$

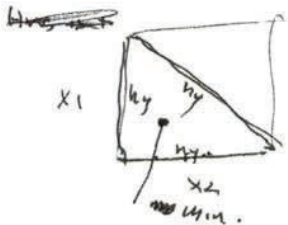
Look at $x_1 = x_2$: $(\frac{1}{x} + \frac{1}{1-2x} - 4)^2$ from $x = 0$ to $x = .5$

$$\frac{d}{dx} \left(\frac{1}{x} + \frac{1}{1-2x} - 4 \right) = -\frac{1}{x^2} + \frac{2}{(1-2x)^2} = 0 \quad x^2 = \frac{(1-2x)^2}{2} \quad x = \frac{1-2x}{\sqrt{2}} \quad (\sqrt{2}+2)x = 1$$

$$x = \frac{1}{2+2\sqrt{2}} \quad \left(2+2\sqrt{2} + \frac{1}{1-\frac{2}{2+2\sqrt{2}}} - 4 \right)^2 = 6.43$$

Because of symmetry wrt x_1, x_2 , ∞ will have a min on the line $x_1 = x_2$.

$(\frac{1}{x} + \frac{1}{1-2x} - 4)^2 - (\frac{1}{2x} - 2)^2$ has min at $x = .297$; ~~max~~ w. value 13.18 and ϕ .



I plotted $4(\frac{1}{x_1} + \frac{1}{x_3} - 4)(\frac{1}{x_2} + \frac{1}{x_3} - 4) - (\frac{1}{x_3} - 2)^2$

$x_1 = .1$ to $.4$ $x_2 = .1$ to $.4$ It looks like + min really @ $x_1 = x_2 = .297$

2nd its value is 13.18 : So the result is that $\frac{\partial}{\partial x_1}$ partial deriv of $D(\Pi) - \phi \leq 0$ and $\frac{\partial}{\partial x_2}$ is positive in all directions

Hm, I shouldn't spend any more time on this! — that is $\in \Sigma D(\Pi) < 0$

is really adequate, it can be easily shown for all radii using Calc's method (I think!)

But just assume Author's proof is O.K.

8/27/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

: Using Notation of S99: § 3.1-3.3: $P_{\#}(D_n)$ corresponds to μ (t. true & d.f.)

operating on D_n (which corresponds to A)
Maybe not use Post Notation! It will cause confusion: Say $P(A_i|Q_i)$ is t. "true"
generator of t. set of A_i 's (w.r.t. t. $\sum Q_i$ set). Then $\prod_{i=1}^n P(A_i|Q_i)$ is indep of order of the Q_i 's.

Also our approximation functions; $O_j^n(A_i|Q_i)$ will be indep of order of t. Q_i 's.

As a result $\sum_j a_j \prod_{i=1}^n O_j^n(A_i|Q_i)$ will be indep of t. ordering of t. Q_i 's.

The conditional p.c. obtained for a A_i part way thru t. set, will, hvr, depend on t. Q_i 's that have preceded it.

Hvr, for a particular ordering of t. Q_i 's Both $PC(1)$ (i. true-generator) & t. wtd set $O_j^n(\cdot, \cdot)$ will obtain a p.c. for each A_i & an indexed p.c. for each bit of every A_i — & a p.c. for t. termination of each A_i .

So we can use S78T3C: (w. radix 3) on t. pair of probly sequences t. bits of $\{A_i\}$ — t. 2 seqs are P.c.s obtained by .07 & those obtained by .03 (PC.1.).

From the argument of S99 p. 259 (col II), we get this con. lower bound w. constant factor 2.
So we get the final constant thru viz S78T3C.

I must mention in t. paper that t. probly of t. stop symbols S are an imp. part of t. Q A system — It is important to know if a particular A_i candidate ends

at a particular pt. or not, so ~~one~~ knows whether to start on t. next question.

I will use ~~the~~ Much of t. Exposition of CPM N8, 01 to desc. t. Q A soln (for $c = \infty$) & just why we are looking for O_j^n w. large $\left\{ \prod_{i=1}^n O_j^n(A_i|Q_i) \right\}$

There is a second eq. that gives lts of O_j^n & we want O_j^n 's w. much wt.

T: CPM Method is $O_j^n(A_i|Q_i)$

In S99 § 3.1, I considered all poss. ordering of t. Q_i elements. If the $P_n(\cdot)$ is order independent (i. it would be well to make it so, because t. "target" p.d. is order indep.) Then we may be able to use just one ordering & not sum over all $n!$ of them

The individual O_j^n 's try to simulate μ , + target p.d. & so they are order independent so the wtd sum of O_j^n 's must be order independent

Looking at ~~the~~ CPM N8, 01: on p. 2 2nd last P I say "t. distribution 4.14.18 is very accurate" I mite refer to Appendix, which shows this

It would be well to give an idea of t. result of t. appendix.

In Appendix ... we show proof: (I'll have to make this approximate)

Since my statement, here, will not include p.c.'s of stop symbols.

I can actually say that $\epsilon \leq \kappa$ because adding in "s" errors makes inequality stronger!

1908



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00 : **SN** Because the QAs are uncorrelated, whenever a new QA is added to the corpus, all of the p.d.'s for all off. ~~Q~~ A's change. Usually the corr. to ~~any~~ A_i is about the same as all other A_j's. This means error is $\leq \frac{1}{n} \ln 2$. If it's ~~not~~ .00-.02 doesn't contradict it.

02

03 Actually, I'm not quite sure of this. I original ~~error~~ $\leq n$ (bits) maybe not. Perhaps when I'm in the throes of writing the proof of the Conv. Thm. for QAs, I can return to this ~~QA~~ \rightarrow 907.20

04

05 (Suppose that there exists an intelligent algorithm μ , that can look at any question Q_i and give a ~~good~~ \rightarrow good

06 probability distribution $\mu(A_i|Q_i)$ on all possible answers, A_i . ~~if possible~~

07

10 "Appendix B" Convergence Theorems for QA induction.

replace with 106-107

08 Suppose that we have an algorithm μ , that looks at questions, Q_i and ~~gives~~ ~~for each such question~~ gives a probability distribution $\mu(A_i|Q_i)$ on all possible answers A_i , to the question, Q_i .
function defined by (1)

09 We also have for ~~probability distribution~~ ~~of eq (1)~~ ~~algorithm~~ ~~that also gives a~~ ~~probability distribution for each A_i in response to an ordered sequence of Q_i's.~~

10 Given an ordered sequence of Q_i 's \rightarrow (no limit), can we give a measure of how close the ~~probability distributions given by~~ ~~eq (1) is to~~ ~~are in those of the generating function μ ?~~

11 ~~There is a way~~ The Kullback-Leibler ~~distance~~ ~~is one such measure that we will~~ ~~use.~~ (Given reference to Cover-Joy on KL dist). If $\mu(x_i)$ is a probability distribution on objects x_i and $P(x_i)$ is regarded as an estimate of $\mu(x_i)$ then

12 Interpretation: the KL distance $D(\mu||P)$ is defined to be: (under uniformity of $\mu(1), P(1)$)

$$\sum_{i=1}^r \mu(x_i) \ln \frac{\mu(x_i)}{P(x_i)}$$

← this is the cost

It is the expected value ~~of~~ ~~with respect to μ~~ of $\ln(\mu(x_i)/P(x_i))$.
will (and costs) ~~decrease~~ ~~as~~ ~~the~~ ~~estimate~~ ~~improves~~.
That $D(\mu||P)$ is always ≥ 0 . ~~is~~ ~~discussed~~ ~~by~~ ~~the~~ ~~text~~.

30

31 Let us ~~define~~ ~~define~~ $P(A_i|Q_i)$ to be the probability distribution for ~~the~~ ~~particular~~ ~~sequence~~ ~~of~~ ~~Q_i's~~ ~~as~~ ~~given~~ ~~by~~ ~~eq (1)~~ for a particular sequence of Q_i 's.

32 Given this Q_i sequence, what is a good measure of the difference between the probability assigned to the A_i 's by μ and those assigned by P ?

33 Using the techniques described by Peter Dacs (PhD 1997 PP, LV 1992 PP) It is possible to show that

$$E \sum_{i=1}^k D(\mu(A_i|Q_i)||P(A_i|Q_i)) \leq \frac{1}{2} \ln 2 \cdot k$$

← since I used $\ln 2$ in LHS, I can eliminate $\ln 2$ factor

34 Here the expected value is ~~with respect to μ~~ and the summation, j is over all possible responses, A_i , to Q_i . k is the number of bits in the description of μ . "Actual bits".

Maybe $A_i \rightarrow A_j$.



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

Equation 722.25 gives a very simple ^{measure} ~~prediction~~ of the accuracy of a given ~~data~~ ⁽¹⁾.

There are no "order of one" constant factors or additive terms. ~~It is~~ ^A necessary uncertainty in the value of k . We cannot ever know its value. If the generator of the data has a

long and complex description, we are not surprised that ~~the~~ ^{we} should ~~need~~ ^{need} more data to ~~make~~ ^{make} good predictions — which is ^{just} what eq. 722.25 ~~specifies~~ ^{specifies}.

The value of the constant k , depends critically on just what universal/machine is being used to assign a priori probabilities to the O_j and to M . Any ~~part~~ ^{part} a priori in formula

that a researcher may have can be expressed as a modification of ~~the~~ ^{the} reference machine — by inserting low cost definitions of concepts that are believed

to be useful in the needed induction ^{usually} resulting in a shorter code for M (a smaller k) and ~~we~~ ^{we} believe that if all of the needed a priori information is put into the reference

machine, then eq. 722.25 is likely to be ~~the~~ ^{the} best ~~prediction~~ ^{estimate possible} → .30

~~(with no ~~assumptions~~ & ~~data~~ or multiplicative constant.)~~ [?] email.

Since ~~we~~ ^{we} cannot ~~use~~ ^{use} the infinite resources needed to evaluate ~~the~~ ^{the} ~~code~~ ^{code} ~~for~~ ^{for} M ~~is~~ ^{is} ~~not~~ ^{not} ~~available~~ ^{available}.

eq 414.18, (last # of CPA NS, 01 Page 2)

SN

Discuss how this is an approach to "STRONG AI": That "Common Sense" is not needed in ~~the~~ ^{the} long ~~term~~ ^{term} ~~and~~ ^{and} once a system has devised good (by ~~the~~ ^{the} ~~system~~ ^{system} ~~it~~ ^{it} can acquire "Common Sense" by long, or rather "Brain Surgery" (≡ ~~program~~ ^{program}). — (it's a much ~~better~~ ^{better} way to acquire "Common Sense". We want to use as little "Brain Surgery" as possible. Explain just Why.

T. report so far: 696.13 → .40

Intro: More or less done. (2 Abstract w. additions) ^{additions.} why getting these eq. (1) ~~code~~ ^{code}.

Section on QA problem. Use CPA NS, 01

Neuron at 719.19 → .33 addresses Accuracy of the "opt": refers to "Appendix B": Also add (722.16-17)

"Appendix B" is 720.10-722.15

Done in 0/6/02 version

.14

At first glance, this result may seem unreasonable: Suppose we ask the system many questions about Algebra, until its ~~error~~ ^{error} margin errors are quite small — Then we suddenly begin asking questions about Linguistics — certainly we wouldn't expect the small errors to continue! What happens when we switch domains suddenly, is that k suddenly increases. A M that can answer questions on both Algebra and Linguistics has a much longer description than one that ~~is~~ ^{is} familiar with Algebra only. This sudden increase in k accommodates ^{expected} ~~the~~ ^{the} larger errors in a new domain in which only a few questions have been asked.

→ 80% for many recent versions of Appendix B



00: : If it takes time \uparrow to generate and test the function G' that solves the problem, then $\exists T_0 \cdot P(G') \leq \uparrow$ and $T_0 \leq \uparrow / P(G')$. The total time spent on search is $T_0 + T_0/2 + T_0/4 \dots \approx 2T_0$. or $2\uparrow / P(G')$

05: Maybe best to copy section on L-srch from OSS 85, 86 or 89!

06: [SN] For f. talk, it may be part paper! Explain that the fact that a simple

probability distribution is able to evaluate

09: First state just what 4.1.18 is able to do: We use the same p.d. z_j for all kinds of problems induction problems, a no-matter what c. p.d. is that generated to A's from to Q's, the system for each case will eventually discover that p.d.

This is not amazing in itself. We just have to be sure that our z_j assigned to the p.d. that generated to A's, is not zero. This can be done by any universal

Turing machine or any Universal Algorithm.

That this system ~~is~~ useful is convergence, to generate an acceptably large z_j to the p.d. that generated the Q's. If that z_j is very small, the error in probability

will be very large for any reasonably sized data set.

The / Amazing property of fact of the system is ~~that~~ that by using a universal

algorithm, information which describes P.D.'s using function definitions that have been useful for describing past successful solutions to problems, we are able to get

24: acceptably high z_j 's for unknown future problem solutions, and ^{hence} vary small errors in response to questions

Number 2 is part of good induction (1) (right value is 4.14.18) (2) good z_j values (i.e. good unc.) $\rightarrow .37$

Getting the best poss. 4.14.18 will not be v.g. ~~the~~ production induction unless z_j of soln is large enof. Bob Note 734.10!

It would be nice if it could say .09-.24 in a very clear, concise way!

There are 2 aspects of the paper! (1) Logical correctness of the ideas.

(2) Writing it in a readable, understandable way.

So best just write. The practice of writing will gradually improve my expression, if I think about it. Also, maybe get Alex to make criticism of ~~the~~ intelligibility, ect.

The 89 paper wasn't written badly!

32: [SN] The idea that L-srch may be an optimum if all info is in P.D.: ~~the~~ "P.D." ~~is~~ ^{in this case} ~~is~~ ^{any} ~~is~~ ^{is} expressible in such

35: a system \rightarrow so we can (in theory) put all info into the "P.D." \rightarrow the system is potentially optimum!

I. optimization of the (LCP, LPST) is discussed in (724.00) \rightarrow (728.19)

Start out by writing just as if I were explaining it to a person!

37: [SN] Re: .24R: was good (large q. is by low ~~cost~~ direct (each time) \rightarrow large z_j by large update cc.

The Q of best allocation is "solved" by "50% soln". (I could put this in talk's give "proof" that "50% soln" was not bad. (But note 732.10!))

9/7/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: (Spec 727.40) Re: LA & Core Cancer type problems: In "core cancer", the usual way is to find a good set of subgoals that lead to ~~the~~ ^{top} goal. Each subgoal is good if it depends on predecessor or other good subgoals. The value of a subgoal is if there are many ways to achieve it.

— if ~~the~~ roughly ~~is~~ open to many goals, there are many ways by PC paths to top goal.
In the Dynamic Pgm / solv of Core Cancer: There are too many paths to evaluate, search. ~~The~~ no. can be reduced by putting many "outcomes" & "responses" in equiv. class; so one only has a few possys at each step of the "LA" problem.

In general: to try to solve "core cancer" consider how a human would go about doing/planning, such a project: Initial epidemic studies to get ideas as to what's relevant.... (to start off).

At some pt. one must decide if problem is or is not within capabilities of present science (e.g. transmutation (Pb → Au) in Middle Ages) ~~v.s.~~ doing it now.

19 (727.29 726.35) Back to 726.32-34: To what extent can all hours be discovered by the system and (equiv/used) by the system? (By "System" I mean 726.32-34)

The idea of the optzn. of (GCPD [PST]) being able to acquire any hour.
I had idea that any hour can be expressed as 2 PST. Not true! Most hours are Modifus of PSTs. As such, they should become part of the Grammar of (GCPD [PST]) ← (A stochastic language). — i.e. a hour can be an "element" of a "grammar".

The "Big Problem" is how this ~~relates~~ relates to 729.00-08 — which is how the Grammar is optimized. How do "good hours" help make ~~it~~ work

"T. System Grammar" (i.e. the system), better?
028 may be the most important possl. Q about the system!

Can hours be used to improve the "score" of GCPD? ← Well, hours are ways noted in observing PST's at work. They do help ~~to~~ predict what "score" a PST will ~~have~~ have in various problems. Why should a hour that tells which PST's are v. in a particular situation, be better (more sought after) than a "hour" that tells which PST's are particularly BAD?

Perhaps it's the Basic Problem in the all-over feedback flow of the entire system, ~~the~~
T. last stage of GCPD, → GCPD₂ is "open loop".

← a ∈ PD₂ is PD if it has functional PD core.

(Spec 729.00)

9/7/02

729



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: (Spec ≤ 28.40): "Stage 1" of data only tries to get good (prod) of G_{PC} for various P_{ST} 's for various problems

~~It is not particularly oriented towards by Gars.~~

In stage 1, TM constructs this ^{stoch} Grammar which contains (perhaps) all past PST's. A hour will modify this Stoch Grammar & its PC ~~parameters~~. A hour that worries about the low Gars region of the Stoch language space can non-fl. less be useful in helping to reject low Gars cand. in stage 2.

Perhaps a hour could be ~~an~~ an other "New" PST, rather than ways to modify old ones. Well, any thing that modifies the Stage 1 GCPD would be equivalent to a "Hour" - a giving reasonably by-PC's to certain PST's w/low PC's - ^{formerly} its Modifn. of the Stoch Grammar.

T. Q is: Say TM (discovers/recognizes) a hour in a set of problem forms. Show in detail how this modifies GCPD₁ & GCPD₂ in a way that uses the hour's we would expect (like) to.

Perhaps regard a Hour as one kind of Modifn. of GPD₁.

Logical Reasoning Can be part of a Hour. Does this Modify the Q of .12-.13 in any serious way?

At present, my idea of a hour is rather foggy! Needs ~~rethinking~~ ^{Reviewing}.

20: Well, as simple example of hour: we have, originally a P.D. on PST's, CC's, G_{ij}. We find (heuristic) that if PST_i is restricted in a certain simple way, we get much higher G_{ij} for small CC_{ij}. This modifn. of GCPD₁ will produce a more "productive" GCPD₂. Hvr, ~~these~~ modifns of GCPD₁ of this sort are made (optimized) w. a defmt goal (not getting a good GCPD₂, so it would seem that the optimization of

30: GCPD₁ would (lacking this knowledge of what it is to be used for) would not be able to work as efficiently as it would if this knowledge were somehow available in constructing GCPD₁. One pt. is that w.o. this ^{previously} knowledge, the construction of GCPD₁ is

legitimate.

35: One conceives of ^{another part} Meta heuristic variables to this system! It seems likely that there are certain improvements (e.g. new PC or corpus) of GPD₁ that would produce very little (burr) betterment of GCPD₂, than others -

40: So one should concentrate ~~more~~ (improving GPD₁), not nearly on ~~max~~ max ↑ of PC or corpus per CC expended, but on ↑ its PC in a way that "feasibility" of GCPD₂.

45: → There is a danger in .33-.35 of biasing the search for codes in an illegitimate way!

230.40
Spec

ID



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

729.40: Hvr., it would seem that there would have to be Logic ways to work on "parts" of GPD, w.o. introducing ~~less~~ undesirable "bias". E.g. say one wanted to expand the part about "Physics" more than part on "Math". (The ~~parts~~ ~~2~~ ~~are~~ ~~nearly~~ ~~related~~ --- perhaps all parts of the GPD, are (potentially) related. — By ~~strategic~~ ~~2~~ (most ~~any~~ ~~algn.~~ But conversing one to work on certain parts of GPD, more than others will necessarily give undesirable bias.

There is a More General Q: In making induction codes for any corpus at all: what kinds of modifications of ~~the~~ ~~code~~ ~~to~~ ~~be~~ ~~such~~ (essentially Logic) can one tolerate

w.o. essentially harmful bias introduction?
— Consider ~~the~~ ~~repress~~: Say we start w. linear to order 5. Next try quadratic w. up to 5 roots. Then cubic to 5 roots. . . . The choice of "orders" biases such: the choice of which quad & which cubic term, forces, biases such.

Bas: 1) pr seq's in the code pres

A big Q (maybe main Q) is Signific. alt. biases: their "Payoffs" —

Positive or Negative. IN SM, bias toward use of linear Models is counterproductive because competitors use linear Models. IN ANN, except for very elementary models, the ~~var~~ variety of model types ~~is~~ rapidly, so that it becomes ~~less~~ less likely that the "competition" is using it. Similarity in a quadratic, cubic, etc, models.

In Using pseudo. Random nos. for Momentum in interpretation, say, we are concerned not nearly w. randomness, but w. "orthogonality" of the p.r. functs that find. being integrated (Whatever that means!) (I think we want zero cross correlation) In some cases, use of very simple, rather regular "p.r. functs" is adequate.

So, in doing GPD's search for "yields" to "maximize" (19.18!) What kinds of biases can we introduce, that will not ↑ GPD's yield. (b) give p.e. across that ultimately ↑ "yields" of TM,

Consider $c = \infty$ for GPD: This is Gal's vid of kind of Bias. For this case, all known have been discovered, at the best PST's are known. (All discoverable known (limited by content of corpus, $c = \infty$) have been discovered).

This would seem to be the best one could do. So approximately $c = \infty$, by getting the equiv. of " $\Phi(19.18) = \text{max}$ " would seem to be one way

Actually, the situation is less understood: All solutions to all solvable problems have been found. Hvr., the PD assoc. w. a given [the A2] set, will still be a PD & have accuracy limited by the convergence theorem

732.23

(see 737.00)



00: 732.40 Search should be organized ("bread") in various directions, must be somehow decided on by "Designer" (?)

03 T. man outstanding Q about TM was 728.19-30 (2 ~~725.32-35~~) (Also see 724.00-08)

The recent ~~work~~ work: (reviewed on 732.00-03) seems to mirror or less "solve" it? Note: 734.00-09
A Great Simplification!
Note 732.38 on just how to ~~discover~~ discover/implement, "heuristics":

06: 724.26 [SU] on Lsrch; In the short Method of 724.17, in which the Guiding PD is changing during the Lsrch. We may have to use a "Random" method of Lsrch of 724.12, because of the way that info about GPD₂ is obtained. If so, we could just continue our brutal Random Lsrch as "normal" - even tho the GPD₂ is changing during the Lsrch. This is not as "efficient" as 724.17-26 because we waste much time working on Cands that already have too much work done on them. Just how inefficient it is ... I don't know. I may be able to ~~comp~~ estimate,

hvr.
A (Guess stabs in the dark): This may be no worse than 2 times as slow as doing it the "correct" way.

An Approach to Analysis: Say we have a certain ~~state~~ (state) p.d. for 1 second and we do Random Lsrch for 1 second (or deterministic // Lsrch for 1 sec.) Then the PD suddenly changes (for 1 second). We can either just continue random or // Lsrch w/ the new PD. Or use fr. "better" method of 724.17-26

How much "Better" is better method?
If the PC's change very slowly, I suspect that there is little difference. The "wastage" occurs when we work on Cands that already have too much work done on them - whose PC's have recently ~~to~~.

24: ~~724.26~~ → A (perhaps) better way to do the random Lsrch: when a cand is chosen, we look at how long ago it was chosen last. (This gives estimate of present PC) we can track this w/ total cc spent this far on that cand. We also know total cc spent on all cands must be so we can estimate whether this cand has "too much cc" Thus far, & skip it if it's "too fat"

[We also may note several recent times at which this cand was "hit" - to perhaps get a better PC estimate. Say last 4 or 5 accesses: (or fewer when we've just started on a problem.)

24 ~~could~~ could be better than "normal" random Lsrch. We can adjust amount of FCC we use to "bring the cand up to standard" - so sometimes, it may have less variance than the "normal" method, even if the GPD is static! At any rate, 24 looks like a v.g. way to do Lsrch!

It would seem that imp. stuff done since 728.00;
724.35-726.05: Exp. of Lsrch → Proby for Paper (at least, my version).
726.06-30 on the Table - Thats why Lsrch may be Optimal: also 727.25; 728.19; 732.00-03 Bibliography of Shot 5.6 off
726.37: Discussion optzn. of 4/18 u.s. on unseen of 2/3. NoB 734.00-09! 218 732.03
Actually: 2 aspects of optzn. of 4/18

9/12/02
ED



00:

All Yes! In Sol 86, 89; I felt that some hours (but not all... Not Quite about, say, or (imp. during Lsrch) could be implemented by reordering trials.

More recently, the idea that modifying, inventing PST's could implement any hour, seemed. Latest idea: that modifying the d.f. over PST's instructively includes all possible PST's:

So any "invention" of a PST simply consisted of assigning a reasonably big PC to that "invented" PST. This means that modifying the P.D. can do all things

- 1) Modify order of old PST's
- 2) Introduce ("invent") new PST's.

So my ~~old~~ old idea, that all hours were expressible as modifications of the P.D.

09

726.25
726.24R
726.37

was literally correct!

[SN] Comments in "optzn of 4.14.18". In 4.14.18, a_j is fixed, for a given problem. When we change problems for 4.14.18, the "j" of a_j may change, but all of a_j remain the same; this is to Appt.

On the other hand, consider the P.D. assoc. with Lsrch. for a PST for a given problem.

It will vary w. TM's "Experience." So how does the P.D. for 4.14.18 differ from the PD used in Lsrch.? Well: they are entirely different PD's: The one used for 4.14.18 is

another an apptd part of GPD. The one used for 4.14.18 induction is GPD.

20

There is a way for us to change the "apptd" for 4.14.18, however. — That is via

The "Summary Machine". — In which the new "apptd" is a summary of the corpus up to now. Strictly speaking, we can't change the summary

"apptd" up to now. The only way to change it is to "backtrack" — which is equivalent to not really considering the "summary machine" as a complete

summary. Equivalent to some degree of "backtracking" would be to retain in memory, a fair no. of non-optimum codes for the corpus, rather than

"significantly different" from the apparently "best" codes.

30

One trouble w. "Summary Machine" is that they are (usually) not universal d.f.'s. i.e. we construct them by fitting a small set of functions to the past data. This makes "backtracking" necessary and expensive.

[SN] I discovered (twice) that all Lsrch is much more easy 2 times faster than $T \leftarrow 2T$ Lsrch. For dec length n , it's 2 a times faster (I think the in efficiency of $T \leftarrow 2T$ is not that

In decoding strings a_0, a_1 , it involves a_0 & a_1 independently. All Lsrch decodes a string containing "0" & also with "1". Both seem share the coding of "0".

How this would work out using random Lsrch, is unclear.

40y



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00 : An approach to writing this paper: I have "outline" of 696.15-.40: Make progressively more detailed outlines. Indicate what parts have been written (or discussed) and where. One advantage of this: I'll be more certain to keep whole system in my mind. So I don't forget. Probly D: RAY.

07 § I Q/A Problems: CPM Nov 8, 2001: The part now down as a Q/A problem. Also, I have Appendix w. conv. form. CPM Nov 8 [201.23 suggests that I explain how Numerical T.S., bag industry, Ency. can be covered by this formalism] 719.34-40
719.19-.33 should be added to 719.34-40
Appendix B: ~~720.10-722.15~~

15 § II PST's: Introductory section on PST's in general - derivation of various sub parts. (T. force should be $\frac{1}{2}$ \equiv 722.25-.28 so that ruffly covers the Q/A problem.)
2) Derb Lsuch: Perhaps ~ to Sol 86: Give 2 or 3 GHT's and perhaps Monotonic T-first Scaling problem - or perhaps how sections on "Scaling". Discuss Optimality. If all info is in P.D. 2 aspects of how P.D. can express hours! TEST T < T marked.
that ii: reordering of trials, ii: creation of new (PST's), by giving more parts to PST's that has $\approx \phi$ pc before. How this sort of implicit Lsuch could be optimum. 736.00
(How) Can we put "all info in P.D."? Discuss updating after each docking.

Low prob method
2 Problems
outline
of
exercises
count, such
align.

22 Some ideas to write Expo. papers about:

- 1) Realization of my work system to "C/C". My system does "C/C" as one mode of exb. solving.
- 2) How many universal d.f. is not desirable - that it would be of no value. (The historical evolution of Apmpd t. bio/ly 1721)
But one normally uses "summary machines" not "simple" UMCS.
- 3) No Free Lunch Theorem. But it does make lots of difference as to where we place UMC you use. This fact was not noticed by Math community because they were interested in potentially infinite $\epsilon \epsilon \epsilon$. I'm interested in small $\epsilon \epsilon \epsilon$, (ϕ 1, 2, 3).
That i. ("Princ. of indet") Uniform d.f. is usually meaningless, since deciding whether a single string should represent has to be decided by user. I use uniformity over UMC input for my "Price of indet". There is still choice of how to wrap output strings to R.R.
- 4) How to improve GA continuous improvement of best/crossover (2) Realization that UMC, cross are just 1 & 2 $\epsilon \epsilon \epsilon$ cases of induction. (3) Ability to process input problem - to suggest virtual population is appropriate mut/cross.

9/16/02

737



00:736.40 : § III ~~CONVEX~~ GCPD : 704, 28-29; 72335

These discussions are rather brief. I will add in some ways to resolve

- 1. C.P.D. (1) T. 3 input CPD (this is easily made into a ^{1.5} McCulloch CPD)
- 2. & output is list of ^(P, A) pairs in \approx P1 order.

Intro to talk : (Maybe refer to a short or intro of paper & Abstract)

I will be going to talk about a system that I've been working on latest developments in a system that I've been working on for some time (Sol 86-88, Sol 89)

It is meant to be an approach to strong A.I. —

A system that will surpass man in many areas of science and mathematics and perhaps art. It differs from ~~other work in this area, it's just~~ ^{other work} ~~in this area, it's just~~ ^{as in how}

I am not ~~so much~~ ^{so much} interested in knowledge itself, ~~as in how~~ ^{as in how}

it is required — how machines may learn. To start off, the system will learn to solve ~~two~~ ^{possibly} very general kinds of problems. Most, but ~~not all~~ ^{possibly not all}, problems in science and mathematics ~~can~~ ^{can} be expressed in these forms.

The second kind ... ~~is~~ ^{is} learning a theory to fit data, surface reconstruction, inductive inference ~~of~~ ^{of} ≥ 11 kinds, over ≥ 100 variables...

Perhaps don't get into "Strong A.I." — this disturbs many people.

Scientists resistant. ^{start at 1.35} It is meant to be a "Scientist Assistant" of great power and versatility in many areas of science and mathematics.

It differs from much other ambitious work in this area, in that ^(then start at buying @ .16)

→ Perhaps have more in this direction in intro.

176 ... "To start off Perceptron system has very little knowledge, other than algorithms for learning, and a very ~~simple~~ compact probabilistic model of the world.

"To start off, the system has a very ^{simple probabilistic} ~~simple~~ model of the world, and a very general algorithms for learning. ~~The~~ ^{It} then learns to solve ~~two~~ ^a ~~very~~ ^{single} ~~general~~ ^{general} types of problems. Most but possibly not all,



20

: [SN] Unclear in my mind as to exactly how system works. Unclear \rightarrow How they do it:

Say a problem comes in: It is either a INV or OZ problem, (if it's a INV problem, it has to be

first converted to OZ). GPD₁ immediately processes it as PD on PST's, and we do something with L_{search}; then update GPD. As TM evolves we modify L_{search} by updating dom L_{search} & by adding new dom L_{search} to it.
A better way: Prob. comes in: we update relevant part of GPD₁ then we use it to solve the problem. \rightarrow (31)

04

09

0

11

20

24

30

31

36

[SN] On a proof of conn. from for Bessy/QT! Perhaps I don't need as much

"new stuff" - Gac's proof for conv. of four measures might be short!

Say I consider $P(U, K_1)$ to include not only "undefined" but stopping.

Wally, Bessy may not help: we do need to stop symbol - even if machine does U at end of "A_i", B_i is 0, 1, but we need to know its set 0, or 1.

If we have a string $A_1 A_2 \dots A_n$ concatenated w.o. spaces between them, then A_i is \supset previous part of a machine producing & say " $A_1 A_2$ " as output. T -space at end of an A_j is an essential part of the output.

O.K. So use Gac's system w. "U" a "#U" comes up, code string "occasionally"

We still have error \leq ϵ (probly 0, error \leq (probly later) $< k$.

We use Gac's proof for 0 & 1 only: we don't have to add in error for U

(A trouble is my "proof" is Prob just as I show pc^2 of $A \rightarrow 0$ rapidly, This would show that pc^2 of $(U + \text{something}) \rightarrow 0$; which is false!

There may be a way to get around .24 - having "and" occur many times in the past

may apply to conclusion in the conversion!

But don't spend time on this now. T -present proof is adequate for time being.

I should get this clear in my mind! (I.E. Just How TM Operates)

One way of thinking about it: As soon as problem comes in, TM notes if it's a INV or OZ problem. (if it's a INV problem, it's converted to OZ by an Input perm.). Perhaps at this point a GPD₂ distro, is created from GPD₁.

Another way: I just do \approx L_{search} every time; we update first (both GPD₁ & GPD₂) -

The nature of update depends on progress in prob that over of prob solving \approx in TM in general. First no update second, update w. simple corpus of successful searches (GPD₁ only) but not quite is a partial corpus used by GPD₁

(3) GPD₁, GPD₂ but no correlation + GPD₂ at end of $T \leftarrow 2T$, look for condt.

9/17/02
5D

"CONTEXT" a new view (looks v. g.) 2.22

739



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 738.90 : ⑤ GPD, GPD2 more frequently - rather ^{Time/for} during Lsrch ~~the~~ ⑥ Updates to
watching current trial - so uses "correlations" to decide on next trial. Use of random Lsrch - possible w.
⑦ No Lsrch, update is always simult. w. trials. Jump to new trial when "expected
remaining completion time" \leq (E, R, C, T) of new found trial if less than E, R, C, T of present trial - screens
⑧!

03
- 1.36 ft is fine, but I need more of a detailed idea of what "PST's" are being ~~found~~ searched, PST's use
In ① T. [PST] sets a spiral given by the Lsrch language. ② we have an out cases of ~~that~~
induction there to begin induction on how to modify "guiding PD". In ③ Our induction facilities must
be good enough to get reasonable P.D.'s of cc as a function of problem domain; PST; (so do this
before we start each problem - since we have to know what the problem is. ③!).

10
In ④ thru ⑦ we need progressively better update facilities - our induction past
problems must be progressively better. \rightarrow see 767.36 for how to learn how PST's
 \rightarrow how to deal w. (18.21)

17
18
Re: situations 1, 2, 7: could I devise samples of problem types & of skills
that TM has acquired?

18
At one time in p2st, I considered adding (ad hoc), a bunch of PSTs that I
felt were v.g. - ~~the TM~~ put them in "factored" form, & I have TM

20
"to enter" them & devise a stock lang. using them as corpus.
Would this fit into large schema - if so, how? \rightarrow spec. \rightarrow 741.00

21
22
Another idea that I'd like to fit into the large schema - the use of "Context"
in constructing trial induction functions. (certain kinds of "Context" are normal
regularities in a corpus - it's a way of defining a type of "Bern. seq. ..."
If something is in several different contexts, we can use that trick of mixing non-producers.
wt. of each context = $\frac{\text{dom}(\text{cost}) \times (\text{SSZ})^{\text{wt.}}}{\text{SSZ}}$. T. prod. dom may perhaps be done by
the no. of legal contexts at that pt. that have been defined thus far.

30
My original Motivation for "Context" was that w/o it, the pc's of corpus ~~became~~ became very small, due
to the number of them (as corpus \rightarrow 20). "Context" ~~was~~ implies small section of
corpus being used for prod. It means small SSZ, but perhaps pc's closer to ϕ 's since
the contexts are defined only when they \uparrow pc of entire corpus.

\rightarrow For END problems, this "Context" ~~was~~ will amount to a modification of how v. a prod is defined -
It is something that can be "added in" later on, as TM develops. Normally the pc of a conc.
depends only on its past frequency. The ability to define contexts gives a different opnd, just as the ability
to define ~~new~~ ^{functions} modified to opnd.

\rightarrow 741.00

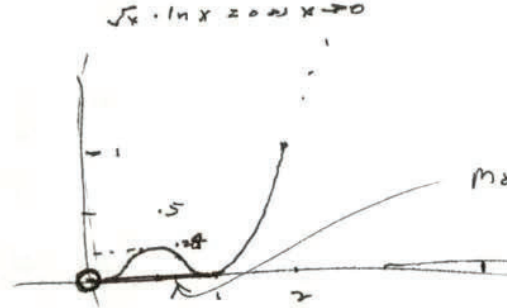


JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

For the spread of roots from $-\infty$ to $+\infty$: [there is considering values of
Coeffs for linear (or n.l.) regress. — ~~the simple prob.~~
We expect \downarrow from 0 as $x \uparrow$, The first "corner" is an imp't param: say it = 1 for coeffs.
Since we expect roots linear coeffs do have 1. (Non-linear may be different —
may depend on rms value of signal/coeffs).

The part from 1 to $+\infty$ can be like $x \cdot (\ln x)^2$.

$x \cdot (\ln^2 x)$	1
0	0
.5	.24
1	0
2	.96
4	7.7
8	37.6



so $(1 + x \cdot (\ln x)^2)^{-1}$ would not be bad.

Maybe may it ϕ for $0 < x < 1$

Algorithm ϕ is ratio of ϕ for large x . This determines values for $-1 \leq x \leq 1$

So int $(0, \infty, x \cdot (\ln x)^2)$ = ?

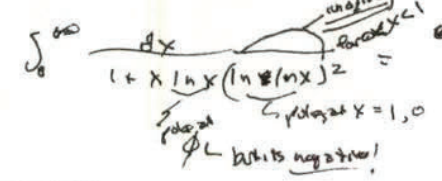
int $(x \cdot (\ln x)^2)^{-1} = -\frac{1}{\ln x}$ $\int_0^{\infty} (x \cdot (\ln x)^2)^{-1} dx = -\frac{1}{\ln x} \Big|_0^{\infty} = 0$ Unlikely!
 in integrand is always 50. $\int_0^{\infty} x \cdot (\ln x)^2 dx = \frac{1}{2}$ (cancel out)

oh! Integrand has poles at 0, 1, ∞ .

int $\int_0^{\infty} (1 + x \cdot (\ln x)^2)^{-1} = 2.804281080$ so $()^{-1}$ is the kernel const $(\approx \frac{1}{2})$ 3.566

So \int_0^1 normal funct. : ≈ 3.566 ; \int_1^{∞} normal funct $\approx .6434$ } (both Maple & whitt to do this S.)

$f(x) = \frac{1}{x \ln x (\ln^2 x)^2}$ $f'(x) = \frac{1}{(x \ln x (\ln^2 x)^2)^2} \cdot (\frac{1}{x})'$



It didn't find plus calc'n! ~~try 2~~ \int_0^{∞} example.

dec: $x \ln x (\ln^2 x)$ ever = -1 ?
For $x < 1$ $\ln \ln x$ is imaginary!

It is my impression that R_{11} $x \ln x$ $\approx \log^2 x$ has lots of wt. in by x values —
so much less on $[0, 1]$ interval.

For non-linear (poly) regress: x coeff or x^2 should be multiplied
~~multiplied~~ divided by x . ~~multiplied~~ rms value of x . coeff of x^n divided by $(\text{rms of } x)^{n-1}$.
 After division this way use a prop of .00 ft $(x \cdot (\ln x (\ln^2 x)^2))$.

 $\rightarrow \begin{matrix} 767.00 \\ 765.00 \end{matrix}$

9/19/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

oo: (5 part 739.21) IT should be possible to "add in" PST's. There seem to be 2 models for [PST] set:

- ① to start with poss. PST decs ... a universal def. on them given by a universal Σ fun.
- ② A finite set of PST's, initially inserted by designer/trainer.

We can start w. ① and "add in" a fact set of PST's: T. w/out. new PST's will have to be assumed relative to PST's in ①. The Σ of all Σ 's need not be Σ — we automatically renormalize w. random or deterministic (or both).

How the "update" would work, is unclear. One way would be to keep 2 sets of PST's, separate & update them separately. — This is certainly NOT in the spirit of

~~The~~ Universal "Transfer lang".

T. relative w/out. 2 sets of PST's could be adjusted so as to make pct of corpus. (?)

On the other hand, we have 2 kinds of ~~prob~~ objects: The O_j^i of 719.18,

& the PST's. The O_j^i 's are supposed to be for the entire corpus. T. PST's are for specific, single problems.

It is possible that we are looking for a single PST that is same for all problems

~~the~~ ENV, $\Delta \Sigma$ & IND — like ENV prob had a single sought for PST.

In 739.18 ff, the search for PST + factors \rightarrow a (Universal) ^{stoch} grammar. . . .

In the idea of GPD giving a PDec to PST's: There is a no doubt here, since

we will "eventually" need only 1 PST, — but we've split the problem by breaking it down into 2 parts — i.e. ~~the~~ GPD is [PST].

{ We can fix the present system either way: as ~~we~~ looking for one PST for all problem ~~or~~ a GPD that assigns \rightarrow diff PST to each new problem.

We can do some kinds of problems (say ~~the~~ IND) one way & 2 (ENV = other $\Delta \Sigma$ probs) the other way: look into this.

Probably it would be good to look at each "PDec" system. I've perhaps put most of 2 GPD on (Prob. decs, PST) rather than on PST indec (or average over all) problems.

A PST chosen w.o. consideration of problem, must itself look at the prob. decs & categorize the problem to decide what method to use to solve it.

In the ~~the~~ original system, the GPD looks at prob. decs, so the PST ~~is~~ has to more narrowly re-categorize the problem for solution. ~~There~~ There is less intellectual "work" that need be done by the PST.

9/19/02

TSQ's .00 continues to 745,23

742



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

74150 : Consider algebra: simple MTM induction: Max of 4 (4-18)

T. way we do it is just look for standard "regularity" in ϵ -corpus.

At one run in w/ ANL, I had " $3+7=10$ " $Q=(3+7=)$; $A=(10)$.

03 $3 \& 7$ are known to be numbers. There are "sum, sub, mul, div" binary operators. "neg" is unary.
"t" & "=" are "unknown" symbols. 10 is a numeric output, so we first look for unary

06. So these solve it. Hrrr, maybe there is "reasoning" being used, & TM doesn't know how to do "that sort of thing". Until TM learns to do that kind of reasoning, it is

"Forbidden"

Computer PC of $3, 7$: pre-comp $3, +, 7, =$, sum, sub, mul, div, neg, !, 0.

"symbols: sum 3,7 $\rightarrow PC = \frac{1}{11} \cdot \frac{1}{12} \cdot \frac{1}{13} \rightarrow \frac{1}{1716}$

since there are 2 ways, total PC of both $\rightarrow \frac{2}{1716} = \frac{1}{858}$

But it has to do $3+5=13$ also:

well, we can let $3, +, 7, =$ be the first, second ... fourth symbols in Q .
In which case it works for any $a+b=c$, but the result is critically dependent on just where symbols are placed. It wouldn't be able to guess ~~the right order~~.

$Q = \square$ about $3+7=$: 5 symbols in Q .

04 Hrr, the idea of translational invariance would have to be hard. No way around it - unless we build it into "TM".

9/20/02 Another way to think about this: T. General Heur implemented by Levin "AM".

Th. Q was, how well they fit together in the "learnability" of new concs/heurs.

A possl. way to do this: to "give" TM the idea of translational invariance (or fix it so that

transl. invariance is not essential (?). - No recognizing a "a+b" configuration does equal "a+b" invariance concept.

Also give it "reasoning" concs needed for .03-.06 : .03-.06 is a "working backward" -

a common heur. [we need numerical result, 10: to get it we need to use fixed numbers (10) as output: list possys. Try simplest possl. inputs to start: (if they don't work, try another level of functions.)

3) My idea is that "all info can be in GPD" - so that if Levin can't find a soln. in reasonable cc, it's because heur was imp. into missing part of GPD. T. main simple ways of putting info into GPD:
1) Defns. of fixed.
2) Defns. of contexts for functions
3) Some I've not thought of.

Very imp. idea re: design!



90: 7+2=9! One way to implement TSPQ's using V. (nacy?) ideas of 742, 31-40:
Establish various "state pts" along t. t.sp. toward same goal. This is \approx like the idea of conc. nat. but instead, we ~~write~~ write these ordered set of state pts. - then we try to bridge ("connect") betw. them.

If it proves diff. to implement, ~~write~~ write these ordered set of state pts. - then we try to bridge ("connect") betw. them.

to GPD has to (be modified/ fixed/ replaced).

[SN] Another kind of PD that TM does is Perm Seq: defined by 6-2 (phases of objects, to no. of times they occurred, Th no. of opportunities for its occurrence - register.

→ A list of ~~state~~ "state pts"

- 1) doing $3 * 5$; then $8 - 7$ ~~problem~~ ... Permutation problem (Perm Seq) model
- Then $7 * 8$, $4 \div 16$ " " " " "

3) ~~Identifying~~ Identifying "relevant context" so that "+" is recognized as \rightarrow Sum, act.

4) ~~Being~~ Being able to do paren-nested expressions.

5) Maybe being able to do Alg expressions w. some parenthesis - using hierarchy of functions - say to learn to predict freq. of a, v.s. b ; then a v.s. b, v.s. act. We write macros & "builtin" conc. -

So TM would automatically look for duplication of symbol seqs, so it could deduce Perm Seqs.

How to do 2) (Perm seq. of $3+7$, act) \rightarrow 3-ization of ...? (corvins?).

We have essentially a new seq of QA's: $Q \quad A$
 $3+7 = : \text{sum}(R_1, R_2) \quad | \quad 8+9 = : \text{sum}(\dots)$
 $9 * 3 = : \text{mul}(R_1, R_2)$

From Perm, it should be not hard to get "+" \rightarrow sum; "x" \rightarrow mul, but details have to be worked out.

From 2) $3 \left(\begin{smallmatrix} x \\ x \\ x \end{smallmatrix} \right) 4$ to QA (nested expressions): is not apparent.

a) ~~A~~ A poss. way: it first learns $3+5=8$; then $8 \div 2=4$ then $(3+5) \div 2=4$.

b) a/o say it learns $2=3+7$; then $2=10$ then $2+2=12$
8! i.e. T. concepts of "Equality".

36 ... from 29: an earlier approach: ~~to~~ teach meaning of "eval": $\text{eval}(3+7) : 10$

Then $\text{eval}((3+4)*7) = \text{eval}(\text{eval}(3+4)*7)$

[SN] Another way to deal w. ANL: How it be "builtin": i.e. its just part of our way of communicating w. T.M.! . Hrr, I wanted "ANL" as an example of a simple learning problem.

Perhaps its not a good example! Actually I could do say 4. decisions as to whether to state w. ANL, by working at higher level More advanced Alg ~~to~~ learn problems.

9/4/00

GA.00

Theory of how mut/cross over $\leq 2, 2$.

747



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:

\bullet **SN GA:** "Mutation" \equiv interaction w. $\leq 2, 2$. How to do it, etc. : from X_0, G_0
 to be mutated: Construct P.D. on X_1, G_1 X_2 is rand. descnt, G_2 is its possibl. G-value.
 chose $X_2 \ni G_2 \ni$ Max.
 For crossover: Given $(X_0, G_0), (X_1, G_1)$ again construct D.f. $[X_i, G_i, P_i]$
 (chose $X_j \ni E(P_j)$ is max.)
 T. above choices are for **Max Good** Modifiable it to D.R. to be "concerns" in no
 account - i.e. **Joint PD** of all possibl. contngs to f. set of $[X_i, G_i]$

0:

(Spec 743.40): In General, I should be able to state each hour I use/heard for a particular
 problem soln., then (express those hours as Moduls. of GFD, 2/0 as Tsp's Dist
 world bring about that GFD modula)
 T. clear statement of X. Hour, is perhaps the hardest part.

15

in ~~the~~ problem of -1.29 (\equiv substitution): Substitution is a rather complex ~~conc!~~ conc!
^{primitive}

16

It may be one of the funct's that Kleene ("that Mathematician") uses to define a universal Alg.

Perhaps I can teach $3 + (9 \times 5) \rightarrow 3 + \text{eval}(9 \times 5)$, \equiv similar examples.

20

A ~~way~~ (Risky) ~~way~~ to teach TM: Have it run rather ~~little~~ to solve rather diff't
 probs., ~~by~~ by inserting hours that humans seem to use to solve those hours
 problems. Then, perhaps TM would be in a position to discover more complex hours
 This is ~~an~~ approach of Landin ~~Evrisco~~, but we know more about ~~the~~
 (sample ~~of~~ nature) of examples and amt of cc. needed to find the hours.
 It's risky because it's harder to tell (if you have missed some (perhaps small)) but
 essential concs.

In line w. -1.32-.37: In English: "Eval of a string, is invariant if
 we (substitute/replace) part of that string by its "eval". We need to conc "Substitution".
 Maybe make it primitive (.15-.16).

30

One way to do -1.29 (noted above, eval.) Make a set of "primitive" funct's,
 that are adequate to do it. Try to make the "primitives" factorable into
 more "basic" funct's that are likely to be used elsewhere.

34

In general, this can guide a TSPQ design pm! For each "state pt."
 in the TSPQ, we write a soln., then try to factor the soln. into a
 smaller set of "primitives".

36

This idea plus the Learning of various "Definitions" by TM, will be an adequate,
 not-to-diff't way to write TSPQ's.
 I may want to learn "Lisp": or derive a lang of eqn't. power.

9/23/02

TSQ
742.00

TSQ's (2nd final Sem. methodology) 1.00

www.fachtv.com

745



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 744.40: Try to do -1.34-36 on ① 743.10-20^{then} ② learning to solve linear eqns
.01 ~~found~~ N.L. eqns. ③ quadratic, cubic eqns.

.02 **→** For Good TSQ: Look at old AI pems like Styler's symbols interpretation pgs.

.03 **█** Express it as a seq. of "state points".

.06 Essentially -1.34-36 is t. way to write tsq's (at least in Mechanics - T. general path
- t. "goals" of f. TSQ over another story!). I have to ~~be~~ descr. each "state pt."
w. complete exactness - enit to be able to express it as a function. If I can't do this,
then I don't understand that part of t. induction. In general, I may be able to write
a soln function for a problem & give it, but that soln may be of excessive CJS!
It has to be noted: If I can't factor it, then I don't understand that
particular induction.

.10 The language I use to descr. my solns, "polarizations": This I can invent. - dense
substantive "primitives" as needed (like "substitution" -1.15, 743.29).

.12 I'd like to spend more time on this now, but I have to get back to t. paper.

TSQ SO: Main ideas: .06-.16: .00-.01 for initial TSQ's & .02-.03 for more
advanced stuff.

.17 When I get these tsq's as a seq. of functions, I can tailor the y
Learning Algn. to fit into this particular mode of TSQ construction.

.20 The present trace on TSQ's starts 742.00: Read this carefully thru 745.09
Reread until it seems clear how write TSQ. In particular, Note 742.31-40 ("all in form GFD").

.23 **→** A good way to write TSQ's (for both TM & M₂), would be to make
a very large ~~seq~~ array of probs & statepts. for t. TSQ. I would then
gradually fill it in. As soon as I had ~~what~~ what looked like an acceptable
path from one pt. to another, I could compute v. CJS. (not know id @).

.30 To start, t. statepts would be in ~~the~~ fuzzy English - then → exact
eqns as t. TSQ Array matured/moved toward completeness.
Keep things in "English" as long as pract. ... it makes "adjustments" easier

JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20:

Since any hour can be expressed in the ~~GP~~ GPD, we want to "optimize" GPD.
Can we do this via a trick corresponding to $GPD_1 \rightarrow GPD_2$?

We have the infinite set of all GPD's (ignoring GPD). They correspond to: Set of all PST's.
Then we have some kind of GPC over GPD's. - ~~basically~~ we want a PF of GPC over all GPD's.
Then we get a PD over which GPD will be "Best" w.r.t. that GPC. (This will be ~~the~~ ^{implemented} in the usual Mt. Carlo method of $GPD_1 \rightarrow GPD_2$).

What is the GPC of a GPD? We'd like it to be for future problems. [This Q may be related to the "Thy's GPC" problem.]

~~the~~ ^{the} ~~discn~~ of 100 is not so "blushy" - while it is probably not practical to ~~try~~ try to optimize ~~the~~ GPD in any very global sense, still, understanding what the best methods are, we'll enable us to evaluate various approximate ways to improve GPD.

We could optimize w.r.t. past corpus of problems. This is kind of the "Action Alg. Evolution" problem!

The GPC for the QA problem automatically looks into the future somewhat, because of the convergence from wandering for apply future Q's! Hrr, this is the common situation in which "Action" is easy if the "Action" is pc eval.!

Would it be useful to compute the GPC of a GPD w.r.t. a particular problem? \rightarrow See +1.19

20

I have some earlier Discuss of the GPD ~~etc~~ (including both [PST] & PLS etc) ~~in~~ 2 weeks!
is perhaps how to optz it. (C. Peter last month perhaps!) or ~~etc~~ 2 weeks!
Drop this for a while!

24

Look at 739.22 ~~ff~~ - 739.40; 748.00 \rightarrow ~~741.40~~ \approx 741.40
also 738.31 ff

25

I looked at 24-25 rats! Q: Just how do I do "ANL" in that framework?
I need to have a custom approach. So I try all ~~gen~~ gens

28

29

30

the IND problem is Q2: I need to have a custom approach. So I try all gens
w. Best express as input \rightarrow the desired optz. \rightarrow output.
So the input is $[Q_1, A_1]$, we want, as output a O_j that is vg. - also as cost we want successive O_j 's of hyper GPC.

Because of the unique structure of INDuction, we can ~~search~~ search ~~any~~ any prob get a reasonable attempted soln. by trying O_j 's in Least order. - This is easier than looking for good functions to wrap from [QA] to O_j . (It does not, hrr, have to pass of ever being optimal. \rightarrow This I may want to "look into" this! - 739.22 seems more relevant!

36

Also, using tricks like learn induction and "outcast" (738.00) we can get not bad "induction".
So we can go pretty far w. pure QA problems (is any other induction problems) we use the ~~the~~ $GPD_1 \rightarrow GPD_2$ device.

SP.



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: So Do M_2 first (M_1 may not be possible).

01 So: Need to stack grammar here \Rightarrow ass: $(\text{prob. dom.}, \text{PST}_i, \text{cc}_i, G_k)$

$(\text{PST}_i, \text{prob. dom.}, G_k)$ If you cc is included in prob. dom. — but it would be worth to have cc as separate arguments, since we have ~~var~~, empirically, > 1 value for each problem.

From now then ~~from 01 to stack grammar~~ \Rightarrow GPD.

While M_2 may indeed, "whither way" \Rightarrow M_1 become to essentially simple ~~prob. dom.~~ PST,

It could well be that the "system as a whole" will remain factored into 4. $GPD_{1,2}$ Likely
Part 2 [PST_i] part. (Likely) $\left\{ \begin{array}{l} \text{9/27/02: See 761.0010: At (near) end, TM will} \\ \text{only have GPD}_{1,2} \text{ that includes all (probably) PST's} \end{array} \right.$

One stochastic final stack/gram on the [PST] is a GPD_2 . (for both M_1 & M_2 — (E think))

But ~~from~~ \exists also to assoc GPD, grammar. ~~AAAAA~~ In M_2 case, there

is a ~~separate~~ GPD_1 / grammar for each (prob. dom., cc i)

One apparent ditty w. much of the ~~prog~~: There ~~is~~ an ∞ of PST's

We are (theoretically) interested in all of them, w.r.t to present problem dom.

We can cut out almost all of the PST's because their small p's give them to small PC of ~~being~~ over selected \Rightarrow GPD_2 . This could be done in a simple MC or so

way! We have an initial ~~strip~~ on the PST's (\equiv "ready to go"). We use this

to (MC or so) choose one, then we evaluate $\&$ E.P.C. of it.

Consider ∞ of poss. PST's: It is likely that there is one (of very complex ~~for this corpus~~) that ~~can~~ solve all of the prob. dom. for, rather quickly (A.H. ~~hvr~~ — which is why its dom. is so complex!) This sounds ~~like~~ \equiv i. SM's strategy ~~is~~

Looks very Sanders!

evaln" problem — exactly!

"Cross Validation" even if it would usually be "not bad" would not work here

One could always find an A.H. strat. that would survive to "Cross Validation".

Def

Suppose I had a soln. to AAE (Action/Alm Evaln) problem: (This for 2 corpus, a

certain elem. has a certain yield. I can estimate how good it will be in future (unseen)

corpus. [kind of vague!] \Rightarrow would this help in the problem of .25 - .30? ~~at G.V.S.C.C.~~

Another ~~way~~ way of looking at the problem: that of d.f.'s on the A of strategies, (AA's & PST's)

is correct, but in selecting the "best" ones, one runs into the SOP problem! So AAE is SOP

Maybe very tightly coupled problems!

9/29/02
I.D.

!! Action Algum Problem. Soln !! (39-750, 14 and following) SOY: Doth of Problems: 26-38 749



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 749.49: This is a disturbing development! I've been worrying about AAE & SOY problem anyway.
Little progress in AAE; some in SOY.

How Critical is this problem? Can I make some kind of "Buffard Derby" Soln?

One poss. way: ~~that's the~~ Somehow, multiply the prob of being max, by a prop of the PST. If the prop is rather broad & flat, the situation is bad. — But

under that condition, one can do induction — so it is a bad assumption to make!

With fixed sample, as size of data grows, it would seem that the d.f.s of the results would tend to be correct.

In the SM problem of Estimating future yield of technology, I considered the maximum amt. of "edge" in a lot: This is its info content. A shot w. n bits of data could at max get a $\frac{1}{2}$ spurious yield of $x \cdot 2^k$. I felt this was overly conservative for the SM. — But now if we consider all the bit strategies, it may be that this is not so conservative. If we selected e. best bit strategy, it ~~should~~ should be able to get a yield of $x \cdot 2^k$, using Kelly's. (I'm not so sure of this, here.)

Anyway, I think SM is easier than the general AAE problem, & I haven't been

able to solve the ~~SM~~ SM what a val. problem! I did have that coding method, but sound o.r. intuitively; yet w. careful analysis, sound not to work! [See SM [1.19-20] (7/27/02) [1.25-29]]

19: SOY ^{Spurious Optimiz. Yield} Problem! Say we have a bunch of AA's! We have a joint PD for the yield n-vector of the bunch. We have a way of finding the PD of the peak G. (Obviously M.C.R.V.s, but also expressible as a mult. dim. S.). We also have a way to get the pc of each AA's ~~being~~ having best G (Also available via M.C.R.V.s or by S.). We can then get the G d.f. for that particular AA; Easy if they are "independent" Marginalize" by integrating over all other AA's (if any).

25: SOY Problem! In "true" SOY problem, we have a joint G d.f. of n AA's. $P(x)$; we have to appoint d.f. $P'(x)$ (due to sampling error). We select a "most likely" best "AA" on basis of $P'(x)$. It has a d.f. of G which is in error due to $n \rightarrow \infty$. How bad is this apparent d.f.?

26: Well, if n is very large, there will be some AA's that look much better than they really are (this is SOY effect). Also, so v.g. AA's that look much better than they really are; but there are fewer truly v.g. AA's, so this is not as important an effect. As $n \rightarrow \infty$ (as in T.M. case) this disparity can be very large! There is probably a tendency to pick one of these "spuriously good" AA's, & the effect seems to get worse as $n \rightarrow \infty$ (particularly as to AA's are a "universal set" — so they can have any degree of "fitness".)

38: So this does sound like the T.M. problem.

39: SOY Problem! Can we (aptly) look upon this problem as a Bayesian estimation in which I had (incorrectly) assumed a uniform prior for each of the AA's being best? Yes! — see + 1.13 for reason



SO: 799.40: This looks like a correct view! What's not clear is just what a prop. is "of":
Is that -1.40 , the spread of each AA being "the best", or should I use some other spread?

At any rate, the Big impl. idea is that the prop. is very unlikely to be "Uniform" over the AA space: More likely uniform over the dens of the AA space.

However, Another Q is: what do I want the prop. to be "of"? - e.g. I can have an spread of the yield D.F. of each AA.

→ I want an "spread" that can be modified by a "Summarizing Machine"

So it is "impossible" w. ~~training~~ "Training".

Apart from the details, this looks like to break thru of the last few Decades!

Seems to be the wrong problem $M \leq H$. (I'm not sure as to whether SOY $(-1.26 - 38)$ solved).

→ T. main idea here, is that whenever you compute a spread, one must do it in a Bayesian way & one must know in advance.

There is still a lot of Q about what spread to use: any spread of G.P.D. or of G.P.D.?

Also! What form should it take - what it is dependent on - what is a "Summarizing Machine" of.

Lets Go back to an (apparently) simpler problem: Stat. eval. in SM!

For each strat, we want a P.R. on G. - This is a mean yield (4 yield; also a ≤ 2 ;

What we want (in ln down) $\geq \ln \leq 2$ for a reasonable ΔT - and so we can get ΔT for longer ΔT 's by assuming \approx "independence". Lets gloss over the details here, say there is

some ΔT that we want to D.F. of G (\approx yield) for.

One guess: for all Strategies spread: $\mu \geq 0$ in ln down: σ^2 is a function of coding cost & strat. Just exactly how σ^2 is a function of code p is important; but not obvious. Maybe $\sigma^2 = \ln 2^{code length}$.

$= -\ln(\text{prop. of strat})$. This σ^2 code length is \approx the amount of mult. yield from a pure for the strategy. Its "Available Fudge". This gives large σ^2 for unlike codes: This is opposite of

What we want! Perhaps $\sigma^2 = (-\ln(\text{prop. of strat}))^2$
 σ^2 is ≤ 2 of total (normally), so σ^2 is $\frac{1}{code length}$ that means $\frac{1}{code length} \approx \sigma^2$

A reasonable approach: Say X is code length of strategy, let $\sigma^2 = f(X)$: $f(\cdot)$ to be determined.

Investigate what happens for various forms of $f(\cdot)$

$f = \text{constant}$ is what could be \approx spread errors. Another wraps! $f(x)$ is not f . This f has

to converge - σ^2 only σ^2 of the spread.

- But maybe I'm barking up wrong tree! For Bayes, Spread of hypothesis = 1 - first coverage.

36 [I maybe thinking about the wrong "Approx".] A simple case with the spread of G.P.D. to spread that
37 [given PST will be the best PST for a given problem.



411

00: 7:50:40: Consider "Uniform" approx for GPD₂. There are an ∞ of PST's to consider & ~~compare~~ one must (if in Mt Carlo mode) try all of them w. equal frequency to get the distribution of "good". It may be just there are an ∞ of PST's that get the same by G values for Uniform ~~approx~~, approx, the Mt. Carlo method would not be practicable.

05: If two are non-uniform approx, what is the distribution of "max" ρ ?
06: A poss. Mt. Carlo data ρ (or Mt. Carlo choice PST (bases of approx PST's). w. ρ rank
06: Mt. Carlo choice ρ value of ρ from that PST. We get a distribution of points: each has ρ (or PST_i) label.
06: What is prob that "PST" label is "best"?

09: An interpretation: Do ρ MTC choices like .05 = .06: ρ count how many times PST_i has been Best, = C_i .
10: Do this many times. Then ρ to GPD₂ of PST_i is $\rho \rightarrow C_i / \sum C_i$. Does this vary much w. k ?
11: Does it \rightarrow a limit as $k \rightarrow \infty$? What is value of ρ for $k=1$? (Monte Carlo) It's simply a prod of
11: PST_i. What is value for large (but finite) k ?

13: Well, say G has a ~~sharp~~ sharper limit G_{max} G has discrete values only. For large k , the population of
14: ~~set of trials~~ having G_{max} will be diff. for any k .

14: There is a "large cut k_0 " so that it's very likely that G_{max} has population of at least 1. for any $k > k_0$
14: we will get a distribution in G_{max} & t. no. of cases each TPT has prod of G_{max} . So
14: it will be indep. of k of t. d.f.'s!

19: So .13-.14 is probly riter! - i.e. idea of .09-.11 wouldn't work.
20: An eq. that I can't exactly justify, but which may have reasonable behavior!

21: Wt. of PST_i:
$$\int_{-\infty}^{+\infty} dw \left(P_i(w) \right)^{a_i} \cdot \left(\int_{-\infty}^{+\infty} dw \prod_{j \neq i} \left(P_j(w) \right)^{a_j} \right)$$

21: a_i is t. a prop of PST_i: $\sum_k a_k = 1$
21: $\int_{-\infty}^{+\infty} dw P_k(w) = 1$ for all k .

26: The former case of: finiteness of PST's was dealt w. by making all $a_i = 1$.
21: does not give same result, if all $a_i = \frac{1}{k}$ (for k PST's)
21: If in .21, I replaced a_i by 1 & " a_j " by $\frac{a_j}{a_i}$. I'd get consistency. w. 26

30: .21 \rightarrow wrong way: should be for .26
$$\int_{-\infty}^{+\infty} dw \left(P_i(w) \left(\prod_{j \neq i} \left\{ \int_{-\infty}^{+\infty} dw P_j(w) \right\}^{\frac{a_j}{a_i}} \right) \right)$$

31: To modify .30:
$$\int_{-\infty}^{+\infty} dw f_i(w) \left(\prod_{j \neq i} \left\{ \int_{-\infty}^{+\infty} dw P_j(w) \right\}^{\frac{a_j}{a_i}} \right) \rightarrow$$
 also 758.15
see 752.31!

Note, for small $\frac{a_j}{a_i}$; $X^\epsilon = \sum_{l \geq 0} \binom{\epsilon}{l} X^l$. $X < 1$ so $\ln X < 0$.
31: $X_j = \int_{-\infty}^{+\infty} P_j(w)$; $\prod_{j \neq i} X_j^{\frac{a_j}{a_i}} \approx \prod_{j \neq i} \left(1 + \frac{a_j}{a_i} \ln X_j \right) \approx 1 + \frac{1}{a_i} \sum_{j \neq i} a_j \ln X_j + \frac{a_j}{a_i} \ln X_j$ see all 756



749.18.19

Tom English
Jackson State U...
Jacob Hiss...

0: 751.40 : In SM if one has a strategy, there is always one "theoretical" way to evaluate it.

Take to past corpus using AHP, make a D.F. on all future corpus of the corpus.

Test the strategy of interest, on this future D.F. of the corpus of the corpus.

T. complexity of the strat. is irrelevant under these conditions.

Could we do anything like this in the AA problem? AA seems quite different!

Consider PST evaln. We want to guess how good PST₂ will be for problem, X, w.

CB = CCo. If, as in .00-03 we had evltn (actually CB=20!), we could

actually determine G as a funct of PST₂; X & CCo - exactly. We could do this for

any PST, independent of domain.

These problems: AAE, SM strateln, PST peaking, SOY, are related but not exactly

the same: I've worked on them a lot, but ~~can't~~ possibly for SOY I have no

usable (or any) solns.

I should: state each problem clearly (for SOY 749.26-38; see)

Keep a record of state of soln. : Look at "How to Solve Problems" list.

List some special cases of known solns.

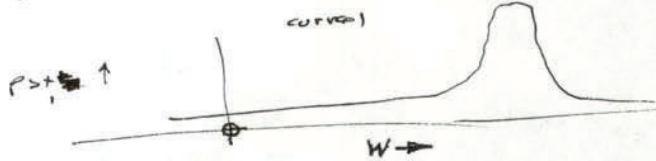
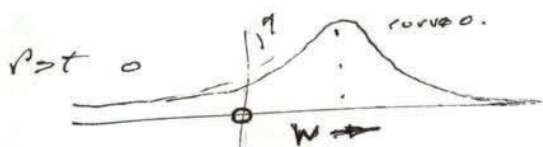
List some non-solns -- : why they don't work.

SN A wild idea! for each PST: we have a P.D. on the G.H.S. CC.

How, if the den is long, the P.D. is less likely to be accurate.

So in comparing PST₀ w/ PST₁; if PST₁ has ~~much~~ long den, then

When PST₁ seems better than PST₀ for many W values,



the prob that PST₀ is worse than PST₁ is very high if the curves were wide!

But we have very little confidence in curve 1

This miter justify 751.31 $\frac{\partial j}{\partial r}$ over our relative confidence in PST₁ v.s. PST₂.

a small $\frac{\partial j}{\partial r}$ means $\left\{ \int_{-\infty}^{\infty} dw f_j(w) \right\}$ is close to 1.

One complaint about 751.31: If one computes the prob of each PST being Max, using this formula: the pc's will not sum to 1. We could normalize, but this messes the point.

0/1/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 752.40: An easy way to test conjectured models: Try them on simple cases w. $n \geq 3$ or 4 PSTs. Try on SM; then compare w. SM in which analysis doing pure prediction (so one knows what to write Soln. is).

01: A ~~poss.~~ analogy betw. probabilistic & non-probabilistic problems: Say we are trying to predict stock A. It "could" correlate w. any of 9k different stocks. Its likely that one would (A.H. ishly). To use this correlation, we have to specify the stock, which amounts to $w \log_2 9000$ bits.

10 In the case of a very large (actually ∞) PSTs, we have analogous problem, but the thing being predicted is different. Similarly w. "Strategies" in SM: Reason even ∞ of them.

11 Another promising approach is 749.18-19: using yield of a strat. as a method of coding the corpus. It seems clear, that if a strat. has yield "significantly $\neq 0$ " that restricts the poss. corpus to in fact a large coding method. If P_i is the a priori on the i th poss. corpus element; q_i the strategy modifies this distribn. to P_i' ; q_i is the trans pd. into corpus; then the coding gain $\leq q_i \ln \frac{P_i}{P_i'}$ bits.

17 Sometimes $q_i = P_i$ or P_i' or nothing. 17 is the difference in KL distance to q_i of P_i' vs P_i .

20 Consider a random binary seq. to be the daily price (± 1). A strat. with yield Y . Is Y constant on the corpus. The larger Y is, the more of a constraint (i.e. fewer poss. codes w. that yield).

26 One kind of strat. Just tells whether to buy or sell — so one has a partitioning of the postfixes of the semi-infinite binary string corpus. W. certain postfix one buys +, other postfix "-"; other postfix, nothing. 3 partitions. So any partitioning of postfixes into these 3 categories is a strat. At all times, one will be int. market: either long or short \$.

30 — we could make this exponential yield if we ~~know~~ know probability or loss — but just consider this simple (non-exponential) betting method.

32 [On the other hand, if one averages a yield of Y_0 per bet, one should bet a certain fraction (function of Y_0) of one's fortune on each bet for "Kelly betting".] — So perhaps exponential betting is normal. The value of Y_0 can be determined by Kelly's formula or Kelly's formula.

34 For a binary seq. of 20, it doesn't make any difference whether we use exponential betting or not. For a given Y , on a corpus of length n ,

0/1/02

754



0: 753.40: There will be a certain subset of strings that have yield = γ . If k out of 2^n strings have yield γ , then the strategy has info content $\frac{k}{2^n}$ or $n - \ln k$ bits.

If $k=1$, then clearly the strategy specifies 1 string or n bits info.

Also note $n - \ln 1 = n$, as it should. So the ratio of the no. of possible sequences before

after the strategy specifies γ is $\frac{1}{k}$ of the no. of bits info.

How this is not a good way to code (normally). For a v.g. strat, we expect

a doubling of fortune/yr. This is 1 bit/yr. If the strat. takes 2 bits to describe, it takes 2 yrs of data to verify that it has positive yield.

This means Very simple strats. Or superhyields!

Passing in .06 ff: if a strat has a code worth of r bits, ($2^{-r} = \frac{k}{2^n}$)

Then info of γ , our real yield is 2^r (i.e. if we did suitable Kelly betting, perhaps

using x wins or r units (as $-1.32 - .34$) to get pc's.

Actually, Pay is reasonable!: [If, every d days, we make a bet w. a certain average expected yield, this is equiv. to reducing the future courses of the corpus by a certain factor.]

To estimate k , $(.00)^d$ the redundancy of the strat; Run

random seqs thru it & see how frequently the yield = γ . The SSZ

for yield of just γ may be small, so if k is a \downarrow func of γ , we can collect data on no. of "hits" of yields near γ , & interpolate

(\equiv smoothing of data).

The large will work also, if we have a Brownian Motion in the market. We simulate it like 19 ff by a random walk with the appropriate step size and ΔT of the market. We apply the strategy to it & see how often its yield is within a factor of 4, say, of γ .

is no interpolation in the n domain. In this case, we want the density of cases per $d\gamma$ or $\frac{d \text{ no cases}}{d\gamma}$. This has to do with

precision we need to specify γ . (I think I've been thru this before, but $(.01 - .05)$ is $(.18 - .23)$ are for a simpler SM model

(i.e. binary seq) & the system does seem to work!

Around 8:30 time, I should look at the simplest argument v.s. the coding method of strategy as a code

} -> Note 755.00

8/2/02

755



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20 (754.39) In Ray Brownian Corps: σ var. "per unit time" is constant. To deal

W. ~~large~~ fine varying variance a bit: Devise strategies so that threshold for next bet, is always \propto threshold for previous bet. woops! This gives constant σ^2 being assumed.

A way that might work: bets depend only on $\Delta \ln$ price history. This helps some, but still, σ of σ^2 per hr, could vary a lot over a week, and so our ~~market~~ Brownian Mkt on \$754.24 would not be v.g. (.18)

.08
.09
10

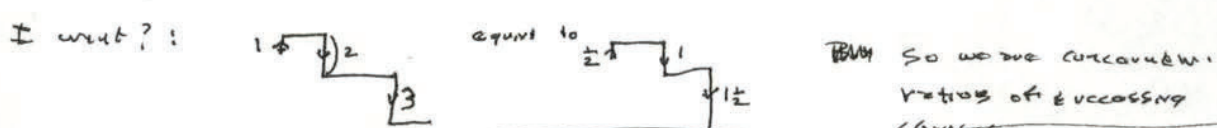
A way to study the variation of σ^2 in the SM signal: plot σ total variation (using some grid size: see what happens if grid size ~~is~~ is changed? double, x4, x8...). Then take a smoothed value of this as "local variance".

Q: would this respond (more rapidly/"better") than smoothed RMS?

To get RMS, we need some pre-smoothing to get σ mean (μ) - so we can get RMS deriving from it. The total variation in Redoff .09 seems to have less "smoothing" in it. A min grid size could be the grid ~~size~~ size used by the Mkt quote process.

.18 (08)
20

A poss way to deal w. "drifting" σ^2 : T. strategies categorize σ post fix (history) in 3 categories (with bins of 753.26-27). These categories depend on ln of changes in σ (ln price). Unfortunately $\ln \phi = 20$ and $\ln(-\sigma) = \text{imaginary}$.



A Q is: How ~~well~~ ^{well} does want to try to simulate SM signal - i.e.

How good should the "default" Diff be? The discussion of 753.20 off on binary strings, is the Basic Argument - first has to be ~~applied~~ applied, somehow, to the more general case.

30
31
33

On one side Say we have a fixed strategy: If we use a default diff w. very small σ^2 , the prob of obtaining yield Y w. the strategy, will be very ^(large) _(small). So a σ^2 off. default def. is very important! But in the real SM, the σ^2 is ~~is~~ always varying. partly not true

How. Some stocks may be invariant under change of σ^2 - re. yield indep. - in fact over of changes that I could into was σ^2 indep. We look at O p.m. first $\frac{1}{2}$ hr. of Mkt. Then this scales our better local behavior for rest of day on market.

On another note: Even using the simple model of 753.20 off, the thing that I'm really interested in (corresponding to the AA problem), is the diff. of the yield, Y of the strategy over future corp. What 753.20 says is: "I've found a short code

0/3/02
59

256



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 755.40 For t . corpus - its v. bits shorter than a usual deriv. - so t corpus yield 2^r if I do Kelly = 1. " I use the strategy to figure out how to ~~do~~ base do Kelly bits - This is normally not the same as using t - strate to boot.

Another interpretation: Since this short code for t . corpus exists, it means that if ALP had been applied to that corpus, it would have had at least

yield of 2^r . Also: It means that "expected value", "an unbiased estimate" of yield for a corpus of length n is $\geq 2^r$

we can look at the problem of predicting t . corpus: we have found a predictor of wt. $2^r =$ recip of predictor \times (PC of corpus w/ t . predictor) -

we have at least 1 predictor w. wt. 1: ~~It assigns~~ It assigns $p = .5$ to 0 & to 1

Consider t . statement. PST_2 had a ~~matrix~~ G of G_{ij} for

Prb j . " Her, we are interested in predicting G for Prb j by PST_2 (G_{2k}) It would seem that this predn. - (i. data of G_{2k}) would be indep of t . deriv length that PST_2 . But perhaps t . D.F. is rather \neq (at), because we did not spend much cc on it (if PST_2 's deriv is very large).

18 Consider t . s-grammar: $P(PST_2, Prb_j, G_{2j}, G_{2j})$: G_{2j} as a Probabilistic funct of PST_2, Prb_j

19 Prb j ~~is~~ G_{2j} are fixed. so $P_j(PST_2, G_{2j})$.

Is it true that if PST_2 has a long deriv, then P_j will ~~be~~ tend to be \neq ?

Well, for rewrite ~~type~~ type Grammars (like CFG's). long deriv, mean lots of substitutions, so low P_j - Tho this is "off to mark"

In .14 we have a continuous def. on G_2 not. discussion of t . structure of PST_2 seems irrelevant.

we have ~~in~~ in (.18) a $P(\dots)$ that were optimized via (4.18).

If PST_2 is complex (\geq has long deriv) P_{22} , will .18 take a long time to evaluate?

30 Back to ~~0.18~~ 0.18 Consider a "cut off" ϵ so only PST 's are included if their ϵ bits are $> \epsilon$. Clearly, for small ϵ ϵ t . system will not work (its too close to $\epsilon = \phi(!)$). So if we had some lower bound for ϵ .

Good! It would seem that in suitable circumstances, one would ^{have to} consider ~~the~~ ϵ of any degree of smallness.

So, w. finite $SS \geq$ - just what is t . limitation on how small ϵ can be?

[Kelly's sharp ϵ limit is like "Rwindow" = I'd like a "soft" limit on ϵ , more in t . spirit of ~~X~~ window]



00 : Hum! In comparing 2 PST's for likelihood of "best": T. ^{P_i} ~~approx~~ can be used in

follow. way: P_i(G₀) is "better" than P_j(G) if ~~the~~ PST_i exists and

-03 PST_j does not exist or P_i(G) > P_j(G). ~~But~~ T. probly that $\left(\begin{matrix} \text{PST}_i \\ \text{PST}_j \end{matrix} \right)$ exists is $\left(\begin{matrix} P_i \\ P_j \end{matrix} \right)$

→ One test of .00 - .03 is, to opt for 751.09 - .19

.05 Another poss. way to do "Max expected G": Pick pairs of PST's w. ~~the~~
probys P_i⁰, P_j⁰ (≡ approx of PST_i, PST_j resp.). Check a G value to reach
in accord w. its P.D. for G. Compare to 2, & give a score of 0 or 1

10 ~~for~~ for loser, winner. ~~For~~

Approach (.05 - 10) may be ~ or ≡ to (.00 - .03). In both, we do pair-wise comparisons.

A poss. approach to "Pc that PST_i is better than any other PST":

14
$$P_i^0 \cdot \int_{-\infty}^{+\infty} \left(P_i(G) \cdot \prod_{j \neq i} \frac{\int_{-\infty}^G P_j(G) dG}{\alpha} \right) dG$$

Replace α by $\left(1 - \int_{-\infty}^{\infty} P_j^0 \cdot P_j(G) dG \right) = \left(1 - \left(1 - \int_{-\infty}^G P_j(G) dG \right) \cdot P_j^0 \right)$

20 For small P_j⁰ $\approx \exp\left(-\int_{-\infty}^G P_j^0 \cdot P_j(G) dG\right) = 1 - P_j^0 + \int_{-\infty}^G P_j(G) dG$

Then .14 becomes
$$P_i^0 \int_{-\infty}^{+\infty} P_j(G) \exp\left(-\int_{-\infty}^G P_j^0 \cdot P_j(G) dG\right) dG$$

This expression is Min of 751.20 - .40 and 752.20 - .40.

26 A rationale of .14. P_j⁰ is the probly that PST_j will occur (≡ "exist").

$\int_{G_0}^{\infty} P_j(G) dG$ is the probly that if PST_j occurs, it will have $G > G_0$.

.29 $P_j^0 \int_{G_0}^{\infty} P_j(G) dG$ " " " PST_j will occur and " " " ~~eventually~~ PST_j will not have $G > G_0$.

30 $\prod_{j \neq i}$ is probly that none of i PST_j's (j ≠ i) will have $G > G_0$.

33
$$\int_{-\infty}^{+\infty} P_i^0 \cdot P_i(G) \cdot \prod_{j \neq i} \left(1 - P_j^0 \int_{G_0}^{\infty} P_j(G) dG \right) dG$$
 is (sort of) probly that
PST_i will exist & be better than all other PST_j's (j ≠ i)

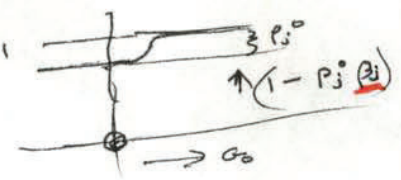
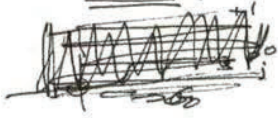
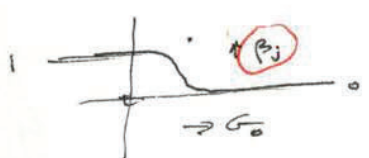
See 758.20
- 759.12
for whitening
Criticizing!
(Looks pretty bad!)



01/14/02 JD
20: 757.40 : -1.35 is "rational" (-1.26 ft) & is a reasonable eqv. to investigate.

Try a test \approx or \sim to 757.09 - .19.

Is there a Mt Carlo way to realize 757.35?



~~It does not take 1.35 may be barriers in a critical place~~

$1 - P_i^0 P_j$ is the probability that PST_j will ~~not~~ exist and have $G > G_0$.

10 $\rightarrow \prod_{i=1}^n (1 - P_i^0 P_j)$ is the probab that none of the PST's will be present at $G > G_0$.

for small P_i^0 's it is $\approx \exp \sum -P_i^0 P_j$ since $P_j \leq 1$ & $\sum P_i^0 \leq 1$ the product $\prod (1 - P_i^0 P_j)$ must converge to > 0 .

$1 - P_i^0 P_j \approx (1 - P_j)^{P_i^0}$ only if P_j is small. $(1 - P_j)^{P_i^0}$ is better, since P_i^0 is usually small.

15 $757.30 \stackrel{is}{=} \int_{-\infty}^{+\infty} P_i(\omega) \prod_{j \neq i} (1 - P_j^0)^{P_i^0} d\omega$

[T. idea of using P_j/P_i^0 was that we considered a further smaller fraction of PST's, w. smaller sum of $P_i^0 P_j$, than it]

If we computed \int

20
21 $F(G) \equiv \prod_{j \leftarrow \text{including } j=1} (1 - P_j^0 \int_{G_0}^{\infty} P_j(\omega) d\omega)$

would give the result
first we find $1 - P_i^0 \int_{G_0}^{\infty} P_i(\omega) d\omega = \delta_j(G)$ for each j
Then to get -1.35

34 $F(G) = \prod_j \delta_j(G)$ Then $-1.35 = \int_{-\infty}^{+\infty} P_i^0 \frac{P_i(G_0)}{\delta_j(G_0)} F(G_0) dG_0$

In fact, to obtain -1.35, we would ~~not include many~~ include only a relatively small set of PST's in the calcus, since most of them would have such small P_i^0 's that they would modify the result very little.

30 $F(G)$ is really not very small. $\approx f(G) = \prod_j \delta_j = \prod (1 - P_i^0 \int_{G_0}^{\infty} P_j(\omega) d\omega)$

Not exactly! $1 - \epsilon$ is $\leftarrow e^{-\epsilon}$ $\sum \epsilon_j < \sum P_i^0$ so $\sum \epsilon_j < 1$ so $\prod (1 - \epsilon_j) > e^{-1}$ and F could be > 1 .

30 -1.35 from .34 $\approx \int_{-\infty}^{+\infty} P_i^0 \frac{P_i(G_0)}{\delta_j(G_0)} dG_0$ (it could be hyper factorate at most)

36 $-1.35 \approx \int_{-\infty}^{+\infty} \left(P_i^0 \frac{P_i(G_0)}{1 - P_i^0 \int_{G_0}^{\infty} P_i(\omega) d\omega} \right) dG_0 = \int_{-\infty}^{+\infty} \frac{P_i(G_0)}{\frac{1}{P_i^0} - \int_{G_0}^{\infty} P_i(\omega) d\omega} dG_0$

If P_i^0 is very small then $\approx \int_{-\infty}^{+\infty} \frac{P_i(G_0)}{\frac{1}{P_i^0}} dG_0 = P_i^0$

If P_i^0 is ≤ 1 , say $P_i^0 \approx 0.5$, is not bad!



00: 758.40: So it begins to look like 757.35 is usually $\approx P_2^0$ (A terrible result)

But 758.36 R is between P_1^0 and $\frac{P_1^0}{1-P_2^0}$. (Probably).

The worst part of this arg. is that $F(G) (-1.21)$ is between $1 \pm \frac{1}{e}$.

$$-1.35 = \int_{-\infty}^{+\infty} \left(P_1^0 \frac{P_1(G) \cdot F(G)}{1 - P_1^0 \int_{-\infty}^{+\infty} P_1(G) dG} \right) dG$$

Since $F(G)$ can vary in size by a factor of $w \in \mathbb{R}$, this expression can vary (at most) by a factor of e from P_1^0 or $\frac{P_1^0}{1-P_2^0}$.

It might be poss. to find out what's wrong w. 757.35, by considering the case of a small no. n , of PST_j 's of = wt. In this case I (think) know the right soln. Take a look at (751.3) again.

$$757.31 \rightarrow \int_{-\infty}^{+\infty} \left(\frac{P_1(G)}{\int_{-\infty}^{+\infty} P_1(G) dG} \right) \left[\frac{\prod_{j=1}^n \left(\int_{-\infty}^{+\infty} P_j(G) dG \right)^{P_j^0}}{\binom{n}{j=1}^{(j=1)}} \right]^{\frac{1}{n}} \left. \right]^{\frac{1}{P_1^0}} dG$$

T. Geom Mean

T. Strongest case for this guess is that it satisfies the constant of .10-.12. (.20-.30 Below is a

Another perhaps more intuitive expression. Cudly base same as earlier one) slightly stronger argument possible

- 1) Pick a random G value from PST_1 .
- 2) " " " " " " a random PST_j (using appr to choss PST_j).

What is pc that $1.13 > 2$?

To criticize $z = .20$ consider the standard of n PST_j 's of $P_j^0 = \frac{1}{n}$ for all. $\&$ $.20$ doesn't seem like a result: It gets to be that PST_1 is better than a "typical" PST_j - not to pc that its better than all PST_j 's.

$.13 (=751.31)$ does seem to be closer to 1 : It does take the product of

all of other pc 's being \approx worse than PST_1 .

Consider the "standard" of $.26$: we compare PST_1 w. the $n-1$ power of the geom mean of the other (PST_j 's) being worse than PST_1 .

Tho. $.13$ is "not unreasonable," it really doesn't have any definite reasoning behind it!

It would probably be poss. to express $.13$ as a Monte Carlo approx., using some other function than $\int_{-\infty}^{+\infty} P_j(G)$ as the comparison function. Doing so might also give some understanding of behavior of $.13$.

See 769.33 for poss. applic. of (13)

0/5/02

REV: 17

Development of "GPD-PST Grammar" - from Early Beginnings

761



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 760.40: [SN] Re to GPD-PST grammar: In early TM, we will have a small no. of PSTs inserted by user. Later, we will insert more, a conventionally def. GPD-PST Grammar. I think these steps can be done "smoothly", w/o basic Module of Introduction.

may be only later

"More" means a finite no. - they will be evaluated as individuals, indep of the details of their deems. Hvr, each PST can have a special way to evaluate its assas, ~~RESEARCH~~ Problems, CS, G d f. - ~~each~~ each PST can have its own stack (Prob, CG, G) Grammar. - Re to Grammars can have

"shared" concs,

Next, we will begin factoring PST deems, so as to ↑ % of overall corpus deems. This factoring may be diff, so trainer may not help TM by hints, or actually do part or all of factoring.

16 . . . Also, ^(while when) this is occurring, trainer may add many more PST's.

17 Review of 746.00 - 761.16 (PST, Prob, CS, G)

746.00 - 747.40: Discn. of whether to use System where GPD is on

or just (PST, G) 747.37 - .40 is reasonable resolu.

748.17 starts diff of "200 PST's" and to Ad Hocness of ^{very} complex PST's
Idea that this is a complex of AArvaln. problem, SArvaln. problem, SOY problem - which I haven't adequately solved. [Later, it appears that the problem seems different]

5 may have "simple" soln! 760.16 - .18

749.26 - .38 V.G. Deem of SOY problem

749.39 is idea that in evaluating PST's, one must consider their April: Just how this April was to be used, was at first unclear!

Some Attempts! 751.05 - .06 A Mt Carlo way; inter proba! (bid .13 - .19 show it wouldnt work)

751.20 was eventually decided on as a "good try" - 759.13

753.11 - 756.10: A good treatment of SM problem of using a strategy & its approach yield, to Code a SM corpus. It does seem to work 756.06 explains what it actually means

756.18 - on connectivity of GPD{PST} grammar. - Deem of problem: (Later 760.16 - .18)

Seems to be correct soln. 756.30 - 757.25 ← trial soln. to (756.18) 758.00 - 759.12: Demo that 75335 will not work

759.13: report of 756.20: which looks "not bad".

762.00



00: 761.40 ^{50%} 759.20-.40 tries a new model, but gets more reason to believe 759.13 is "not bad", which ends up N.B.
759.35 - way to analyze 759.13 ; 760.00-10 other ways to cost 759.13

01: ⑤ 760.13 : Suggestion that I review AAE, SM strategy, PST evaln, PST Grammar Construction, SOY
⑥ Empt., apparently Related, ideas. ^{PSF evaluation}

④ 760.16-.18 Idea that Goro for (PST, Prob, CC, S₂) grammar, is just Goro for standard QA problem : Max = 414.18
followed by design of how to "relate the" soln to other solved & unsolved problems of 760.13

Discussion: t. main, most critical ideas reviewed.

- 1) Dorn of SOY problem 744.26-32
- 2) "soln" of evaln. of SM strategy via coding 753.11-756.10
- 3) A poss. soln. to the PST amptzn problem: 751.20, = 759.13 This is a more exact form of GPD, \rightarrow GPD₂.
- 4) Final soln. to GPD; [PST] Grammar construction problem - Part ~~that~~ was same as Goro of 414.18. (760.16-.18)
- 5) suggested review of relations of AAE, SOY, SM Strat., PST evaln, PST Grammar Goro, STEIN effect

0/11/02 (Later) ① I'm not completely convinced that 760.16-18 is a soln. (or "t. final soln"). One very serious "objection": It does not seem to answer to Q of 746.00 i.e. given \$ 2.5 Gramms for PST's with \approx 414.18 Goro.

One could be "much better" in merit has PST's (or just one PST) that gives (casually) much "better" ^{less cc for some G} solns to problem.

God \rightarrow A poss. Approach: That t. ~~the~~ 5-Gramm simply presents the facts. Finding a "good" Grammar is a separate problem: Yes. finding t. "Best" PST under these circumstances can be quite diff, and ~~perhaps~~ may need a "creative" soln. - see 769.20 for more discussion.

But t. "Best" GPD, is still t. one w. t. Best 414.18 Goro.

While it is in theory, poss. to evaluate all of the PST's, Prob, CC, & obtain exact G, we are limited in time available & must use approximate PC values.
(It's like to pd. for the 10¹⁰⁰ bit of IT: Given a finite CB, t. pd. will depend critically on t. CB.)



20: 262.40: As the system seems to work now: It can function as a QA system only.
As I'm thinking about it, the update algm. for O^1 , would be relatively simple.

It would use ~~previous~~ functions & subfunctions of previous solns. Also, it would use context of 'linear' & 'super kinds' & it grows in knowledge.

However, the Update system would be static. It would not change.
I don't (yet) see any need to change. All kinds of regys would be accessible to the system (I think!).? (vary variables hrrr!)

10

~~It~~ If I'm lucky (but it may ~~not~~ be optimum) it might be "fairly good" — good enough to do updating of IND's other or problems.
Good enough for trainer to introduce new PST's & eventually have a FPD [PST] grammar. And eventually, good enough to ~~improve~~ improve (propose new PST) PST for Inductive problems.

20

21

Look back at 696.13-40 (outline of paper). There is no discussion of the special IND update! It's not exactly Lurch, but it's almost Lurch, because world looking for short paths. However, the Pd (the cartesian) ~~main~~ is invariant! Also, I probably wouldn't be using much Backtracking, so I should retain several || solns. that are significantly different.

I haven't yet been able to define this idea.

The reason for preceding remark: That the update almost always only add to the existing "soln" as new problems are solved. It is possible for the update to "correct" the soln a lot by observing a new conc. that has ~~large~~ large case count, but unlikely since new conc. (should be) discovered by OSL (One Shot ~~bug~~), error will be at least case count ≥ 2 or 3.

30

One way that a great (revision) can take place! A Large CB search is done, which is able to detect regularities previously unobserved because of such small CB being used (i.e. no specific or or INV facilities).

As a "stand alone" system, the QA system can be improved by (adding/modifying) update techniques! E.g. GA (GP) or RANN. (or regular ANN.) — specific tricks like dividing up corpus (at first, using indices put in by trainer)

0/6/02 ID

ANN and RANN: How to Get ALP values



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20

Say t_i wts for a ANN are betw. $\pm \frac{1}{2}$. So for a N# wt. system, T_i h volume of dens is 1. Then, if t_i h vol. of "acceptable" wts is V_i , its pc is just V_i . We can balance this against the "precision" of t_i output in prediction.

23

If t_i error is σ^2 , for k observations pc per obs. is $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x_i^2}{2\sigma^2}}$

I guess its nice to have some kind of spread for σ .

10

More exactly, t_i pc of all k obs is $\prod_{i=1}^k \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x_i^2}{2\sigma^2}} = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^k e^{-\frac{\sum x_i^2}{2\sigma^2}}$
so $\sum x_i^2 \equiv K\sigma^2$ so $\dots = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^k \cdot e^{-\frac{K}{2}} = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^k$

If σ , a priori we know only that x is

betw. 0 & A , then t_i first obs. costs — No!

Because of the predn. in t-ANN, all we need to descr. x_i data is t_i .

Seq. of x_i 's which are distributed about σ of w. var $\sigma^2 \equiv \sigma^2$.

Then even more exactly, say each x_i is known w. precision $\pm \frac{1}{2}\epsilon$.

23

Then t_i pc of each obs is $\frac{\epsilon}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{x_i^2}{2\sigma^2}}$

so t_i total pc of data is $\left(\frac{\epsilon}{\sqrt{2\pi}\sigma}\right)^k$.

The ϵ^k factor enters in, however how we code the data,

so it is not relevant for predn.

In comparing pts in wt. space $(.00-.03)^k$. Each pt will have an associated σ i.e. a σ^{-N} factor.

Say we have z pts in wt. space distinct local maxima;

20

How to compare them: One way: (z) : Take h volume, V_0 (or V_1)

31

→ in which σ^2 is w. min $\frac{1}{2}$ of t_i minimum. (see 36)

32

Compare $V_0 \cdot (\sigma_0)^{-k}$ with $V_1 \cdot (\sigma_1)^{-k}$ or $V_0^{1/k}/\sigma_0$ with $V_1^{1/k}/\sigma_1$.

or $\ln V_0 + k \ln \sigma_0$ with $\ln V_1 + k \ln \sigma_1$

34

One could do all of this above (31) by integrating over t points in wt. space, rather than "t. h vol. within $\pm \frac{1}{2}$ power pts".

space
265.00

9/6/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:264: ^{Spec} ~~to~~ : The problem at 764.00 seems similar to how many coils to use in regression. (L in or N. L in.)

The trouble is in comparing V_0 & V_1 (-1.31) for t_1 , 2 different nos. of coils.

For t_1 "new" coil(s) we need an \geq prior \leftarrow or a set of values between which the
prior is uniform. In this case we can compare V_0 & V_1 in a different way: Say V_0 is
 N dims $\geq V_1$ is $N+1$ dims: V_1 is a N dim space \wedge (comparable directly w. V_0) mult
by ≥ 1 dim space V_1' . It is V_1' that we need to \geq prior for. If a certain segment
in V_1' space, of length l , marks out " $\frac{1}{2}$ power pts"; then we want the ratio of l
to the a pri "legal" region of V_1' , same

Actually, the analysis of .01 if $(V_1, \delta V_2)$ isn't v.g. — The regions of interest

are not like those described; they are more apt to be hyperspheres.

[N.B. for high dimension, most wt. of a hypersphere is close to its "skin"]

I did write a bit ~~more~~ within last month or 2 on the slope of $\left(\frac{1}{1+x} \cdot (\ln x)^2 \right)$ SES 790.00-90

~~It is a bit more complicated than that.~~

In linear regression, t_1 coils are ~ 1 . In non-linear regression (say
coeff of products of x & \log of x , or t_1 coils scales as powers of t_1 in values of
 t_1 variable to be predicted.

Expo. 10 on $G(x, t)$



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

0/6/02
ID
spec
20: 763.40

After I write t-report; I will want to review recent (last yr or so)
"ID" notes for ideas relevant to various problems.
Make list of Diffys in TM: for each \square - tie diffy to pfs. in report
Where they are relevant - \square Make list of places in ID notes where
those diffys had useful work on them.

Review of GSB to present: Listing of parts written, ect.

09/019/02

Expo: for introduction

Hwr! Not strictly upper Lt border of "Introdu" & it discusses great variety of Optim. problems
On t. Gen'd. optm. problem.

Inversion problems and Time limited optimization problems can be regarded as special cases of a ^{more} general optimization problem.
Suppose that we have a function ~~for $G(x, t)$ that stays to reals.~~
real function $G(x, t)$. x is an unknown string ~~of reals.~~
We must trade value for x , which we present ~~at~~ at time, t ,
so that $G(x, t)$ is as large as possible, after G will take

the form $G^*(x) \cdot H(t)$, where $H(t)$ is a decreasing function of t .
An example of such a ~~problem~~ problem; Suppose we had to design an automobile satisfying certain requirements. We want it to be of low cost as possible, yet if we take too long to design it, our competitors will beat us to the market. We have some sort of function $G(x, t)$ that relates the cost of the automobile and the time ^{for the} design ^{to be} completed, to its ~~market~~ comparative value to us.

Inversion problems can be regarded as having a $G(x, t)$ function, in which $G = -\infty$ unless x satisfies $F(x) = 0$; otherwise, $G(x, t) = F - T$
we want to maximize G
In time limited optimization problems, $G(x, t) = -\infty$ for $t > t_0$ for ~~the~~ $t \leq t_0$ $G(x, t) = G_0(x)$, the usual time (limited) optimization function, $(-)$

In ~~the~~ section on updating, we will show how to solve ~~the~~ this general optimization problem.

9/10/02



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

0: 767.40 : So, in ~~revising~~ 738.36-739.03! TM consider that the "update" always means
update of the complete (EPD, [PST]) Grammar - but as with previous of updating, the
way its done, the amount of time spent on finding "new" parts - This will slowly
change as TM "updates".

.06 Ent. for 88, IND problems, for a long while ~~the~~ "era" are always done using
the same PST: using the dist of concs ~~as in A214~~ (as in A214) & the dist of contexts (739.00-90)
Eventually, when ^{many} open probs have been solved, & we have good PST's for
open, we will begin to try using different PST's for induction - eventually,
the PST's for induction will be part of the (EPD, [PST]) grammar.

BUT it may be possible to have a stand-alone pure QA system & stay in it.

'initial era' of .06 and so fairly far. It may be gooden of for SM!

~~There was very old historical data, TM could discover names before trust of SM says does. - It could be~~
a useful kind to TSQ!

.06 So, w. P. considering of 767.36 - 768.11 & 738.36-739.03 it becomes possible
to decline the whole system

Perhaps have section ~~the~~ devoted to .16 stuff could "putting it all together"
(a summary)

SN for QA (is possibly ordinary Sequential predn) we can simply code sequentially -
looking for extensions of the present "Best" code that & corpus pc.

We will retain several "Branch" codes at each pt, if we can, by oversampling -
At any pt. in the such, we ~~will~~ ^{way} have several codes for the "re-cut Branch". Assoc.
w. each is a pc & a cc. We (use solns of hypc & lowcc, so we will tend to remember"
(or "Backtrack to"), trials of hypc $\frac{pc}{cc}$, (= lowest Least = $\frac{cc}{pc}$).

This will work in a somewhat, if we use the right language.

When definitions are made, they are ~~put~~ ^{start} put at the head of the code, so all subsequent
code can use them. After we have the definition set, we can try to re-prepare corpus, to
maximize pc. This re-paring may be easy: we can have varying units of "LA", & try all
parsings that are possible for the next "k" symbols. "k" tells how much "Look Ahead" we want
we have a partially ordered set of partial codes ~~of various~~ routines of various lengths,
we could use the (hypc/symbol coded) as a force to compare a (re-prepare trial) parsing.

This sounds reasonable Put it "On the stack" ~~the~~ 6.11.35 Later very relevant
ideas: SUM(ACS):
739.22

.06 [S10] 511.26 - 512.21 V.G. Summary of TM ideas! (i.e. ideas since 589.
512.22 - .90 Decn of "Such")

0/11/02
JD



00: 768.40 : Continuing to review,

739.0007 **CONTEXT** in searching for END codes v.g. Very imp.

740: on to a prop $1 + x^{(nx)^2}$: Its properties! How to use it.

How to scale it for diffmt zplcng.

742.00 - 745.00 on TSS 7.

→ 746.00 ff this may be still be a serious unsolved problem! A hour finds some new PST's, say.

How do we "reward" this discovery? Why is t. GPD that implements that hour Better

Plan t. ~~GPD~~ GPD that didn't implement it? Did I solve the rancidity?

→ 761.17 is a review of 746.00 - 761.16 (This is t. problem of 746.00)

At 762. ... ± that ~~is~~ 760.16-18 sold it (cost of 414.18 Gov for GPD [PST]) ~~drawn~~

Now, I'm not so sure! Where is t. Motivation to find a PST's of hyp. expected G?

760.16-18 may gain v.g. Av. for all poss. PST's: so if we ask for a v.g. one,

we should get a honest advice? (Hrs, I don't see what looks like

"Reward for discovery of v.g. PST's" ! (This was t. original problem of 746.00 !)

If we satisfy t. Gov of 414.18 wrt (GPD [PST]), will we get t. kind of "Good PST selection" that we want from GPD₂?

→ Well a possible soln. to this problem: t. Gov of 414.18 does get a

good GPD₁ def. for t. PST's, but t. form of t. soln. is usually not directly

usable to find GPD₂'s (PST's) of hyp. "expected" G. Finding a v.g. (expected)

PST from t. GPD₁ def. ~~is~~ is a key to a major discovery process.

It is a diff. optzn. problem. Given GPD₁ is Prob₁, cc₁, find

a ~~PST~~ PST₂ → it is, out. average, better than any other PST₂.

We have a certain time (C) to do this in. We could first wrt. all other PST's, so t. Gov. is really not entirely satisfied (yet).

See 762. 27-40 for more discn. of t. soln. of (20).

In general, it sounds reasonable, but I'm not completely sure of it, by any means!

→ 759.13 (Re Multiple \prod integral) may be a ~~more~~ reasonable approx: It may be realizable (perhaps) by Monte Carlo approx.

We compute $\int_{-\infty}^{\infty} P_j C_j dx$ for PST_j's that are reasonable candidates,

Then we may be able to do 759.13 via M. Carlo.

Or, compute $\sum_{j=1}^n$ part of 759.13 ($\equiv \delta(G)$), then find $\int_{-\infty}^{\infty} \frac{+100 P(G)}{\dots} \cdot \delta(G) dx$ is here



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University

53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 769.90 : Wall! t. "Soln." of $\begin{pmatrix} 760.16-.18 \\ 769.20 \end{pmatrix}$ may solve 746.00 (fr. "Discovery" problem), but
 Actually, fr. "Soln." of 760.16-.18 is v. similar to proposed "Soln." of 747.21-.22 - which
 I rejected at that time! I think. Main criticism is 748.25-.306 - ~~etc~~
 i.e. if all strats are considered, some very A.H. strats would have very success
 on past data, but would be unlikely to be good in future

118
 Hvr, GPD has to consider fr. Q! Given to devn of PST₂ & its past record,
 & fr. records of all other PST's that have been tested, what is D.F. of PST₂ for
 Probj's CC_j? (which has not been tested).

11
 Well GPD₁ can perhaps SHOULD look at fr. complexity of PST₂ as one param to help get
 fr. desired D.F. GPD₁ is a single function that has been evaluated (in theory) maximized for P14.18
 over a very large data set of $\langle \text{Probj's, CC's, G's} \rangle$ triads.
 Remember this GPD₁ function has to operate (along w. GPD₂ creation) in
 a very finite time, so it probably should take less of fr. PST into
 account in trying to predict its future behavior

20
 Does 11-18 ~~data~~ ^{X data} into SM strat. evaln.? Int. corresponding problem
 - We have a set of strats that have been revaluated on various sections
 of a time series of SM quotes. (This includes fr. same strat on different
 time intervals.) I guess what we want is corpus is: ~~etc~~ This is a bit
Muddled of
"Cross Correln"

22
 [strat devn's, section of data, G_j, ~~no. actual data~~ ^{no. actual data} (very bit of new sup. of data), G] we want P.d. of G.
 P.d. is for a corpus of nextup (like 22). We want a function that gives
 good G values for th. actual data (fr. first 5 strats are "Q", th. last G is "A").
 To what extent should our function consider fr. length of strat devn?

30
 A funt that would soon to work for SM, would use fr. μ & σ^2 of strat
 on past data to predict future d.f. of ~~next~~ yield. ~~etc~~ This would give
 by 414.18 ~~etc~~ on most strats that have been considered. (particularly
 if we confine our corpus to somewhat successful strats of short devn Q)

23
STEIN What about STEIN? For strats of "long M" L, T. overall average yield
 will have $\mu = 0$ & σ^2 will \downarrow w. L. ~~etc~~ say $\sigma^2 \propto \frac{1}{L}$, this gives fr. compd. So if fr. yield
 of fr. Strat. is $\mu_0 \pm \sigma^2$ try average w. $\mu = 0$ w. $\sigma^2 \propto \frac{1}{L}$ so $\mu_0 \rightarrow \left(\frac{\mu_0}{\sigma_0^2} + 2L \cdot \phi \right) / \left(\frac{1}{\sigma_0^2} + 2L \right)$
 $\approx \mu_0 \cdot \frac{1}{\sigma_0^2 + 2L}$ for $2L \gg \frac{1}{\sigma_0^2}$ | actually $\sigma^2 \propto \frac{N}{2L}$ (for periodicity)
~~etc~~ so $\mu_0 \rightarrow \mu_0 \cdot \frac{\sigma_0^2}{\sigma_0^2 + 2L}$ which looks pretty bad! - No matter σ^2 does
 \downarrow as rapidly as $2L$?

5000
 774.12
 on STEIN
 Approach



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20: 770A03 It seems that I've not solved the problem!

21: Review of 769.11 - 770A03

761.17 reviewed 746.00 - 761.16!

At ~~769.11~~ ^{769.11} I was concerned that the apparent soln. of 760.16-18 (the general GPD on PSTs using the $\Phi(4,18)$ core) did not address 'discovery of new PSTs' —

which was highly important. 769.20 ff deals adequately w. this (769.20-30, + 762.27-40)

How, main the main problem of 746.00 - 761.16 was not solved — the prob is very complex (A.H.)

strats would work to problems rapidly w. big G, yet not extrapolate in to new problems. See 30

Some of the ideas in 746.00-761.16 addressed by Q: i.e. ϵ_3 methods

of 762.11-15 (1) SOY, SM strats evaluated by coding, (2) A function of GPD \rightarrow GPD₂ (integrals \approx product)

762.20 suggest review of AAE, SOY, SM strats. PSTs, STEIN.

770.00 - .10 discusses inadequacy of GPD drawn in 4 (4,18) core.

The rest of the page discusses persl. approaches that don't seem to work.

Except ~~770.33~~ 770.33 \rightarrow 40 on the STEIN approach — which might be used to work. (along w. the ϵ_3 of 762.15 \rightarrow (753.13))

SN (Standard)! Re to ST eq. 753.13! Could we get something like it by

using, for GPD, opt not raw 4,18, but for each probly assoc \approx assoc w.

(PST_i, Prob_j, G_k) we take ϵ_3 power \approx being to apip or PST_i. This means that for small ϵ_3 the PC will be ≈ 1 is not very dependent on ϵ_3 PC resid

by ϵ_3 . the function being opted by "4,18". This means the function works harder to ward opting data involved w. short PST's.

As a result, the function would be less accurate for small ϵ_3 PST's!

How, Ray is given less wt. in 753.13 \rightarrow (GPD \rightarrow GPD₂).

It's not clear that 19 help 753.13 any!

I want to keep the problem in mind for .07-.08 — That optimizing 4,18

would still give dirty data .07-.08 That Bad, A.H. PST's would not get

v.g. G distribns wrt. to past data, not extrapolate poorly into future.

Is the PST evaln. problem equiv to $\Phi(4,18)$? Given a corpus $[Q_i, G_i]$

to find d.f. for G_j given new (or old) Q_j ? Here, $[Q_i]$ would be the $[PST_i, Prob_i, G_i]$

in a corpus. If PST_j were constant in the corpus then we have an

only $[Prob_j, G_j]$ as Q_j .

Another view: Go back to corpus $[PST_j, Prob_j, G_j]$. say we would be able to find the best fit to corpus well. (1) Could we usually find a \sqrt{PST} that would give



00: 771.40 : 5 previous results for that function? (2) Say the funct. was a universal pd. — Like prior

3 input var ~~... and only need 2 inputs: Q and R~~ $n = +1$

T. fact that its 'universal' makes no difference for two input's

say to 3 inputs were S, Q, R . \rightarrow descrs f. stock function. output is A.

If we sum over all S_i , there ~~are~~ ^{are} probably several S_i that get by scores, ~~and yet~~ }??

extrapolate badly for certain (long, A.H.) Q's (??) $\odot \rightarrow$

Go back to 771.34-35 : say to corpus is $[Q_1^1, Q_2^2, G_i]$ (G_i is ^{always} 'real')

We fit a good funct. viz 414.18! Given \geq ^{PST, Prob, G} 'new' Q_j^2 , can one always find

a Q_j^1 (A.H.) \rightarrow to ~~match~~ stock. funct. ~~found~~, gives \geq 'purposely' by ~~...~~

G d, fi ?

perhaps that's supernatural: PST's are functions. We have a universal checker over

~~...~~ Prereq. So any describable function is legal.

In 08, Q_1^1 is to function, Q_2^2 to ~~prob~~ G_i is f. G evaln. (G_i can be open or closed) ^{function}

Say we optimized via 414.18, a 5 funct/over corpus (PST, Prob, G).

There will always be many A.H. PSTs that will get by ^G scores. Now O^j could look at the lengths of the PST's; BUT I don't see that it has any motivation

to do this! In particular, if the PST's in the corpus have \geq always been

short, ~~...~~ O^j has no reason to believe that long PST's are Bad!

Could this problem be due to "No Free Lunch"? I.e., if a ~~method~~ induction worked, working well on one ~~or~~ sub-corpus, it must work bad on another sub-corpus. The "other" sub-corpus is the one w. "long in puts".

25 [Is it possl. that a very smart ~~...~~ O^j could look at a PST & realize that while it worked well on the past corpus, it is unlikely to work on the new problem.]
E.G. I could do this: I would probably recognize an A.H. PST! But perhaps some PST's — while quite A.H., are not easy to recognize. T. only clue is their length.

30 [It may be that if TM has enough experience in various parts of logical reasoning, Math, & Statistics, — that it could be suspicious of excessively "long" PST's.

How, TM may not be able to become Very Smart before it tries to deal w. this problem — i.e. Being able to solve this problem may be necy. for much smartness!
In this case, we must start w. some perhaps "ruffin' dirty" ways to deal w. the problem. E.g. In the $GPD_1 \rightarrow GPD_2$ Xform,

0/12/02
ID



00: 772.40: I could find a way to cutoff PST's of "excess long M". — perhaps use some trick like that SM method of using strategies to code for corpus.

Another way might be to eq. of 753.13 (S' of T; S ...).

So we use Passo ruff up date methods until TM is smart enough to figure me out until TM is smart enough to figure out something better

This assumes that TM will at this point, be smarter than me ...

So it must be a fairly advanced system at this pt.

Our Approach: Take ^{10 or 20} PSTs of highest mean G. Restrict to not have long codes for PSTs: Use some perhaps use simpler M code which to each diff. ~~code~~ for G PD₂ (equal weights for all 10 or 20 PSTs)

➡ A reasonable approach to the MAIN Problem: In QA induction, after a uniform set of Q's; say we suddenly start an arbitrary different kind of Q sequence! The ms error will be a lot. The convergence function will be slow, the ~~func~~ function needed to generate "the answer" is much more complex than before.

Similarly, if our Q's for (PST, prob) → G induction included only short PST's of a certain kind, then it ~~will~~ ^{should} not want well w. long, A. N. PST's.

In QA induction as well as in PST evaln. problem, whether a new Q is "new" to old Q's so that we ^{can} expect good answers as in past corpus ... I don't know how to deal w. this.

E.g. Say we have an Encyclopaedia - type QATM and Transvers have been ok. for Q's about Physics. We don't know if it would do well in Biology — whether Ency has ^{much} Biology or whether TM knows how to Access Biology info from the Ency.

ABCD Pilot
ABCD Extra
SSO

In General, we Don't know what areas the QATM will be "Good" from knowledge of its "Corpus of Excellence" — (This is true of Humans as well as TMs). In the case of Humans, because of "common sense" behavior,

Humans, we can often guess well, by "Empathy" & our own Experience. However, in ~~finding~~ "finding the best" of a set of PSTs (or strategies SM),

We seem to automatically range over a vast number of induction in accuracy, & we pick "the best" a PST (or strat) ~~that~~ ^{that} looks v.g., but which

TM has very low predictive accuracy. The result is that simply taking what looks like the "best" ^{strat} PST, is a catastrophically poor choice. (QATM says is)



00 : 775.401 One "rule of Thumb": Don't consider PST's that are longer than 1/2 length
Corpus PST. Have lots of corpus PST's & now ~~we~~ ~~over~~ max length.

"Cross Validation" of choices of "Best" will "reduce" error (at cost of ϵ),
but there will always be sufficiently A.H. docs, that will give catastrophically
bad results. So "Old Faithful" will not save us!

It would be nice to have some way of Estimating how much accuracy we
can expect for QATM & A's that are "distinct" from the corpus

How to measure this distance of a Q from ϵ . Set of Q's int. corpus.
Looks like a "Clustering" problem

112 A **STEIN**ish Approach: (Say SM stat Evaln): Say l is doc. length
of strab: For SM Data batch size of S . (No. of data p's).
Get D.F of yield of strabs of length l on corpus of size ϵ . This can be done
by Mt. Carlo. Mean will be ϕ . - we want σ^2 .

119 < It is poss. that this is essentially meaningless - since a stat. can
be of any size σ^2 , if it bats out. We ~~might~~ get around this by having
bat be Kelly-type; i.e. pc of yield be mean since strab of sample. -

20 Some "staring ditters", perhaps.
21 \rightarrow We can perhaps restrict our strategies so that ~~we~~ ^{They} will not have enormous σ^2
.19-.21 is one direction to ~~try~~ try. \rightarrow

for PST evaluation, this is a QA machine, for ϵ corpus, a O_j will give a
pc for each A_j ; we can look at ϵ . μ & σ^2 of Inpc of a O_j -
evaluating functs.

26 Hrs, ~~most~~ most ~~will~~ will assign $pc=0$ to at least 1 of A_j : so $\mu = -\infty!$ **No!** (35)
We could simply ~~discard~~ discard those O_j . - Still, perhaps most O_j would assign very small
pc to one or more A_j . Say we put data in pc domain & reject O_j 's w
any zeros.

Actually, what we want is ϵ . d.f. of pc's of O_j 's - μ & σ^2 may not be
best way to talk about ϵ . d.f.

It would seem that a "Pure STEIN" approach might work!

35 (26) ~~But~~ Actually ϵ . d.f. for G can always go from $-\infty$ to $+\infty$ - say Gaussian D.F.
or may other D.F.'s; so $pc=0$ can't occur. Also, we may want to consider
In pc because we ultimately use products of ϵ pc's for "evaln". (\approx "pc of corpus"),



00: 774.40: Anyway, back to STEIN: $\hat{\mu}$ PST evaln.: output for each PST, Prob. $\hat{\mu} \approx$ ^{confidence usually} P.D. on G , from $-\infty$ to $+\infty$.

So for all PST's of length l , for the k ^{empirical} problems, we get a k vector of PC's $l \mu_{PC}$'s. We know, by MC Carlo estimates get a d.f. of $l \mu_{PC}$, for random problem in list, over all PST's of length l . This gives us a needed σ^2 of PST's' μ of G .

For individual empirical μ_i & σ_i^2 of each PST, we can ~~estimate~~ ^{be} combine w. a. ensemble μ & σ^2 to get a "true" μ for that PST. — we then pick the "peak": "say" is relevant only ^{when} t. ordering of expected mean μ is affected. \langle I forget just under what cond μ_i is true. \rangle

Anyway $l \mu_{PC}$ domain: say t. $\mu_0 \equiv \mu$ of ensemble, σ_0^2 its var.

If μ_i & σ_i^2 are data for k corpus of k problems, for PST_i:

Then "true μ " $\approx \left(\frac{\mu_i \cdot k}{\sigma_i^2} + \frac{\mu_0}{\sigma_0^2} \right) \left(\frac{k}{\sigma_i^2} + \frac{1}{\sigma_0^2} \right)^{-1}$.

If $\frac{1}{\sigma_0^2} \gg \frac{k}{\sigma_i^2}$ then $\mu_i \gg \mu_0$ $\Rightarrow \mu \approx \mu_i$

$$\frac{\mu_i \cdot \frac{k}{\sigma_i^2}}{\frac{k}{\sigma_i^2} + \frac{1}{\sigma_0^2}} + \frac{\mu_0 \cdot \frac{1}{\sigma_0^2}}{\frac{k}{\sigma_i^2} + \frac{1}{\sigma_0^2}} = \frac{\mu_i \cdot \frac{k}{\sigma_i^2} \cdot \sigma_0^2}{1 + \frac{k}{\sigma_i^2} \sigma_0^2} + \frac{\mu_0}{1 + \frac{k}{\sigma_i^2} \sigma_0^2}$$

say $\alpha \equiv \frac{k}{\sigma_i^2} \cdot \frac{1}{\sigma_0^2} \equiv \left(\frac{\sigma_0^2}{\sigma_i^2} \cdot k \right)$

$\alpha \ll 1$

$$\approx \mu_i (\alpha + \alpha^2) + \mu_0 - \alpha \mu_0$$

So if we are comparing many PST's compare $\mu_i \alpha = \mu_i \cdot \frac{k}{\sigma_i^2} / \frac{1}{\sigma_0^2}$

Since k & σ_0^2 are same for all PST's, we compare $\frac{\mu_i}{\sigma_i^2}$ only (μ_i, σ_i^2 are pure numbers, so there is no "dimensional" consideration).

In \ln domain. This, of course, is if we have Gaussian D.F.'s. Since they are not, we may want to use some kind of M. Carlo or mixed MC Carlo, analytically methods. — Perhaps mix w. my method of get PC of HK winners!

Note that we have to take 'ln pc's, because we are interested in product of pc's. This result is disturbing, because for large l (PST length), we can



0/13/02

00: 785.40 : 2 ways to find A.H. PSTs w. v.g. $\frac{\mu_i}{\sigma_i^2}$ that are actually no good!

OH! This is comparing PST's of same l . — But for different l 's
we have different μ_0 & σ_0^2 . For each l we pick PST's w. best $\frac{\mu_i}{\sigma_i^2}$:

Then we compare them w.r.t. $\frac{\mu_i \cdot \sigma_0^2}{\sigma_i^2}$ ~~is overall same~~ so compare

For large l (A.H.), μ_0 will ~~be~~ but I don't know about σ_0^2 . Actually,
 σ_0^2 may not be so small. ~~may~~ μ_0 may be small.

Perhaps $\alpha \equiv \frac{\sigma_0^2}{\sigma_i^2}$ is not small so α is approximately -1.20 ff. are

wrong! — But it's not really important. — T. comparisons can be
done w.o. approxn.

Also, I may be able to get empirical M.C. D.F.'s for

the ~~est~~ ensemble d.f.'s for various l values. — Then do some
kind of M.C. picking of best PST's.

we compare

$$\frac{\mu_i \cdot \sigma_0^2}{\sigma_i^2} + \frac{\mu_0}{K}$$

$$\frac{\mu_i \cdot \sigma_0^2}{\sigma_i^2} + \frac{\mu_0}{\sigma_0^2}$$

μ_0 as a function of l

$$\frac{\mu_i \cdot \sigma_0^2}{\sigma_i^2} + \frac{\mu_0}{\sigma_0^2}$$

10
3

20

30

It would be good to go over to general meaning, interpretation, of the STEIN model,
to see just how relevant it is here.

In the case of Baseball Batting averages: when sz per player is ϕ , then
estimator is ensemble mean (if available). If $sz \geq 1$, clearly ensemble
mean is quite important. As $sz \rightarrow \infty$ / player $\rightarrow \infty$, we are all wt. on individual
player.

For Model, we assume we know ensemble μ_0, σ_0^2 . Then
to get d.f. for each player: It is product of (a) a prior from ensemble, (b) empirical
d.f. from times at bat. For this problem, the β distribution is v.g. because
batting averages are on (0,1).

I had some trouble deciding on an A-prior for the Batting Average Problem —
Don't remember just what it was; but small data changes in it seemed
to give large changes in Batt. ave. estimates! So, we should then
be able to use v. empirical / data to get a good a-prior!

Which a-prior gives the best approx of the data?

5/13/02
FD

No. Params \geq no. of Data Pts : 778.00



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 776.40

Large l corresponds to very large no of param, & gives nutty results!
 So we should at least limit size of l to be $< \infty$ no. of problems (or whatever) or whatever

03

(The l may be able to deal w. no. of coils $>$ no of data pts in a v.g. way
- T. wt of a model is \propto ~~volume~~ \propto h. vol in h. space of pts that bring us with a h. ball of radius \propto about h . data. ~~Actually~~ in this case, l time

one should get better predn. than using fewer param, & data pts.

This is because a greater no. of models can be considered.

~~It's~~ Our way to do this h.v. volume is to assume each param is on (0, 1); or scale it so that it fits. This enables us to compare h. volumes in different nos. of dimensions.

Hvr, in order to do predn, sample from space of pts. that satisfy constraints.

14

is it possl. to use forgo. to evaluate ANN - why no of coils wts \geq no. of data pts. \rightarrow 778.00

19

For report: Since it's "just a report", it doesn't have to be "complete" in any sense. I don't want to write sections on QA predn.

Early IND update using defined functs & ~~then~~ defined cuttimes (sequential coding w. backtracking).

Then, updating: ~~IND~~, ~~OZ~~, IND

For section on Lsrch: copy $\&$ 575, 86, 89 or whatever, then add in part on ll Lsrch & Mt covl Lsrch.

Do discuss PST's, ~~EPD~~.

Do have section on "Putting it all together".

29

But before (or when) I write report, be fairly certain that I have a good picture of the entire system in my mind, & in these notes ().

I had a picture of the system as a whole - involving: the different stages of sophistication of Lsrch, updating code. - look in various to find it. [738.31 - 739.17] (767.26 has vts $\&$, & discussion.

The General Alp method samples from all param space & it wts. & Gene. By sampling near posterior, we save time, & get cost accuracy. A huge trick. Pick random pts that satisfy constraints exactly. Then go to Gaussian random distn. from this pt for sampling. Mult by Gene over copies. \rightarrow 778.00



00: 777.40 ^{spec.} On "Curve fitting" : Say we have an array of poss. params & an indep input.
(777.00 - 90 R)
777.18 spec
each. z_i is i^{th} param; it has a prob. $P_i(z_i)$.

03 If we have 5 data pts & want to use fit first ~~2~~ 7 params: We obtain t_i for
next data pt, by ~~chose~~ α : chose z / pt. in z space using approx. It has $G(z)$
of fit for t_i ~~is known~~ corpus. Use \tilde{z} to predict t_i next pt & error $G(\tilde{z})$ loop $\text{mod}(0.32)$
Do loop .03-.09 until t_i is 0.2 of predn. or acceptable.

04 Unfortunately .03-.04 is too slow; very slow param \tilde{z} will have useful wrt.
09 Go instead! β chose z coeffs as random using their marginals: Compute t_i after 5
coeffs for best fit. From t_i resultant pt, distribute many trials for $G(\tilde{z})$
using z Multivariate d.f. about that pt. width of M.V.D should be enough
to allow G to \downarrow by factor of ϵ say. ~~make~~ ~~predn~~ w. wrt. $G(\tilde{z})$
loop to β (.091).

For linear coeffs (of possibly N.R. funcs): This involves solving d ^{linear} ~~eqns.~~ ~~in~~ ~~quarks~~;
- No big ~~issue~~. I think it is good approx of .03-.04 is much faster.

19 However, I have to "Stretch this out" in more detail: E.g. to case in which
 t_i ^{large} ~~curves~~ was generated by \int ^{T. dist} z coeffs, but we are using z . When we solve z eqns
 z coeffs; we get dirty finding z soln. We can usually find one, but t_i G \downarrow
slowly as we deviate from t_i soln in z certain z space. Orth to that z space,
 t_i G \downarrow rapidly w. distance.

26 Getting Back to "Main Problem": T. ~~to~~ STEIN soln is 776.13!

28 Pick $PST_{i,j} \Rightarrow \frac{k \sigma_{i,j}^2}{\sigma_{i,j}^2} + \frac{1}{\sigma_{i,j}^2}$ is Max. $\sigma_{i,j}^2$ are k params of G d.f. of
 $PST_{i,j}$ - which is t_i i^{th} PST of d dim z space.

30 ~~we assume~~ $\mu_{i,j}$ are params of ensemble of $PST_{i,j}$ of d dim z space.
There are k problems that t_i PSTs have evaluated in past,
we assume $\frac{k}{\sigma_{i,j}^2} \gg \frac{1}{\sigma_{i,j}^2}$, but t_i surprise is not much different (if this is not true,

30 ~~then~~ We are picking for max expected ~~the~~ $\mu_{i,j}$ of PST. Because of 50% effect.
~~the~~ t_i $\mu_{i,j}$ we can expect, will be $\ll \mu_{i,j}$. - But I think it will usually
be t_i "Best bet" for "last PST".

0/13/02

729



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 778.40 : T. Supp. "Evolution" to the PST optimization problem is by no means complete:

What is lacking: (1) A method of testing all PST's of length l , that is not terribly wasteful. ~~WASTEFUL~~ I.e. we want a method \rightarrow t. ^{fraction} of illegal or "obviously bad" PST's is small. (2) A method is needed to find PST's for which -1.28 is large. This could be a big discovery process. - Needing real creativity & ingenuity. (3) Not only "incomplete" but rather various! Evaluating even a few $M_2, l, \sigma_1, 2$ would take too long! - So we are interested in PST's w. $CC = c_0$; Plan such procedure that $CC = k \cdot CC_0$ and then. Actually, we will use a smaller k to save time - which gives a certain undesirable bias for results such!

STEIN

But, the idea of the STEIN analysis may be correct! That testing k problems is always a statistical test of how good the PST is ~~is~~ out-averaged (long run). As such, ~~it is not a good~~ if we consider enough PST's, some of them will have apparent M_2 's that are very far from known true M_2 's, & STEIN seems to be a way to deal w. this.

What I'm looking for is a grammar that will ^{fit thru} ~~fit~~ the known data of ~~the~~ PST applies to known problems & never gives P.D. on any of untested PST's. T. p.d. it is is, for the most part not A.H., since few of the proposed PST's have been ~~tested~~ tested (directly) on the known corpus of problems. (The it could be A.H., in the sense that one can see it will work well on certain of the corpus problems).

[SN] T. Stem Analysis of (774.12-776.40; 778.²⁶⁻⁴⁰) may be o.k. for ~~the~~ SM strategy evaln.. We can often evaluate large batches of stats analytically - Do "hill climbing" on params of them, etc. The idea of fixing l value, getting a "best", then try another l value and getting a "best" seems good. Gradually l until the stats start getting worse.

In general, the PST optimization problem doesn't seem to be in v.g.

Steps: (21) It is about as good as I've gotten. T. problem of 746.00 is 771.00. It's not really been solved: Pto to work with on 761.17 is to work of 771.00 - 779.36. Is probably useful in giving insight on the problem, & solving some special cases.

8/14/02

780



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00: 779.40 : 777.19 - .40 Has a discussion of T. overall "Report". One of the main ideas that needs development/explanation is the sequential coding for induction, (using some backtracking) = "Context".

If the soln to an induction problem can be expressed as a part of a conc. net \rightarrow was done in 589, then it would seem that we can learn it by "sequential coding".

HM: ~~One~~ One objection I had to the TSQ's I wrote in Spark, was that the concs used ~~for~~ a subconc of trials, were always ~~complete~~ complete solns to problems. I think that this may have to be true if I write P's based on conc. nets — in which each needed conc. is actually a soln to a problem (not a sub-function in a soln.)

Some poss. concs that I write try to get TON to (run) use having them be ^{complete} solns. to problems: x^2 ; translational invariance of evaln. of funcs.; It may be that the CJS for concs (and in this way are usually larger than those that are always solns to problems).

Also note: I noted this verid in CJS as I progressed in the Spark. TSQ's: This note has been unhydated if I read "I shut (ing)" properly, to a pc of concs "used again" for the first time. — The I'm not sure of this — Int. Search \rightarrow P's,

I always "decided" each problem soln. automatically. Perhaps if I didn't do this, I'd get a large pc of a ~~conc~~ problem soln. that was "re-used". One way to deal with:

Any ~~conc~~ soln. of a problem is a "tentatively decided conc" — which means we will try to use it for OSL. In general, any sub-net can be used ^{for} trial OSL, but ~~conc~~ funcs that are solns to probs are "more likely" to be useful. We do not know, want to make this, \rightarrow T. optmpt never changes
"More likely", part of δ . prop: It will be a "heuristic".

[SN] In 589 when I ~~had~~ had that conc. net — as a model of how it (ing) took place: ~~One~~ One of the things I was missing, was to "Recognize Alg": the part that realizes that a certain function

0/14/02
IF



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University
53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

60: 796.90 should be applied in this problem. The conc. net as I drew it, was applicable to solving individual problems once you know that it was needed a new function to solve it. ^{once you recognized it.} So first stuff is write about "Recognition functions" (Obs)

was really messy! → But see (29)

As a heuristic trick, I divided the problem into 2 parts: (1) Recognition of how the new problem differs from the old (2) How to xfm Q into A (29). So look at that stuff: One of the old "CPM" trial docs write about it. Also, look at original ID pp. — There is much contradiction, revision, bugs found ect, (perhaps!)

The CPM papers (Sept 8, 01; Nov 8, 01); also in P.C. \ Ray
"CPM" (Con prob. text)

The stuff in ID around 431.00 of a probably better is v.g.: Much of that stuff can be used for "Expo" as written.

Re: 1. Recogn (R) functs: The correspond to "Analysis" in Science: Breaking down down into subdomains that are ~~coherent~~ "coherent" units — how we stuff solve.

Merging of domains corresponds to "Synthesis" — realization that the domains can be integrated by having common soln. to their problems.

Initially I will try to have just one R. (Essentially no R!)

Then teach other R's w. indices — Plan various indices, etc.

By 2 time TM has to update R's using v. 5 methods of (CPM NS, 01)
it will be rather good at induction! (I hope!) — Do say this in the report

Re: Conc. Nets: T. learning of: ~~the~~ appropriate "R" funct could be part of the conc. net; Both R & P have to be (r)net. I don't have a clear idea as to how this is to be done, hvr!

Anyway, I want to write now: 777.19-40, ~~780.00~~ 780.00 ff: T. overall way TM develops — sequentially — what the T-SQ's look like, etc.

My present idea is that one could go pretty far w. @A TM only: No very advanced opt. systems: try to stay w. Con-net ideas! Try to do (29) s. that Recogn functs are also "integrated" in the conc. net

Also, the idea that all heuristics can be expressed as Modifiers of the Guiding PD.
(The ~~the~~ Physics idea is not used in early induction)

20: 781.40: PLAN!

T. Primary Plan: That if one can make a conc net for a problem, then the problem must be solvable by simple conc. data, only. That use of context should \uparrow PC: Also use of R (which is a kind of context) should \uparrow PC.

First: QATM only: Use only "R"; Use Conc Definition: Tentatively define a conc soln. found (780.29)

~~Next~~ See if I can get it to discover concs that are not solns. to problems.
Try context for ^{then R} conc data. Various ways to express context. See if it helps much.

[A heuristic for ^{Context} ~~data~~ data: Find conc. that occurs often! Try to find context that would \uparrow its PC. — ~~Remember~~ Try R is for many different concs.]

DEF Def "A context is any identifiable procedure that may influence the PC of a conc." (1)

A problem type (R) can be a context. Recent concs inserted, can be context's

See how fast — In PC is increasing w. n, # no. of probs. Plus per. See if (context) helps.

Perhaps ~~maybe~~ start w. (4) Induction, then try to teach (5) induction.

First, use Bern. model only. (Perhaps later introduce idea that ~~plan~~)

Def. of A_i can be any func of Q_i — How to do this is unclear at present.

Too maybe it's clearer than "How to get TM to ~~learn~~ learn Bern induction."

To some extent, Bern induction is "built into" the α prpd. — but this is just the α prpd of the O^j functions. This is different than relating to PC of A_i to Q_i .

{ I really have ~~nothing~~ (not/written) about ^{long} Bern induction or induction ~~at all~~ } hardly at all!

At a certain pt., ~~introduce~~ (after "context" is introduced perhaps) I introduce a few R's with indices in the Q's. Then try to get TM to discover relation of Q's to R's so I don't need to put in indices.

Then see if I can get the 5 Heuristics for updating, to work.

Perhaps write up the section on QA (ing. so it is completely understandable is it. It will not ~~yet~~ contain the updating methods of GPD₁ \rightarrow GPD₂ using PST's. These will be introduced after it has been used to solve FNV & OZ probs. (Generalized Opten probs).

I think this is a much better presentation method than previously — in which the reader would be reading lots of stuff w/o understanding how it fits to the book. \rightarrow (785.30) spec.

17

00: (Spec 30) 784. ~~17~~ Re: \$ 754.00: I may be confusing ^{original} f. ~~original~~ w. d. "Summary ~~applied~~".
 T. enters corpus: code, optimally wrt to original ~~applied~~. Only to ^{at least} ~~new stuff~~ need be
 coded wrt to "summary ~~applied~~"

03 Another perhaps relevant pt.: In "sequencial coding of f. QA corpus"
 I have all these definitions: when a new QA comes in, I can try incremental coding of
 the previous corpus in some of its alternative codes ^{2/0} (a) this has to be done
 w. reasonable frequency) reparto all or part of f. entire corpus wrt to existing
 definitions. This usually shouldn't take very long: T. time required
 depends on ^{f.} "Look Ahead" limit. used. (22)

01 P 67, column 3. She is disturbed by Wolfram's claim that
 the principle of computational equivalence "has vastly greater
 implications than any single collection of laws in science" — Wolfram's claim
~~is that~~ I find ~~it~~ ^{Wolfram's claim} to be quite reasonable —
~~to me~~, computational equivalence ^{implies} that the same ^{tech requires}
 for inductive inference can be applied to all phenomena in
 their guts and sciences. This idea has yet to filter through the
^{consciousness}
 collective ~~mind~~ of the ~~sciences~~. □ □

22 (10) A law can be ~~often~~ expressed as a definite function over a "Context" ^{often} ~~data~~.
 — So it's just a Big regularity in f. data. But ~~sometimes~~, a law is specifically
 a way to do search faster.
 So, superficially, it would seem that there are 2 pd's: ① f. original ^{invariant} ~~applied~~ ~~applied~~.
 ② T. pd used to guide search for good codes. Initially (w.o. laws) they are f. same.

0:782.30 D Note 782.30!
 31 2) Using only our R (w_N=1) for a long time, using ~~of~~ ^{of} ~~func~~ ^{defns} ~~is~~ ^{is} ~~from~~
 context defns. Then when a new QA comes in, I do .04 → .10.
 This may greatly simplify early TM — yet have it to solve some
 interesting problems.

I really should go back to my ideas on Z (4) but I wrote ~~about~~ Wolff
 about! Perhaps first try it on English (spaceless) text, to see if it
 works; even see if I can get it to work w. QATM
 786.00

0/19/02 ID

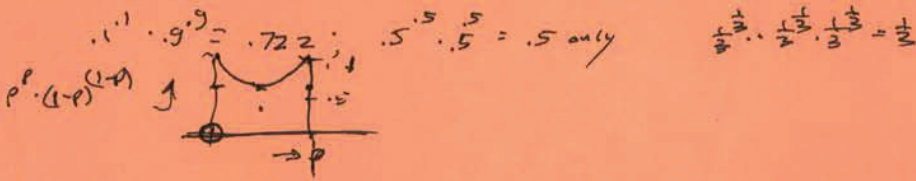
BACKTRACKING .20

00:785.40 : On this "Sequential lang for QA": I'm uncertain as to what I was thinking about!
If I have "all of t. data" then this includes ~~the~~ sols to all Q_A 's, $2 \leq |A|$. If I had to
"reparse" — this would be meaningless — T. last problem soln. would work for all Q_A 's ($2 \leq |A|$).

On way to ~~reparse~~ On only. It has 2 tree structures & its derivations & data.

A "reparse" could take to tree of primitives that represent O_n , & try
to reparse it: conceivably using somewhat different deltas, or t. some
deltas in different order or different frequencies. (i.e. more ~~uniform~~ non-uniform
d. f. over t. deltas gives hyper PC.)

$\{ \in SN \}$ Does $P \neq NP$?
Perhaps (probably) there
are cases that are
undecidable.
It "has to do w. Poincaré"
at $n \rightarrow \infty$, which often
 $\leq \frac{1}{2}$ $\frac{1}{2}$ $\frac{1}{3}$ "undecidable"



Concretely, I will rewrite ~~or~~ one subtree as an equivalent subtree of hyper PC.

T. final functions same but its parts have hyper PC's so different future PC's will be induced.

16 → So far, I don't see any guarantee in this "reparse" !! → 295.25 — Also note RHS of (295.25 - .37)
Good!

BACKTRACKING! Probly best way to do this: Say PC of O_n is P_0 ; we need to
code (Q_n, A_n) ① was best try continuing O_n . ② If that don't work, we move back on O_n ,
& pc "distance" Δ (i.e. we consider all continuation of t. prefix of O_n that is of PC $\frac{P_0}{\Delta}$)
~~try~~ $\Delta = \frac{1}{2}$; so $P_0 \cdot 2$; If nothing found, do $P_0 \cdot 2^2$; $P_0 \cdot 2^3$...

Another possy write to start w. prefix PC = $P_0 \cdot 16$, say, and try random random
continuations; This would have to be done "w.o. replacement" in t. usual L search
manner — otherwise, it's very wasteful.

The idea of " $P_0 \cdot 16$ " ~~is~~; T. code for O_n is a sequence of symbols; linear
to get to $P_0 \cdot 16$, we simply remove final symbols (by 1 until t. product
of pc removed is 16^{-1}).

Unless generated random continuations has some special values, normally
it's easier to just try all code of pc betw P_0 & $P_0 \cdot 16$. These are
generated by "branching", so we always retain t. machine state of "root",
before we branch. This assures us that we never waste time by
repeating part of computation.

"QA soln" could be
several QA's that have
same O_n soln here

Another way to backtrack: Go back 1 (or 2 or ...) QA solutions on O_n .

Normally, in solving a problem, we don't take soln. of best cost, but construct
maybe t. "Backtrack index" should be no. of O_n revisions back

Very IMPT → **SUMACS** (Summary Machines) **.12**

00: 787.40 : Induction only, it may not be diff to search for good s-functs.
But I now certainly not much worried about details of search for better models. This kind of thing is automatically done in the implementation of the algorithm via A2141 - but I need to apply it to creation of s-functs: which is a harder matter. One way is to define objects and contexts. This covers much (but not all) from induction.

So: An outline of the report is 782.30 (785.30) Plan 787.10 - 788.04

784.20
785.20 - 787.09
on backtracking.

They will be 2 what it will treat in the introduction.

So in intro., mention that the QA induction system could work as a stand-alone system for a custom class of induction problems. Subsequent sacking deal in SNU & OE probly is how to bring the ability to solve these probs can be used to improve the QA system, & improve the overall updating system.

12: 787.33 : "Summary Machine for QATM": The Backtracking makes the system close to an ideal "Summary machine". Perhaps there are a few improvements to the system that make the backtracking of summary machines. For one thing, it should save up to 10 or 20 codes to get PD on A. (use just MDL). Backtracking is a way to effectively store many of the best codes. By improving the backtracking using more stored info, it may be possible to speed it up: make it more attractive (as a hobby, have more codes available).

14: Go over QATM very carefully: be sure it's a Summary Machine. This is a very nice way to look at it: it may be possible to put other induction systems in the form of "summary machines" as well!

SUMACS = Summary Machine

16: Another interesting/important property of approximate SUMACS sumacs:

For a good sumac, random input will give exactly the output distribution for future. If it's only an approx. sumacs, (say it has only a finite no. of codes), then random input will give not the future distribution - but it's really not so bad, was concerned to see voice anyway (this is done in practical, RW induction), just how this is done, is unclear. I'd like to be able to put in random input & get not really nice output, but for long and output tracking, it would keep on figuring out replies into sequence that were normal to original sumac.

One way to do this: the prob approx sumac has say 10 short codes for context ("options") - each with its own wt. We could add another channel using the original (from T=0) algorithm. Eventually it would begin to give correct products.

18: 27 It is very uncertain as to how (or if!) it works! → 784.00

Write a lot about how QATM will work (Mainly for me: But Practical!).

For "Pathy is all together" write now a sequence of capabilities of the system and sequence of enhancements to the system: (I have a rather brief list) - but expand it: Also all sort remarks on diffs & how to (perhaps) overcome Problems.

0.0 (788.33) ^{sec}: A basic Q is: for various ^{Potential} "Sumacs" (in induction systems): How & how effectively do they summarize regys in \cdot curves? How well do they enable detection of regys in future?

What they (typically) do not summarize is frequency of ones (but have not yet large enough "definitive" status. Being able to do "OSL" is a special case of ~~the~~ non-yet definable ones.

→ One possible Sumac: Just store t . machine states of M host codes. m can be very large: say 10^6 or 10^8 . We can do this all in RAM or keep it in HDD w/ \approx RAM $m = 2$ parts for swapping while processing. To properly update t ~~state~~ ^{bits} set of machine states, they ~~all~~ ^{bits} have to be tried on each new problem. In a sense, it's "Thorough", but even so, (discounting t skew speed) it is "good enough"? - i.e. could there not be some much faster Sumac methods that would do much better for certain kinds of regularities?

Very P. AZIAT "unc": keeping all of t reasonable alternatives at first, may not involve ~~many~~ ^{much} storage; but t no. of really final modifications will \uparrow rapidly w. n . Perhaps by putting codes in 2 equiv. classes & just store ~~one~~ ^{one} from each class, along w. bias (\equiv wt.) of each class. "Equivt." means they extrapolate ~~to~~ ^{about} t same.

In large, I've been considering discrete codes only. For continuous values, stay discrete data, location of peak (as of precision to isolate it) is wid_r (\equiv wt.) of ~~peak~~ ^{peak}.

SN Most ~~probable~~ ^{probable} \leq predn. will be in form of Bernoulli Basic Sequences. This rule is best stored ~~as~~ ^{as} raw case nos., since updating & prediction can be easily done ~~exactly~~ ^{exactly} from Regy data.

~~At present Sumac (MAYWOOD and ...)~~
A Good Model of Human Predn. is that it is Sumac w. Backtracks & perhaps a lot of strings of alternate codes. Ideally, t unc. of Backtrack vid (its "DepR") will depend on cc available.

SN I had (in t distant past) considered that unity to TSC was t most diff part of t TM project. I have found various diffys in ~~my~~ ^{my} basic understanding of TM. At t present time, the Sumac model idea ~~is~~ ^{is} AZIAT (as a form of unc) may be adequate way to construct a proving TM. that could perhaps learn useful things!

So now, open t TSC design! I hadn't before (I don't remember do in so) any consideration of Recognition functions in t Concnet. In General, it may be that "Euclidization of Human knowledge" is Quite diffit - even for very elementary Algebra!

I had some recent ideas I wrote on TSC Plus second to simplify t problem - try to find Phase notes

- 78 → Perhaps it is possible to show that t Sumac Model using AZIAT & backtracking & codes saved, is about t Best one can do in induction! - (Tho it would seem that other optzn. techniques would often be better.) → 790.00 sec

Spec 759.40: walls, perhaps not! T. model of 759.38 will, w. a certain CB, get a certain set of "Go (short) Codes". It does this by exhaustive search in ~~the~~ ^{the} ~~code~~ ^{code}. It is certainly conceivable that there are better ways to look for "short codes". E.G. An Alg. could look at corpus to be coded & decide that it has a certain class of regularities (or is likely to have a certain class of regys) — then it would look for these regys.

An alog. v.s. .02: In order to do that, the system would have to know about a certain class of regys being "relevant". Discovery of such a class could only be by earlier discovery of certain "short codes" — which could be done by the 759.38 Model.

A NOTE on why L-fact (Post. ~~is~~ INV, OZ, (IND?)) ~~and~~ ^{and} ~~that~~ ^{that} ~~Probs~~ ^{Probs} can be solved in ~~the~~ "not much clever than know in it, the no. of probs".

I expect, that with adequate context, discovering needed abs/solns of new problems should be about the same (unless the new probs are, in essence, "much harder"). Evidence for this: Post it is ~~is~~ ^{is} ~~trivial~~ ^{trivial} for humans; I assume that humans use some general context based search for solns as QATM ~~uses~~ ^{uses}. Note that "Context" can be very general.

Discovery of certain contexts can be ~~done~~ ^{done} ~~in~~ ⁱⁿ ~~the~~ ^{the} ~~notes~~ ^{notes} in Concept notes.

T. imp. idea is that identifying proper context usually will adequately reduce the search space. If this is not possible, then a human would also have trouble solving that problem.

I may be using the TSO in different ways: Below. QA; & INV/OZ:

In QA, I have this concept to (partially) order rules, needed for prob solving.

— Maybe not to differ! In QA using Sumac; we modify the guiding pd so that search for solns. ~~of~~ ^{of} the next problem is much simplified. This is true in INV/OZ as well!

⊙ → N.B. At the end of the section on QA induction, the reader should know how to use L-such (& other refinements) to ~~use~~ ^{work} a TSO: As is, the EXPO I've written is not adequate!

The Abstract gives the general idea of ~~the~~ ^{the} ~~what's~~ ^{what's} ~~being~~ ^{being} ~~done~~ ^{done} & roughly how they work: (Hence, the actual paper will start out w "a special kind of OZ problem", i.e. QA.

Perhaps at end of intro (as ~~is~~ ^{is} written):

submit { Before we tell how Invariance of operation problems are solved, we will describe a system that

Because of the critical role of induction in the general updating process, we will first (section) describe a system for ^{a fairly general kind of} induction based on L-search. ^{Though} This system can operate as a "standalone" system without ~~any~~ ^{any} ~~facilities~~ ^{facilities} ~~to~~ ^{to} ~~work~~ ^{work} ~~in~~ ⁱⁿ ~~various~~ ^{various} ~~for~~ ^{for} ~~induction~~ ^{induction} only,

We expect that

00 790.90:

~~of optimization problems~~, ~~the ability to do induction~~ ~~will be~~
much improved after it has worked many problems in inversion and optimization.
be can applied to ~~more general problems~~ ~~in inversion and optimization problems~~,
the updating of

After discussion of updating is "Analysis & Synthesis".

Footnote P3 of Oct 17, 02

IP In the foregoing we have allowed $R(Q_i)$ to be 0 or 1 only; ~~it is~~ ~~not~~ ~~possible~~ ~~to~~ ~~have~~ ~~a~~ ~~system~~ ~~in~~ ~~which~~ ~~$R(Q_i)$~~ ~~is~~ ~~a~~ ~~real~~ ~~between~~ ~~0~~ ~~and~~ ~~1~~
It is ~~possible~~ also possible to have ~~a~~ ~~system~~ ~~in~~ ~~which~~ ~~$R(Q_i)$~~ ~~is~~ ~~a~~ ~~real~~ ~~between~~ ~~0~~ ~~and~~ ~~1~~
and overlap of R 's are ~~possible~~ ~~we~~ ~~will~~ ~~not~~ ~~discuss~~ ~~such~~ ~~systems~~ ~~in~~ ~~this~~ ~~paper~~.

In the early history of the system, we ~~did~~ allow only one R . - So there is only one
"kind" of problem. It is ~~possible~~ possible to do a fair amount of ^{useful} ~~learning~~ ~~with~~ ~~this~~ ~~within~~ ~~this~~
constraint. Another simplifying constraint, is that ~~of~~ ~~the~~ ~~values~~ ~~0~~ ~~and~~ ~~1~~ only.

To start off, we give you Q_1, A_1 , our first update problem. we want
function $F^j(\cdot)$ such that ~~it~~ ~~is~~ ~~as~~ ~~close~~ ~~as~~ ~~possible~~ ~~to~~ ~~have~~ ~~$F^j(Q_1) = A_1$~~ .

This corresponds to a kind of discrete induction. It can be used for training sequences
in parts of mathematics - in which there is only one ~~correct~~ ~~answer~~ ~~even~~ ~~for~~ ~~a~~ ~~problem~~.

In this case $A_i = F^j(Q_i)$, and we want to find an F^j such that for all Q_i ,
 $A_i = F^j(Q_i)$ and ~~it~~ ~~is~~ ~~as~~ ~~close~~ ~~as~~ ~~possible~~ ~~to~~ ~~have~~ ~~$F^j(Q_i) = A_i$~~ (5)

and 2^j , the a priori probability of $F^j(\cdot)$, is as large as possible \rightarrow 792.03
To start off we ~~only~~ have one Q_1, A_1 pair - we can use L search to find a F^j

in the following way: Using the technique of Appendix 1, to assign 2^j values ^{first} ~~we~~ ~~add~~ ~~to~~ ~~the~~ ~~"~~ ~~argument~~ ~~list~~"
We do an exhaustive search of all possible F^j such that ~~it~~ ~~is~~ ~~as~~ ~~close~~ ~~as~~ ~~possible~~ ~~to~~ ~~have~~ ~~$F^j(Q_1) = A_1$~~
and test F^j is less than 2^j . To ~~be~~ ~~able~~ ~~to~~ ~~execute~~ ~~one~~ ~~computer~~
[The value of T is not critical] operation. \wedge If no suitable F^j is found, we double T and repeat the search.

This doubling is repeated until we find a F^j such that $F^j(Q_1) = A_1$. If we ~~can~~
a hard time, we ~~keep~~ continue doubling in the hope that we will find a ~~substitute~~

F^j with larger 2^j . If we do, we ~~will~~ preserve in memory, all of the F^j functions that worked - though
when the ~~next~~ Q_2 arrives, our guess A_2 is ~~not~~ ~~the~~ ~~best~~ ~~one~~ ~~found~~.

When the true A_2 is revealed, we don't have to do any updating if $F^j(Q_2) = A_2$ for one or more of the F^j 's.
If this is not true, then we have to update. ~~if~~ ~~we~~ ~~have~~ ~~any~~ ~~F^j~~ ~~values~~ ~~in~~ ~~memory~~
we try them first. ~~if~~ ~~any~~ ~~work~~, we ~~will~~ use the best (largest 2^j)
for predictions, but keep the others that work, in memory. ~~Any~~ ~~F^j~~ ~~that~~
don't work for both Q_1 and Q_2 are removed from memory. ~~both~~
If we find no F^j in memory that work for Q_1 and Q_2 , we start a new
L search for a F^j that works for both Q_1 and Q_2 .

IP When the true A_2 is revealed, we don't have to do any updating if $F^j(Q_2) = A_2$ for one or more of the F^j 's.
If this is not true, then we have to update. ~~if~~ ~~we~~ ~~have~~ ~~any~~ ~~F^j~~ ~~values~~ ~~in~~ ~~memory~~
we try them first. ~~if~~ ~~any~~ ~~work~~, we ~~will~~ use the best (largest 2^j)
for predictions, but keep the others that work, in memory. ~~Any~~ ~~F^j~~ ~~that~~
don't work for both Q_1 and Q_2 are removed from memory. ~~both~~
If we find no F^j in memory that work for Q_1 and Q_2 , we start a new
L search for a F^j that works for both Q_1 and Q_2 .

omit. If we have any F^j values in memory that work for both Q_1 and Q_2 , we start a new L search for a F^j that works for both Q_1 and Q_2 .

If we find no F^j in memory that work for Q_1 and Q_2 , we start a new L search for a F^j that works for both Q_1 and Q_2 .

L search for a F^j that works for both Q_1 and Q_2 .

is greater
Then we keep
may have several predictions: Then we take the probability
of the prediction $F^j(Q_2)$ is 2^j , the probability
of a prediction F^j is 2^j .

20
21

25: 792.06

It is no longer easy to keep lines in a factory

33

792.07

20: 791.10 : I should ~~have~~ mentioned that (which looks like a function for Q_1, Q_2, F^j) because a "defined function" in Prolog AZ-191 function does ~~nothing~~. (For myself at least, I should the Note Pad 5.15 practice isn't exactly right - it's an approx. way to do one kind of OSC.)

23 : 291.21 To start off, we have one Q_1, A_1 pair. We want a function F_1 , such that $F_1(Q_1) = A_1$ and a_0 , the a priori probability of F_1 , is as large as possible.

Using the ~~trivial~~ ~~method~~ ~~of~~ ~~Appendix~~ AZ language of Appendix 1, we first add Q_1 to the argument list, so the output will always be a function of Q_1 . → 291.25

07: 791.33 To update, we ~~remove~~ Q_1 from the argument list and add to the list the best F_1 that we've found. We then do an L search to find a function with A_2 as output.

If we do, we replace Q_2 by Q_1 in the argument list and set of the output is A_1 . If it is, we have a successful F^2 ; if not, we continue the search, with Q_2 in the list rather than Q_1 . We continue to search until we find a F^2 such that $F^2(Q_1) = A_1$ and $F^2(Q_2) = A_2$.

This enables us to build the functions we've found so far in the past.

This routine continues as we update $Q_3, A_3; Q_4, A_4 \dots$. Each time we find an F that works with all of the examples up to present, we put that F in the argument list.

In updating, we "backtrack" as little as possible. Suppose F^n works for $Q_2, A_2; \dots; Q_{n-1}, A_{n-1}$, but $F^n(Q_n) \neq A_n$. We do a search for a function F^{n+1} that will work with $Q_1, A_1; \dots; Q_n, A_n$. If we can't find one, we try a search with F_1, \dots, F_{n-1} in the argument list, instead.

FN T. says: Sounds a bit unreasonable: why should I content of a list enable to discover a better F^{n+1} ? YET I think it has to! Obvious in backtracking; ~~no~~ ~~close~~ simulation of SUMAC! ↓ arg test ↑ pc's of things a lot, because individual elements of list have their pc's A. This may be discussed in section on "Scaling".

→ If unsuccessful, we remove F_{n-1} from the argument list and try again, etc. This has an effect similar to that of backtracking. Keeping many alternative, non-optimum solutions to problems in memory. They both give useful alternatives when ~~current~~ the present ~~trials~~ ~~failures~~ fail. Of the two techniques, backtracking is much slower, but more exhaustive of possibilities. Having many alternative non-optimum codes gives us a ~~better~~ ~~estimate~~ ~~of~~ ~~away~~ to estimate the probability of a production.

33 Section Scaling. As we increase the number and/or variety of our problems, the search times for a simple system just described, increase very rapidly. Suppose we ~~have~~ ~~just~~ ~~begin~~ ~~with~~ ~~k~~ functions and constants ~~in~~ ~~our~~ ~~code~~ when we begin our training sequence. If we add one function F^k to the list each time we solve a new problem, the probability of ~~any~~ ~~a~~ ~~particular~~ ~~single~~ ~~code~~ ~~element~~ ~~will~~ ~~be~~ ~~about~~ ~~roughly~~ $(k+n)^{-1}$ for n problems. In L search, the search time is proportional to the

time needed to generate and ~~test~~ ^{test} an acceptable solution, divided by the a priori probability assigned to that solution. For $k=10$, $n=20$ and S as little as 5 symbols in the solution, we will ~~have a~~ ^{have a} search time of about $(10+20)^{25}$ ~~times~~ 2.4×10^7

times the processing time of the solution. ~~This would be the factor of~~ ^{0.035} ~~1000~~ ¹⁰⁰⁰ This would be the factor of 1000 if we just ~~had~~ ^{had} that many successful in searching using all 20 previous problem solutions and didn't have to backtrack. Backtracking one problem ~~is~~ ^{is} squares the search time to a factor of about 5.8×10^{14} , which begins to get impractical. Backtracking further would be completely ~~is~~ ^{is} impractical, with the factor of 13.8×10^{21} .

There are several devices we can introduce to deal with this problem. First, the use of ~~context~~ ^{context} recognition functions, R , effectively reduces the value of n (the number of problems) by a factor ~~of~~ ^{of} about 2 .

Another great reduction in search time can be obtained by context discovery. As we have described the AZ language, the probability of a particular symbol at any point in a description is independent of the recent previous symbols, and is independent of the kind of problem or subproblem being solved. By introducing "context recognizers" we can take advantage of the fact that in certain contexts, certain symbols are much more likely than others. Suitable designed

contexts can greatly increase the probabilities of symbols and greatly increase the probabilities of correct problem solutions — thus decreasing the associated search times.

How do we define contexts? The simplest kind of context is a preceding symbol. In prediction of English text, the "space" symbol is very important because the probability distribution of symbols following it is quite different from the overall average probability distribution.

More generally, any ~~subset~~ ^{subset} of positions of a string of n symbols may be regarded as ~~the~~ ^{the} variety of context for the following symbol. Examples of contexts: the last symbol is "sum"; the last 2 symbols are the same; the entire sequence has an odd number of symbols; the number of symbols is a prime number, ...

perhaps discuss English predn. so context dependent.

Context can be expanded beyond the immediate code string. We can

include the problem being solved (Qnt, Amt) and/or any number of previous QAs. We might want to ~~also~~ ^{also} include previous solutions to problems, ~~if~~ ^{if} this information is available in the ~~earlier~~ ^{earlier} part of the ~~the~~ ^{the} function being coded. In human problem solving, a more global context is quite important. We want to know ~~the~~ ^{the} how and where the problem arises — from Chemistry, Physics, Math... These kinds of context are in part in the learning ~~is~~ ^{is} well; and we must first ~~find~~ ^{find} ways to give the system access to such information.

submit

00

003

08

00

29

30

00: 794.40 Training sequence that is adequate in this respect, we still have to augment our system with techniques to discover useful sub-functions. This involves the discovery of repetitions of sub-functions in correct solutions, and the calculation of how much the total a priori probability of the code is changed by defining such functions.

.04R

First try program, it's not interested in defining recursive functs. → 796.07

06 **SN** If we "oversearch" in L such, we would find OSL solutions, but they would not be given to my pc's Ray detector. Whenever we find a soln, we should examine it to see if it has any functs or subfuncts that have appeared before in the soln: function, but were not given "definitions".

0 **SN** On BACKTRACKING: Its cheap to save "state of machine": consists of lists of datum / contacts & case counts for datum. Since datum are ~~data~~ formulae common lists, we just give address of datum, + case counts. Not much RAM needed! Also, because of the structure of search, each node has a list of definitions up to that pt. Branches add datum. (May be some times delete datum?),

15 Super finally, it would seem that there wouldn't be much point to "oversearch" saving results: For backtracking, it would be just as cheap to go to the node that was used, & pursue the branches (i.e. "oversearch"), but only when "future" events dictate.

306
172
#34

20 One advantage of "immediate over-search" is lost - i.e. good pc values for ~~production~~ production. Also, one could find cleaner codes in those worlds that better conds. for "production" for continuations of the codes.

BIG Q: Suppose is searching for a code & nothing fits (or only very poor fits) thus far. At what point do we abandon this branch & do "backtracking"?

23 Proxly much has been written on TREE SEARCH - - but I have important ~~code~~ parameters: (pc & cc) is "conduct of the" that combines well w. (rc & cc) .36

From purg. it looks like there isn't very much one can do to improve the search.

29 796.16 One pogy is to change order of problems. Another is to (sometimes) "reparse" but the presence of the "complete solns" in the definition list, makes this pointless.

30 Maybe "reparse" means earlier parts of corpus, using datum, & pc's of datum, from later part of corpus, but omitting certain datum that would solve the earlier problems immediately. This sounds a bit A.H., but I may be able to justify it (clean MSP, definition & better way).

31 I like the idea of "re parsing" because ① It does seem to be a commonly used scientific discovery method ② I know how to do it, program-wise. In view of ①, I should be able to fix it up (.30-35)

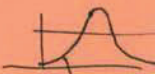
Tim: reason "never" re parsing wouldn't work, is the presence of those datum that would "work on the whole" corpus (or on large parts of it).

.35R

36 .25 When we figure that we've spent much more time on a problem w.o. soln. then we would expect we can do backtracking. The decision to backtrack is how much to backtrack (How many problem solve back) can be "rnd". There are 2 alternatives to Backtracking: One could try "re parsing" 796.06 ff : See NIPS 184.25 ff on Backtracking

In normal Sci Inquiry; we don't have this constraint: so we should be able to explore earlier functions using codes that are later. Actually, parsing using datum & pc joins, is something I'll try to do cause it's better than the current one.

Pin would seem to be very close to what AZ does in evaluating APi's. - It meter for condition & context definition.



Handwritten mathematical notes: $x^3 2^{-x}$, $2^x = x^3$, $2^5 = 32$, $x^3 = 2^4$, $5^3 = 125$. Includes a small diagram of a curve with a shaded area under it.

Transmission Louisiana Ohio.

00: 795.40: -1.10-.40 (for backtracking, reparsing, etc.) looks v.g. It looks like the L search phase at f. QA problem could be done not so badly!

The -1.10-.40 certainly needs some work; it will probably "Mura-ontas" work w.o. backcheck for some v.g.: TSG's. dIND -> SIND Partis (ing "Barroulli Induction"; Ram more general type. Inval to work on this Mess. Hvr, Efficient Backtracking (Note = 1.30-.40) is imp of probly

06: Reparsing is very imp. (-1.25-.35) (-1.32-.35 points way to solution).

07: -1.04-(-1.35) R (on reparsing): Actually, defining "Reparsing" isn't so hard!

Say we have a soln to an early problem: we express it. soln in primitive constants & operators. We can try to reparse this soln, in terms of all of the data used in it whole (f. Juvra) compos. This would only be useful if I had defined sub-functions that were not complete

Solns. of problem. So this "Reparsing" can't occur until I have defined enough sub-functions to make it possible.

14: So, at this point, I have the problem divided into 3 "Eras": The first Phase is f. Early QATM: L search in an essentially Sequential manner. The refinements are discussed in 2.00 ff: Tree Search, to depth of k. $\frac{5k}{p_2}$. (k = depth index).

Some Hours "reference" Mentis: Backtracking, Reparsing, Possibly oversearching to find extra codes (Bursac -1.15 for "objection" to this). Re ordering problems. [In this Era/Phase I. L search is blind: It doesn't look at its "target"; f. sub corpus to be coded. Era II is use of induction to build up. Guiding GPV for INVA "OZ" probs (by "OZ" means "Opten").

20: There is a sequence of refinements on how this is done.

Phase Era III is when QA induction begins to use General PST's to do very general search of function of 14.18, rather than the simple QATM Tree L search of 14.14 ff.

Phase Era II & III may be mixed, hvr. (??) - well, there is a gradual transition from II to III as

the update algorithm for QA is progressively Modified.

Grac has 791 - 794 + 797

26: 797. ... starts the maps vacant now in "finishing" section QATM;

which needs be done: (a) Appendix 2: for conv. from proof.

(b) Editout 794 = 6-795.05) on detecting sub-functions that occur > 1 times

(c) some text by my @ into the paper (I think I wrote this already!)

30: 766. 10-40 Expts on General QATM problem (the varying G(X,t))

735. 07-21 Biblio outline of stuff from 7221 (720.10-722.15) -> Appendix B (or Appendix 2)

I probably want to rewrite the beginning of Appendix 2, explaining just what $\mu(A(Q))$ is.

722. 25-.28 gives a biblio review of "to report thus far". I probably want to change the intro, hvr.

Instructions for Typing

: Sections that need to be typed to fill out discrete QATM:

~~783.20~~ "Practical Induction" 783.20 → .40 ~~3~~ ~~Primer~~ Maybe don't type Re's yet: I
don't yet know ^{where} ~~where~~ to put it!

~~791.05~~ following # on "Analysis and Synthesis":

791.05 - .21

792.03 - .06

791.25 - .93

792.07 - .40

793.00 - .40 (Insert (7.94.19 - .20) ~~at~~ at 793.30) (insert 794.24 - .34 at 793.035)

794.00 - .18

794.35 - .40, 795.00 - .05

~~important~~ but
don't know ~~where~~ to
put it!

The file to which this stuff is to be added is Comprobi.Tex 10/17/02



(703.90) : If the index, j , is a continuous parameter, the sum is replaced by an integral (often multidimensional). If j is continuous and

$$P(O_j) \prod_{i=1}^n$$

→ T. Entro. of t. CPM, Nov 2001 paper has a v.g. explain of t. QA problem. → (15)

0. Make list of imp. docs that I want to include in either the report or my version of t. report. Perhaps do this on computer, so I can insert lines, & move lines.

13 (10) → Also, it explains (to some extent), why 4.14.18 is a good approx of it. If t. pems. in CPM Nov'01 are, indeed, indep. then their wt'd sum being used as approx., means that t. wt's are " $\alpha_i \beta_i$ ". — Which means, t. larger $\alpha_i \beta_i$, t. more confidence we have in it ref. to other pems. (Maybe rather to So 1642 "I-21 Pems Method".

19 Since each "Pem" must have a diff. code, t. codes for diff. pems **cannot** **overlap** (I'm assuming "2 part codes"), so we do, indeed use $\alpha_i \beta_i$ to tell wts. of 2 diff. pems.

→ In t. report, I could discuss this a bit. General discuss. of "Optimization": "Just how Good is t. System?"

So: Present state of report:
Outline of Entire report: 696.13-90 : I should have put in Section of **GCPD**.

T section of **GCPD** could explain about diff. types of CPD's
How index nos. are used in Quirios.
The commonest modes of use of t. GCPD : 1) problem → t. on PST's
2) expresses soln. of external induction problem. 3) $BEPD_{g1} \rightarrow GEPD_{g2}$
used in Lsrch's Modibus of Lsrch. 2 Types of GCPD
3) Guiding Lsrch's Modibus of Lsrch. 4) Guiding other PST's. o.f. WON-tion
Evolution of Proc System optimality of Lsrch. → see (723.35)

Perhaps write
PPS of various
Generalizations,
Extensions,
Expanded discuss
of system.
with t. PPS —
Maybe put in "Report"
But do put in
my "Report".

34
35 → We will describe the initial configuration of Proc system, the contents of its GCPD and AST data structures, & show how Proc system works, using Proc programming primitives & describe problem solving mechanisms. Cognitive devices, and discuss the limitations of Proc devices. → (705.20)



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

001 ; Discussion of "report for" 722.25 ; 2. general report outlined 696.13-40

The outline has detailed discussion of parts of the system. The Abstract & Introduction Give a general idea of how the all works, but I'm not sure a reader could put it all together. - There is no section telling how the system as a whole works.

104 I did suggest § IV an "Evolution of the System" → Maybe "Operation of System"!

Here we revert to the operation of the system as a whole; Starting ~~with~~ as a baby is working Para 1.1, Adaptive Lsrch, ~~and general PST's~~, Scaling problems, TSO's (?)

perhaps make TSO design a sub-section of Section V

I will also discuss OZ problems & INU problems. (The these will be discussed in the section on GCPD - how GCPD differs for various kinds of OZ/INU prob's, as well as TM)

How TSO guides system design (The manual/lecture from the TSO) Various PST's are designed by user to deal with certain ~~problems~~ prob's into TSO.

This initially training of TM is interaction betw TM, TSO, trainer. Trainer modifies TSO in response to ~~TM's~~ TM's response to TSO. Also heurs are discovered by Trainer - But need to be realized (trained into TM) by the TSO.

→ Perhaps write more on optimality of 4.14.18 (for my copy) at least, give ^{Biblio} references.

Biblio: (Ends on 701.19): T. whole dis. seems to be 697.03-699.40, 701.00-19
701.19-19 may be best here, but the dis. in CPM NS, 01 page 2 may be better; however, the idea is that 4.14.18 is U.S. (maybe good as well as), & we want to get as close to it as possible. T. longer but total wt. is the closer we get to it.

121 { SN 703.00-21 discusses QA in which j of O_j is a continuous parameter. (A's are continuous ~~diff's~~ & t. ^{Actualy, f. A's need not be continuous; Only j needs to be continuous; Int. Bern. seq. j is f. continuous param, but f. A's are usually}) (Hessian Approximation)
I may want to include this some how (maybe "Footnote" @).

Re: Expo T. Abstract & Introd. give a general idea of how the system works; But before the detailed sections on QA, PST's GCPD is updating, I think the reader should have a better idea of what the system does before reading that stuff. It may be possible to have "QA" intro after - in tradn. However, w/o understanding Lsrch, it's hard to do the early system.

I may have said that Lsrch uses a P.D. to guide its search. Then explain that after a problem is solved, this P.D. is ^{improved} (modified) in view of the newly solved problem; mention that this P.D. is meant to reflect the experience of TM to date. Insolving the problem types that this part of the GCPD is designed for

35 However, as it is now, perhaps the reader can follow it. T. sections on QA, (Lsrch's PST) are understandable by themselves. T. section on GCPD (704.29-39) is a kind of summary of previous mentions of it. Actually, 704.29-39 doesn't say very much. I write to make it an analog of Google. - ~~From idea~~ An up idea to get across is that there are many P.D.'s that TM uses, but they can be all regarded as one big P.D. (GCPD), which enables ~~smooth~~ Transfer (only between various kinds of P.D.'s).



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

20

[SN] Some time ago (4 months?) I worked on folp problem; I think I solved it!

In updates, what one wants to do is optimize f. "Grammar" first assigns "pc's of best" to all poss. PST's. How is this done? Well first we get a grammar first gets to branch

poss. PST's \geq pd on G & cc for ϵ i. problem in question. This is a genus of $(GCPD, [PST])_1$ to eat to $(GCPD, [PST])_2$ We do standard integration to get a D.P. over

PST's that will be "i. best" — **T-nature of t**, integration depends on t. problem depends on how

— whether its INV or its OZ

$G(\epsilon)$ is a funct. of ϵ cc as well as X.

This entire approach omits diff. of "Corollary" Betw. Cands

O.K: so Nexts

Section on PST's: see 696.21 for Cont. outline of entire paper

1) Defines Verb. "simplest" Lsrch: $(T \leftarrow 2T)$, using p.d. or "codes"
Verb. II Lsrch (time sharing)

"fl" random.

2) Use of Lsrch for ϵ p.d.s. since $T \leftarrow 2T$, II Lsrch, II random Lsrch.

② Look up close to $\frac{cc}{pc}$ order. Some times best.

3) One trouble w. Lsrch: pc is too small: How to fix! "Adaptive Lsrch". \uparrow pc of soln. Update p.d. after each T , spend T on updating. Guiding p.d. than Lsrch. CWS "entire work"

Kind of Adaptive Lsrch: Weaver in a hurry to update because Cands may be correlated.

Time share revision of P.D. w. Lsrch. "50/50" same: Do II Lsrch; I think I worked

out to Mechanics of. P.d.s. work on all cands w. min $\frac{cc}{pc}$; when it is no longer work.

min, work on 2 cands w. min $\frac{cc}{pc}$ keep working and recruit in new Cands into

Set as soon as its $\frac{cc}{pc}$ it fast to convert sol.

Say we doubt $\frac{cc}{pc}$ off. work on set; Run we hunt for cands w. $\frac{cc}{pc} < \epsilon$ new

level. In our random(?) search, we will keep track of track (in B's)

of $\frac{cc}{pc}$ of each cand we examine, so we can know when to recruit it.

We could do random search on search in P.d. order. "22 may not work"

P.D. (P.d.s) are continuously changing to "BIX" idea of

Impehable in 733.24, in formula 732.26 See 733.06 for Notes

This is "Base" in my working of "T. comp"

28

Another method of searching: TM initially selects ϵ "Best" (via GCPD)

PST for t. problem. As it's working on that PST it's (Time share) working on revision

of GCPD, ϵ in view of t. Success or failure of ϵ current PST.

When (and if), ϵ expected additional time for soln. of t. present PST, becomes lower

than expected soln. time of a diffrent PST, it will switch to that PST -

It will then work on that PST until (or if) a switching criterion is met then

sends it to a new PST.

34

35

[SN] In last section of General Conclusions, ack: Concluding Remarks "prob solving method"

We have described a very General Framework that we expect can be developed into a system of very great power. The problem solving technique that we will start with, Lsrch and adaptive Lsrch, will certainly solve many problems, but

9/1/02
ED



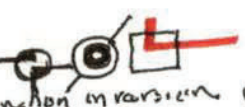
JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

eventually they will be limited by the scaling consideration of section . The most important feature of the system is ~~not~~ its flexibility — The particular methods used are able to accommodate any possible PST — be it added by the trainer, or ~~discovered~~ discovered by the ^{universal search of the} system itself. The development of this system will never reach a "dead end" unless the ~~process of machine inference, itself, is necessarily~~ ^{The limitation being only that of the hardware — which may} ~~limited~~ ^{limited}. It such a limitation exists, it is far beyond what was as humans have done in the past. ~~certainly a tolerable limitation!~~ ^{expans} will continue to double in speed ^{about} every 18 months ~~for quite a few years.~~

12V
200 Hz
.05 uF
= ~~area~~

Levin first describe his Universal sequential search in 1972 (Lev 72).
Though his paper didn't tell how to ~~do it~~ actually do it, his description was enough to suggest several ~~possible~~ implementations.

The problem was function inversion: Given a function $f(x)$ from strings and/or numbers, to strings and/or numbers: Given s a possible range object, S , to find argument x such that $f(x) = s$. We assume that evaluating $f(x)$ does not take very much time, it is a slowly increasing function of the number of bits in s .



For the function inversion problem described in the introduction, Levin proposed a universal ~~sequential~~ sequential search algorithm (Lev 72). Though he didn't ~~give~~ give the algorithm in his paper, ~~his~~ ^{his} description was enough to suggest several possible implementations. ~~he~~ ^{he} described the algorithm in ~~didn't~~ ^{didn't} define the algorithm in any detail, his description was enough to suggest several implementations. (Li. Vi. 7 93 97)

Theorem on the lower bound on how good a search algorithm can be.

1. Consider the set of all functions G , that (output $f()$ and the object s , and give a ~~total~~ total x as output. Let us assume a probability distribution $P(G_j)$ on all such functions. Then, for a small time T_0 we spend ^{maximum} time $P(G_j)$ on each G_j and testing its candidate x_j on $f()$. If we are successful, we quit. If not, ~~otherwise~~ ^{otherwise} it will take us $\frac{T_0}{P(G_j)}$ to test all of the G_j in this way. Since $\sum P(G_j) \leq 1$ our total time will be $< T_0$. If we are not successful, ~~we double~~ ^{we double} T_0 to $2T_0$ and try again on the whole set of G_j 's. We keep doubling ~~the~~ ^{the} time limit and searching until we find a solution.

8/4/02 IX



JOINT CENTER FOR URBAN STUDIES of MIT and Harvard University 53 Church Street Cambridge, Massachusetts 02138 (617) 495-7908

00:695.40: We like approximate optimum path to ~~████~~ -1.37 by \approx "Lsrch". We move along successively cards, keeping track of total pc thus far, & we regulate \Rightarrow cc so to be below a certain multiple of that pc: \uparrow multiple \uparrow binary w. time, it was as if doing "Lsrch".

[Packer -1.37: Say we were able to list cards in /pc order, in ^{best} source that ~~████~~ After this trial, total pc of will may may occur all poss in total orderings. How. \uparrow optimum in sequence is probly not a prefix of \uparrow optimum (kth) sequence \odot !

- Any way, probably various approx. solns to \uparrow card ordering problem will be made clearer by \uparrow consideration of -1.37 - that \Rightarrow an optimum card sequence.

Fact is, hrs & conditions that define \uparrow soln -1.37 do not obtain in "Lsrch":

One does not know all \uparrow pc's & cc's in advance - still, \uparrow concept of -1.37 can be useful.

In designing Approxns. to be used in modeling and implementing GCPD for "Lsrch".

13:694.20 Ok - Back to T. Report.

Main section of report: Introdn: 693.22 - 694.12

15 **I** QA problems: How many are O2 probs: Guro 4-18: Guro convergence Perm. & Significance. Explain that to define a QA. problem, we need only give \uparrow ST corpus $\Sigma Q_i A_i Z_i a_i \dots, Q_{n+1}$ Also \uparrow method of giving inputs to ~~████~~ ON. This defines \uparrow problem: To solve problem one needs much more: \uparrow R.

20 defn. of c. problem tells what edm. looks like - not how to find it. \rightarrow 699.20
N.B. EM OJ(AQ), Jozsa be contrasted param: see also 703.00-01 (Guro to Similitude Paper on Harfotras)

21 **II** PST's \odot Descr. Lsrch & its variations: \uparrow imp't properties: $\{CJS\}$
Also \uparrow 3 GHT's. ~~████~~ \uparrow $\frac{pc}{cc}$ order is best \odot CJS \odot Lsrch is close to best: is \approx best if all info in GCPD is only blind search is allowed.

Descr. Adaptive Lsrch: \odot Change PD for each problem. \odot Change PD at $T \leftarrow ZT$ every card. \odot Change PD every several cards. \odot Change PD every card.
Note $P_1 \neq P_2$, is no longer Lsrch; \uparrow $\frac{pc}{cc}$ \approx $\frac{pc}{cc}$ \approx $\frac{pc}{cc}$ is not a priori known. \rightarrow see 724.10-24 for \approx 696.15-28

28 - even when soln is known. \rightarrow other PST's: They can contain such as; they can be any Pcs; \rightarrow 761.00-16 on mixing P553 and G-PD, mixing

30 ~~████~~ Maybe discuss GA's; RAINN's here. \rightarrow Main PST! They can call other PST's or call \uparrow Main PST! \rightarrow 723.04-10 on P55

III Updating: \odot Simple updating: 695.15-25 (not a diff of \rightarrow 625.25-29) \odot I think some kind of naive version of this may be somewhat workable. It is certain that people will question those to discuss them.

\odot updating using GCPD_{1,2} w/ M GCPD₂ being binary indic pc's of cards.
 \odot " " " allowing "corrections" (dependencies of pc's & cards)

IV The General soln. of 695.00-07 \rightarrow Evolution of system (704.35) \leftarrow possibly discuss GA & RAINN here? \rightarrow (see 723.04-10 on P55)

V Discuss. desirability of TSO design. Refer to 585, 89, 94.
Buy \uparrow give general form of "conceptual" (collide on "abstraction" next)
 \uparrow used by counter-try Pcs. \odot

Functional form of G will cause search invariant. Some? This must be true of any optimization search system. It is not enough that G is a transducer. (Quality of GCPD)

Expo

Abcdefghijklmnop

is expected
be able to
may not be possible

00: 798.40

The system described thus far, is a pure induction system and ~~can't handle~~ ^{some} worst problems of some difficulty. We will show how it can be extended to

work Inversion problems and ~~some~~ ^{both} some limited optimization problems.

~~Both kinds of problems can be solved by Levin's~~ ^{However}

Universal Search algorithm (Lsearch).

can only solve ^(small) problems, since the solution time is exponential in the ^{length} description of the solution.

Lsearch uses a universal probability distribution to guide the search. 812.00 ABCDE ABCDE ABCDE
ABCD ABC AB

SN I need Section on Lsearch: Oz Lsearch is Random Lsearch have not been ^{disc} "Entire Literature". I don't know if Lsearch has been:

Also to (X:dep R) or they factor for Lsearch. depends how one does it, hvr.

SN for NIPS Talk Give reasons why "it" should work:

- 1) The convergence from: what it means (counter ex.).
- 2) QATM: The SUMAC algo: How we should be able to write good sumacs w. good update algos

3) For advanced system: to 3 Gumb Mouse Terms: Details of how they work, & why relevant. - to "factor of 2" algo.

But do counting
re adding P
379.26
378.00
- 379.26

Re Sumac: I ~~dis~~ ^{consider} considerably in some detail some time ago: There were a ^{completing} systems. This doc was heavily referenced: Shouldn't be hard to find; Main corpus summarized at top of first 10 or so lines of page.

If one saves "all" of the O^i , v. R. or wts, to be ~~not~~ ^{not} easily updated for each new Q_{int}, A_{int} . $w_i^m \rightarrow w_i^n \cdot O^i(A_{int}, Q_{int})$. If one stores only a finite set of O^i , w_i^i 's they tend to not contain the ~~needed~~ ^{needed} but are critical in future comp. If finite set are used gets "used up". The needed O^i 's that are excluded, were excluded because of low w_i^i (~~or~~ ^{or} small α / β)

I do want to go back to the earlier reference. It is in an ID. - [which ^{found!} since I was at IDSIA, but I'm not sure. (378.00)]

A way to do a Sumac w. an ~~set~~ ^{set} of O^i ~~is~~ ^{is} to make ~~the~~ ^{the} set of O^i ~~is~~ ^{is} a Grammar (\equiv PD) over the O^i . One such Grammar is the set of ^{function} ~~words~~ ^{words} & their wts, "up to now". - Also to set of cont. "Contexts" & their wts.

The reasoning at .25-.39 (only a finite C^i set eventually ~~set~~ ^{set} tails) may be wrong. If ~~the~~ ^{the} finite set could ~~all~~ ^{all} ~~be~~ ^{be} very low C^i (A_{int} / Q_{int}) yet, if they were not zeros, the finite set would survive. If C^i is, would the new set of wts have a reasonable distribution (set of ratios to each other)?

00: 794. to wa under deal w. 299.38 - .40, by putting a small exponent on t. O^2 work!
 i.e. update factor = $O^2(A_{nt} | Q_{nt})$: δ would be small if, on average
 Very few of O^2 's did well.

kind of prob!
 581.06-13

03: **SN** On Back tracking: if one had rite "wts" for all component concs,
 & one did an L search over t. small concs, using Bernoulli's
 partners it would automatically Back track properly!

2% key
 ± 5%

— Say S_n is some acceptable conc to $\{Q_2, A_2\}_{2,1,1,1}$.
 we have its set of S_j ($j=1/n$) in t. "Bernoulli alphabet" along w.
 all of th. other functz that have been defined, for $Q_{A_{nt}}$, one would
 try S_n as a component, but S_{n-1} & S_{n-2} would also be tried...
 S_{n-1} & S_{n-2} amount to "Back tracking".

It may be imp. to note that many of t. S_j components have been "used at least"
 twice; by being used as a component of a later S_j . — But S_n has lost: its presence
 in t. arg. list ("precorpus") is by "courtesy" only — & ultimately justified only
 by OSL. — Hvr O_n does have fair likelihood of legal use in OSL, so we
 put it in t. "precorpus" to facilitate its use in a new funct. concn.

03
 04
 05
 A poss. way to deal with get around low app of t. needed concs: In hunting for Bern concs,
 if we see a certain conc. is repeated a large number of times, we know that even w. low
 app, its a viable conc. One Efficient way to do such a search: Use a sample
 section of t. corpus to see at the concs how often t. conc. repeats. If we multiply P_{cs}
 by t. ratio of sample corpus size, to the whole corpus & sorted, & it looks good we can choose
 env't to deal w. t. low app, we do t. search for the rest of t. corpus.

— Actually, we don't have to do active "rest of corpus" we can do as much as is
 needed to decide about "going further".
 ∴ So the Mike be one way to express t. "SUMAC problem": How to recombine
 the "followed" C^i 's, new C^i 's that have very small apps? Again, I think I
 wrote about stuff like this. (378.00ff) — These approaches are different from
 what those being considered now!

0
 In 378.00 I was concerned w. "contraction/narrowing" of t. d.f. of O^i 's, .19 ff is
 one way to help. We do have to continually bring in new low app O^j 's, if no want
 O^i 's system to work. (I think it should be able to be made "workable"
 since humans seem to do it! — i.e. Humans use SUMAC'S)
 Humans do Backtrack.

20: 800.4 : In 37800ff, I was concerned w. gradual (or/sudden) degradation of the "Set of Good Codes". Let us look at what happens from the pt. of view of the TSP.

In order to "collapse" the "Set of Good Codes" must not have adequate subfunctionality to represent the solns. of the current problem.

Its an inadequacy of the TSP. ... wrt to the "Set of Good Codes". T. Set of Good Codes is mindful of the current G population in GA. $S \supseteq C \supseteq SAC$

For the (at least) problem to be solvable, there must exist a soln. that is a reasonable CJS for at least one (preferably many) In General, its very (instructive/revealing) to look at the TSP & its effect

on the SAC:

Perhaps the "SAC problem" (of its degradation) only occurs at that pt. in the TSP where the Trainer has lost detailed outside knowledge of the SAC

(what codes TM has "acquired") is ~~what~~ what codes are adequate for soln. of the current problem. Trainer may know what codes are "adequate solns" but he will not know how to factor these codes so that TM can acquire them.

SN A possible impl. Aspect of this problem that I've not considered:

"A code B is for" has 2 aspects: its pc is its mean cc per problem or (mean cc per problem). So, we will retain in SAC all solns. w.

pc w. in a certain constant factor, of the best thus far. But we will be warned about codes that have bad $\frac{cc}{pc}$ (a bad CJS), Re Play will be good

pc's, (good for prodn), they would tend to take lots off cc for future problem solns.

The way it looks now! That in the early phase of the TSP (Before it has occurred) There isn't much of a "SAC contraction problem": There is always a solution that is an extension of a possibly single acceptable code. By oversearching in Lsearch, we can have a 2 acceptable code - something w. higher pc's but larger cc's. If we use a CB (agreed to find what the Trainer knows to be "an adequate soln", then TM will at least acquire that soln.

If TM often finds a better soln (factor, higher pc) than those expected by the Trainer, then TM has entered Phase II (all-14) is the continuation of the TSP will probably be made up by the Trainer based on Local Models. (These are approx models of TM & can be sometimes rather "gross".)

In Phase I there need be no "backtracking": Hvr, by oversearching we can have many codes. Some of the codes can be expressed completely in terms of common code additions - perhaps by (parity) by cartesian products of certain sets of solutions.

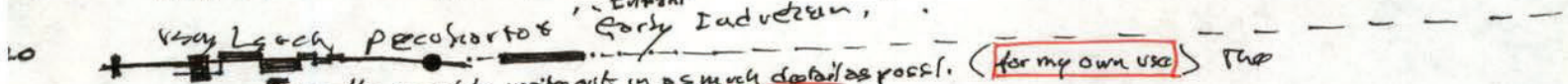
ABCdefg

00: 801.40 : If solns are found by Lsrch I think we want to retain all solns. found. They will tend to be of hy pc a/o of low cc ∴ "useful".

02: ~~mechanics of when~~ ∴ how to back track : Certainly **UNCLEAR!**
I want to Model it More of Human Backtracking - Theory Revision.

In Theory Revision: we back track to previous or earlier problem soln. - But would Cor could be "made to work w. minor changes" - but at that time, was of low pc d/o by cc. (low $\frac{pc}{cc}$). Certain of pr problems since Pr are dependent on Pr soln., so we have Lsrch on these problems, to develop / Post Pr Theory revision Pr. Certain problems will not be dependent on Pr older soln., & we can use the same solns. as before.

During (the lengthy) Theory revision process. We continue to incorporate using the old Theory, perhaps in some "local" revisions (whatever that means!)
It would be good to examine carefully, the process of Theory Revision (in the Sci/Computer) Then Correlate it w. Processes in QATM. It is conceivable that it may involve a kind of inductive process that are needed in a QATM that updates using the full open model. (rather than to update model)



20: I really need to write out in as much detail as possible. (for my own use) The Operation of the "early" QATM. Call this Phase I & II QATM! Phase III is where it uses full QZ techniques to Max 4.14.18.

25: See 801.11 - it has data, discuss of Phase II (this is where TM's response to TSCQ search = "update" → 803.00) is not completely understood by the man.)

One unclear thing! We can use a fixed (CB) time for each ~~update~~ = time spent on previous search(?). or time spent on for finding such? : any way; another way ~~to do it~~ would be "how ~~fast~~ by a value of 4.14.18 downwards better way? In the case of induction it is perhaps easier - often we are so his fixed w. the first soln. we find, because it is of hy pc a prop : (But if it has small cc, we may be able to find a better prop soln. w lower cc - which justifies "over search". It may be nice for the Trainer to decide ~~what~~ level of 4.14.18 to accept - as part of the "update problem" Definition".

1674
210
74 calc
174 Probr
E

00: (802.25)^{spec}: The Main Idea is that if we have an adequate one net for a problem, we should be able to write an adequate rate TSC. "Adequate" means that to recognize as well as to "solve" a given one in the context.

Most ~~hours~~ hours can be put into a net, since most hours are modifus of the (Guiding PD (on a simple level) is more generally can amount to "discovery" or new PST's by assigning "reasonable" PC's to those PST's. — In this case a PST is a soft (probabilistic) reordering of trials.

ncost = "surprised".

How is Impf. during a Run, done? Is it really necessary? Also, of course, Quick Abort — not being served by simple modif. of PD.

Remember the PD for simple QA is for O_n^* 's not PST's. — But (12-14)! functions can (to some extent) implement PST's

2 Hrs., we are looking for "functions" as a function can be defined as a such. On the other hand, such functions have large CC, it would be used only if they

14 had very large PC!

"wired in a bit"

There may be a min set of hours that we have to "unwire a bit", before the TM can learn enough hours to get to Phase III.

20 **[SN]** Image Reconstruction: A blurred image should be expressible more compactly by a high res. image + neural + noise — but it's not clear just how!

T. original ~~image~~ (high res.) image may have much usable redundancy (?).

30 Reason 12-14 seems very imp! Can any PST be regarded as a "function"? Hrs. "functions" can't really implement PST's in general, because ~~they don't know what~~

PST's know what the problem is: T. functions I'm using in QA, Don't. E.g. mt. QA problems a PST would know A_{n+1} as well as Q_{n+1} . It would also "know" that a hypc O_n^* relating Q_{n+1} to A_{n+1} is wanted is not also, we want to by O_n^* (A_{n+1}/Q_{n+1}).

T. "functions" I'm using only know Q_{n+1} & we list/trials in O_n^* PC order.

→ A hour may be able to change to diff. guiding trials, but not it!

30 In Phases I & II it's "blind" in a sense that it doesn't really know what it's looking for! — or that it has a bias (or hypothesis). in these (perhaps) small CC —

So Phase I, II is only able to try funcs of Q_{n+1} in O_n^* PC order (or O_n^* order).

If we modify the PC d.f., we are no longer solving the correct problem: Except that

I to PC is to P.D. for A_{n+1} (not to PD for the entire corpus) — it's to PD for A_{n+1} in view of Q_{n+1} & all previous QA's. — It's a PC of a SUMAC. T. Sumac ~~knows~~ knows Q_{n+1} but not A_{n+1} .

2 propd = "pna"

So it looks like the best views of Phases I, II is as a Sumac: Discovery of functions & Context for functions is simply a definition of the context prior

Could any heuristics be regarded as discovery of context "Contexts"?

I would SEEM that a Sumac approach, using a conc. net, should work --- Res would be certainly Phase I, & possibly part of phase II. A by Q is how to get reasonable pc's for t. concs & debug. t. latest soln. Context is one step in this direction.

I don't remember having worked on this problem from that pt. of view, there is the "recovery" context. w. an "x window" for recovery {or Rundo... which means I'm only remembering stuff used \approx 1 seconds ago} "softness" (Gray) at x window sounds Much Better than Rundo. (20)

So, at present, in Phases I, II, I want Sumac Model: defining functs, \approx OSH, & Context for Funct defn. & the Bernoulli model of A241 are the only things I can find/justify for Phases Phases.

Whether it to which extent, Backtracking is viable/feasible, is unclear.

I could do a fair amt. of "Oversearch" to get extra codes, making Backtracking a bit easier. T. discn. of Backtracking at 802.02-1.20 is good, 800.03-18;

We have a really enormous no. of concs. We really have to assign large amt of pc's to them so prob. pc of entire soln is reasonable. perhaps .20 is one of the major problems --- other run designing t. T5Q itself.

Actually, the scaling problem is, along w. T5Q writing, 4. 2 major diff problems in TM. In doing induction (past stages I, II), we are looking for hyper functions that map Q into A. (Partly: we also need a suitable "R" funct) --- but any way, it is the early A.I. problem of finding a seq. of xvals that brings our thing into another in this case Q into A. The standard methods of sub-goal (deriving), backward & forward chaining & GPS (difference reduction / vector geom), & other common methods. Probably these techniques can be extended to S-induction.

AH! Context Discovery can be implemented by putting in defs. of Contexts in (a/the) "Concept net". One version of "Context" is a structure out of previous corpus (i.e. "retrieved" post fix).

SO: When Peter Baymunk (at Szarb) was getting pc's \approx 10-15 for his problem solns, I think that he was not breaking problems down into sub-functions properly: but while that may have been true, he was probably also getting very low pcs because of a simple "scaling" & not using "Context" which Context allows (could vary "extended context") will be out to get reasonable CJSS --- I'd unclear. Ultimately, we want to see incremental cost of each problem soln. to be a good constant.

PS file: Nov 2, 2002

Generalization of ALP to N dim. lattices

Chris' Question:

Say we have data from $t=0$ to $t=100$. What is the probability distribution at $t=-100$?

One simple answer would be to relabel the data, so $t \rightarrow 100 - t$. We then have data from $t = 0$ to $t = 100$ and we have to predict at $t = 100 + 10$. - which is a standard ALP problem.

However it suggests a more interesting problem: Suppose the data ($t = 0$ to $t = 100$) was really taken from a string of data from $t = -\infty$ to $t = +\infty$. i.e. it really didn't start at $t = 0$.

In this case, to get the probability distribution for $t = -10$, consider a universal Turing machine with no stop symbol. Its input is a unidirectional random binary string. It has an infinite bidirectional work tape. Its output tape is infinite, bidirectional: Initially blank; starts on square zero. This tape is bidirectional, but it cannot erase, only write.

There is a certain probability that the machine will print all 100 known bits and 0 for bit 110. There is also a probability that it will print all 100 known bits and 1 as bit 110. These two probabilities give a semimeasure on the 110th bit.

This model can be easily extended to n dimensional lattices. The "known data" can be any set of lattice points: They need not be contiguous - many data points can be "missing". The predicted bits can be any points on the lattice. The "output tape" will be the n dimensional lattice and we use a special turing machine that can move on the output lattice in $2n$ possible directions and print. The input is as before an infinite unidirectional random string. The work tape is bidirectional and infinite.

We can also use a Turing machine with infinite bidirectional random input tape, infinite bidirectional work tape with both write and erase allowed on the output lattice. Start at any point on the input tape. What is the probability that the machine will, at some point in time, have on the lattice a certain known configuration of data bits, and will have a 0 at a certain specified lattice point ? --What is probability that it will have a 1 at that point ?

This last is of interest because we are using an Unrestricted Turing Machine -- this model is Very General.

allowed on "output tape".
Erasure
No erase ever allowed on output tape

"Needs work" \odot
The machine may or may not have a "stop" state.
As this is formalized, it may be that probability of any particular config being eventually produced is 1.
Hr, as $T \rightarrow \infty$, the no. of times a given config gets produced, may have a fixed ratio to T .

Infinite bidirectional random input tape: ~~the~~ analog of "prefix set": No string is a substring of another in the set. Hr, no "length inequality", we could have an ∞ of ~~contiguous~~ contiguous, non-overlapping, strings of length z . ($z > 1$). $\sum z^{2^i} = \infty$. But, every string must contain at least one prefix set.
So maybe Kraft's still holds. It does hold. And $\sum z^{-2^i} = 1$ for any "complete" set.
Let L_i, R_i be left & right halves of an input code.
For any particular L_i all legal R_i, j ~~is~~ R_i has a prefix set.
So total wt. of all such codes is $\sum z^{-|L_i|}$ and \sum sum of wts of all codes is $\sum z^{-|L_i|}$ which is ≤ 1 .

N DIMENSIONAL ALP

orig. 0th ALP report Systems 804.23-29

10: 804.001 At present, there is plan for early QATH (Phases I & II) ~~is~~ is reference —
 I have to conc up for functions & constraints of functions, ~~the~~ (including "R" functions)
 How much they will slow down to of PC's of solns, is unclear. I may (probably) have
 to try correct new means to 1 PC's of solns. Can I get a good enuf soln to be able
 to go on to phase 3 (use of Induction for updates — is soln of QA as a
 general optzn. problem, using only those 2 tricks: which are simply ways to do finite & prod.
 cost is pretty much same as to Binary QA problem — videtur.
 I will want to write about how both funds & constraints must be used to conc matter
 ∴ int. T & Q.

SN

N dim, bidirectional prefix sets: A Kraft inequality!

Consider all some subsets of pts on n dim variables $\vec{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n, x_i \in \mathbb{R}$

$x_i = 0$ or 1 . Cor. then have k & k elements alphabet). Consider $k=2$, here to start.

Associate with each subset of $\{1, \dots, n\}$ a subset of strings of which it is a substring — finite max total wt.
 (since each string wt. 2^{-2M})

If all we have a set of strings (M, n) & each pair of them has at least k
 in common at which they disagree, then they have no common associate strings. ∴
 so the total wt. of this set is ≤ 1 . ∴ wt. of each string is 2^{-S} , where S is
 the total no. of pts. specified on that string.

To a good proof (which needs to be "cleaned up") addresses well if $k > 2$.

This idea applies to any subset of subsets of points; To be a "prefix set",
 each pair of sets must differ at at least 1 pt. So every pair of the bunch must have at
 least 1 pt. in common, so they can disagree about at least 1 pt. If $k=2$ is never, then we
 can assign a semi-measure to the "bunch". $\sum p_i \leq 1$, $p_i = 2^{-\text{no. of pts in subset}}$.

Other than $\sum p_i = 1$; How can we tell if a prefix set is "complete"?

Q if we have a prefix set on pts, then for any subset of pts: ~~is it true that~~
~~every~~ subset we can express as a sum of prefix sets in only 1 way? I suspect that the B. R. may
 may not be so simple!

SN

Where we are now: 1) 798.00 - 799.10 base to beginning of a new method. Needs to
 be finished.

2) 720.10 - 722.15 (APPX B) - rewrite beginning of exam just what μ is (A/R)

3) SOP 796.26 - 40 for review of what needs to be done.

808
809
810

806.40 **SN** It would ~~seem~~ that "in Recog"; If we have the complete search set, that is AZ141 "Sumac" should be able to discover any "PC increasing" trends that are desirable.

Could ~~the~~ Sumac use Lsrch to discover "Lsrch"? As is, probably not: T. have that I would use to desc. Lsrch, would first be to write that $\frac{CC}{PC}$ is some measure of how long it should take to solve a problem. — in this case CC would be some kind of max CC off trials. But this CC should be to CC of the "soln. trial" is a surprise!

Remember: Lsrch is ~~the~~ a base, if "All in do is in P.D.": But a AZ141 search is not Lsrch. — it just puts soln trials in L cost order, but there need be no heuristic. in to in "guiding PC" (which is fix — is it. @prpl). As a Sumac, hvr. So most coding of the post may take in form of a heur. (here "Context").

In evaluating
t. rate success of
Person's paper
on found
normal (N.B.)
to induction:
T. Q is: What
is the reward
to Jacobson?
This gives t
amount of
"independence"
from me same.
Space to another

Rewriting of "Appendix (B)": 720.10 - 722.15

At only "proof" I'll do is 720.08-10, 721.09-14; (Give 578, T3 corr, & Hutter's Genz, to reader 3.)
I also need to show that $\frac{P}{u} > 20^m$, i. a prop of u w/ a same rate of P usage.

So this only in (page plus maybe 722.00-.15 a/o other discussion of significance) → 808.00

19

20: 720.07 **SN** Thinking about ~~the~~ this again: we actually do predict t. Az sequentially. I want to show $\frac{P}{u} > 20^m$, i. a prop of u . (Actually, I find I should this in t. 599; "2 hundred ...")

It's in § 3.3 of "2 kinds of Prob. - induction" (599). I could just refer to that & then say that it shows $\frac{P}{u} > 20^m$. ∴ I want to show that the analog of 599 eq (3) is true.

Maybe I should define ALP for QA on the basis!

En 599 eq (2) does indeed imply eq (5) (i.e. 8) — but how to get from there to predictions?

It would be easier to prove the convergence ~~to~~ Run for 599 type eqs: in (i.e. 8) eqs (1) & (5).

So in 599: Just do PC's of t. D_n

$(P_0 \equiv P_{D_0})$ we want to compare

$P_0(D_0) = 1$. $D_0 = Null = \Lambda$
So it is TRIVIAL!

$$\prod_{n=1}^k \frac{P_u(D_1 \dots D_n)}{P_u(D_1 \dots D_{n-1})} = \frac{P_u(D_1 \dots D_n)}{P_u(D_0)} = P_u(D_1 \dots D_n)$$

$$\prod_{n=1}^k P_{PC}(D_1 \dots D_n) = P(D_1 \dots D_n) = \prod_{n=1}^k P(D_n)$$

29
30

808.00-.03
809.18-.40
810.00-.29
810.34-.30

Appendix 2

or equation (1)

close up. no space. 809.30 looks ok!

spec. (907.19) Start Appendix 2 again!

Appendix 2:

We will show that for an adequate sequence of (Q_i, A_i) pairs, the predictions obtained by the probability distribution of $q_i(1)$ will usually be good.

P/2m: Express f, A_i seq as $Z = (z_1, \dots, z_m)$. Then show $q_i(1)$ is (and is) conditional

conditional probs for each bit & space - conditional on all previous info in Z .

So P assigns cond. probs to each symbol of Z ; as does U .

Then show that \prod product of these cond. probs for $P > 2^{-2m}$ product of cond. probs for U .

Invoke Solovay-Kempner Theorem. SFH Theorem: SGT

Perhaps state SFH Theorem in direct history as a footnote.

No! Let $P(A_i^j)$ be the conditional probability of the j th symbol of A_i . A_i has $k_i - 1$ symbols. $A_i^{k_i}$ is the special symbol S . Nearly/less!

No! Then $\sum_{i=1}^m \sum_{j=1}^{k_i} E (P(A_i^j) - U(A_i^j))^2 < \epsilon$

$P_{i,j}(1)$ is the probability that the j th symbol of A_i is 1, conditional on all previous symbols in A_1, \dots, A_{i-1} and the A_i 's before it.

A_i has k_i symbols; the k_i th symbol of A_i is always the special symbol S (space). The zeroth symbol of A_i is Λ (null).

then $E \sum_{i=1}^m \sum_{j=0}^{k_i} (P_{i,j}(0) - M_{i,j}(0))^2 + (P_{i,j}(1) - M_{i,j}(1))^2 + (P_{i,j}(S) - M_{i,j}(S))^2 < \epsilon \ln 2$

The expected value being with respect to U

Instead of this 'E' I'd like to make clear just what E is. Hvr. look at 78 TB corr! Turay E is done makes .23 look very wrong!

Is .20 to sum of terms like $U_{i,j}(0) (P_{i,j}(0) - M_{i,j}(0))^2$?

Probably I could write as $\sum_{l=0,1,S} \dots$

More generally, let A_i 's have a larger alphabet than 0-1. The first Σ will mix and match the symbols plus the special symbol.

I guess the main trouble with the expression .20 w.o. 'E' is that the 'Z' sequence ranges

over all possible binary strings: $\sum P(Z) \dots$ is correct.

perhaps write .20, then express it in $R_{i,j}$ form. $m = \sum_{i=1}^m (k_i + 1)$

Z_l is seq. of m symbols. Row no 3^l of Row $l = 1 \dots 3^l$.

So in .20 insert $M(Z_l)$ is sum over $l = 1 \dots 3^l$.

00: (808.03) spec.

02: ONWALT

To do P12, we ~~express~~ express the sequence of A_i 's in $\mathcal{A}^{\mathbb{N}}$, and the spaces between them, as a ternary sequence, $Z \in \Sigma^{\mathbb{N}}$. We ~~will~~ hypothesize a very good

probabilistic question answering algorithm, μ and, using Hutter's generalization of Rice's theorem of §78, show that the ~~algorithm~~ symbol probabilities assigned by the algorithm of eq(1) must ~~converge~~ converge rapidly ~~to~~ toward one another.

[5P] ~~Let~~

Let μ be a ~~very poor~~ predictor of the A_i 's; $\mathcal{A}^{\mathbb{N}}$ will also "converge to μ ".
Perhaps sum over any first ≥ 2 bits have just n "s" symbols. If a seq doesn't contain n "s" it $\epsilon=0$ from here μ is μ . It is easy to compute how many

Legal (n -s's) $\leq \sum_{i=0}^n \binom{n}{i} 2^{n-i}$ (total no. of bits out, i) \times no. of ways you can place n s's in n places

[5N]

If μ is a very poor predictor of the A_i 's; $\mathcal{A}^{\mathbb{N}}$ will also "converge to μ ".
I think we also assume that μ is the "creator" of the ~~actual~~ A_i sequence that was actually observed.

So .02 doesn't ~~really~~ really ~~depend~~ depend!

18: 808.03

To do P12, we hypothesize that the actual ~~the~~ sequence of A_i answers $\mu(A/Q)$

that have been observed, were created by a probabilistic algorithm, μ and

that μ can be described by ~~our reference machine~~ ~~with~~ k bits, using the reference machine that assigns ~~appropriate~~ probabilities to all ~~partial recursive~~ functions.

Any probability distribution that assigns probabilities to every possible A_i , must also assign probabilities to each bit of A_i . Suppose that S_r is a string of the first r bits of A_i . Then the probability given by μ that the $(r+1)$ th bit of A_i is ϕ is

$$\sum_j \mu(S_r \phi X^j) / \sum_j \mu(S_r X^j)$$

X^j ranges over all finite strings. which we will call P_j . (italics, v.c.)

Similarly, the algorithm of eq(1) can be used to assign a probability to every bit of every A_i .

We will represent the sequence of A_i 's by a string of these A_i 's separated by the ~~space~~ symbols, s - denoting "space". Z , then, is a sequence of symbols from the ternary alphabet $\{0, 1, s\}$. Both μ and P are predicted algorithms for the symbols of Z . Using an argument similar to the foregoing, it is clear that both μ and P are able to assign probabilities to the ~~space~~ space symbol s , as well as to ϕ and 1 .

See how this comes out (BAD)

Hutter, M. Optimality of Universal Bayesian Sequence Prediction for General Losses and Alphabets <http://www.idsia.ch/~hutter/mhutter.htm>

20: 8:40: We have, then, two probability distributions on ~~the~~ binary strings Z . The first distribution, μ is the creator of the observed sequence, and the second distribution, P represents eq (1) ~~is~~ and tries to predict the symbols of Z .

For two such probability distributions on ~~the~~ binary strings, P is a corollary of Theorem 3 of Sol 28 applies: The expected value, with respect to μ (the generator) of the sum of the differences in probabilities assigned by μ and P , is less than $\frac{1}{2} \ln 2 \cdot k$; ~~where k is the number of bits in the description of μ~~ ~~such that $P > c \cdot \mu$ for all arguments~~ ~~all arguments~~

Hutter (1) has generalized this corollary ~~to~~ beyond the binary alphabet. So it applies to systems with any finite alphabet. We are here concerned only with ~~the~~ binary alphabet, but it implies a corresponding theorem if the symbols of the A_i are from any finite alphabet.

33 ~~Superscript.~~ $\sum_{i=1}^n \sum_{j=1}^{h_i+1} (\mu_{i,j}^R(t) - P_{i,j}^R(t))^2 < k \ln 2$

A_i is A of Z assigned by P

$P_{i,j}^R(t)$ is the probability that the j th symbol of A_i will be t , conditional on previous symbols of A_i as well as previous A_j 's in the sequence Z and the corresponding Q 's.

t takes the values 0, 1 and S

$\mu_{i,j}^R(t)$ is defined similarly to $P_{i,j}^R(t)$, but it is independent of previous A_j 's in the sequence.

h_i is the number of bits in A_i . The (h_i+1) th symbol of A_i is always S .

The total number of symbols in Z is $\sum_{i=1}^n (h_i+1)$

0, 1, S

no 2

\rightarrow this is a, c, S

~~The total number of symbols in Z is $\sum_{i=1}^n (h_i+1)$~~

$\rightarrow P \rightarrow Q$ sums over all strings Z that consist of n finite binary strings separated by S symbols (spaces). no (S)

$\rightarrow \mu(Z)$ is the probability that μ assigns to Z , in view of the ignorance of Q 's. k is the ~~length of the~~ (shortest) description of μ .

29 This implies, that the expected value with respect to μ of the ~~error~~ ~~error~~ squared "error" between P and μ , ~~summed over~~ summed over the individual symbols of all of the A_i , will be less than $k \cdot \ln 2$. $\rightarrow (34)$

30 SN eq 13 is a bit funny because it sums over ~~the~~ Z 's of different lengths. This makes no difference in ~~the~~ result: $T \cdot Q$ is 1 is to "proof" ok.

The theorem has been very ~~far~~ far μ 's that generate strings of length n .

but

For essential proof use method of Gacs on Li-Vit. But he used P r. the same measure convergence theorem.

34 $\rightarrow 29$ Since the total number of symbols in all of the answers can be very large for even a small number of questions, it is clear that the error per symbol decreases rapidly as ~~the~~ n , the number of ~~the~~ QA pairs increases.

Gaul. Theory of Problem Solving

mm

00: 8/10.40 : **SM** A General Theory of Problem Solving:

Steady state operation: Old Ideas $\equiv OI$; New Ideas $\equiv NI$,

• New ideas are formed by combining old ideas in ways that solve problems.

Methods of combining ideas are in themselves, ideas; so we can do all kinds of combining by simple composition of ideas

To solve a problem, we have to have not only the necessary old ideas to be combined to solve the problem, we have to have good ideas on how to narrow down the search for an adequate combination of old ideas. ^{so it takes place w. as least as possible.} So a good way to solve problems would be to have a idea that looks at the problem and looks at available old ideas and is able to map out a search strategy to find ~~the~~ a new idea (or old idea) that solves the problem. ~~This~~ Ideas of this sort, I called PST's.

10

Some ways to forgo might work: L such is often a v.g. search method in some cases \approx best. One good auxiliary PST, might be for a PST to look at the problem's given a PD on ^{old} ideas that could be combined to solve the problem.

I really have to work this out in more detail: it is essentially the way TM is supposed to work!

20

A basic idea: Best is "All into is m.f.P. from L such is \approx best".
Hrr, the PD. must include many PST's, search functions, etc.
So if we use ~~that~~ Normal L such: we want an algm that looks at the problem, and the set of available ideas, it implements a search strategy. — I think this amounts to a PST.

perhaps every time a problem (or sub problem) is solved, ~~all~~ all relevant PST's are "updated". T. many of "update" may vary betw. PST's.

30

31.

SN IMPT is perhaps related to forgoing: We write a short, elementary ^{in not a long time} TSD, that we expect should be learnable by a human (we may be able to test this empirically!). We then have to discover ~~the~~ a set of hours that are adequate for this performance. If we can't do this: we don't really adequately understand the performance of .31 (Note 8/16.00ff on "Scaling")

799.00-40, / 812.00-
799.00-10

00: 799, 10 : By modifying the guiding probability distribution, in view of previous experience, the solutions to future problems are given higher probabilities

~~Since the time to find a solution by Lsearch is proportional to the time to generate and verify a correct solution, divided by the probability assigned to that solution, it is clear that the performance modification of the probability distribution, will decrease the search time needed for a solution.~~

FF ~~Does this technique modify the probability of solution enough to permit the solution of difficult problems in acceptable times?~~ We will be using several methods to modify our probability distribution. The adequacy of these techniques is discussed in section 3 on "Scaling." Not really!

In our initial training of the P, A induction system, we used a particularly simple form of Lsearch that is appropriate to induction problems. We search for production functions in order of their a priori probability divided by the time needed to generate and test the trial. This search is "Blind" in the sense that it contains no information on the nature of the induction problem to be solved.

omit ~~The search algorithm may know about the past, but it does not know what it is looking for - except that the trial must satisfy a certain test. Superficially, this would seem to be the same as to inversion problems - in which we are searching for a string that satisfies certain conditions. The difference between the two problems is that in inversion, we want only a string that satisfies the conditions and we want to find it as rapidly as possible. In the case of induction we want to find a function (represented by a string) that satisfies certain conditions, but we want the description of the~~

~~This search algorithm is "Blind" in the sense that its trial functions depend only on past data and not at all on the target. And that it is trying to generate. As the system has had enough experience with certain kinds of induction, we will be able~~

to employ and update more effective search algorithms for inversion problems (which are a kind of time limited optimization problem) and induction. These are not "Blind" algorithms.

799.40 seems more relevant

780.40 ← ?? seems irrelevant very unlikely!

00:

812:?

The system described is ~~not~~ ^{this far} ~~not~~ ^{pure} inductive system, ^{and} ~~no~~ ^{an} exception to work some problems of some difficulty, it has a serious deficiency — it uses a "Blind" search algorithm. When it is ^{looking} ~~searching~~ for a prediction function ~~to~~ ^{to} fit the data, it's trial functions depend only on summaries of past data, and not on the target. An ~~inductive~~ ^{inductive} inference problem ~~is~~ ^{is} characterized by ~~the fact~~ ^{that} ~~that~~ ^{it is characteristic of} that they ~~are~~ ^{are} often solvable by ~~this~~ ^{this} kind of search.

We will next ~~describe~~ ^{describe} two kinds of problems in which blind search is usually ~~not~~ ^{inappropriate}.

10

We can, however, do much better. After our induction system has been applied to the updating of search algorithms for more general problems (inversion and time limited optimization problems), it will be possible to use it to update more general (non-blind) algorithms ~~for~~ ^{for} search for good induction functions. One ~~motivation~~ ^{motivation} for working inversion and optimization problems is the expected improvement in the general induction system, obtained by sharing concepts used in updating. The induction system itself can be viewed as a kind of "Bootstrap Dance" in which any improvements in induction translate into better updating and more improvements in induction.

20

30

00: 816.40 : Section 2 : ~~Induction~~ Inversion Problems.

In Inversion problems, we are given ~~a function $G(\cdot)$~~ that a string S , and a function $G(\cdot)$ from strings to strings and we are asked to invert G with respect to S : to find an x such that $G(x) = S$. We want to ~~do this as fast~~ as possible.

One way to solve such problems is to use Lsearch. (Appendix C describes different ~~types of~~ kinds of Lsearch).

In one simple kind of Lsearch, we look for a function, F_2 , that maps the description of $G(\cdot)$, and the target string, S , into a ~~trial~~ string, x_2 .

$$x_2 = F_2(\tilde{G}, S)$$

\tilde{G} is a string that describes $G(\cdot)$. It is usually a program. We can use the language \tilde{A} of appendix A, to obtain a probability distribution, on all possible $F_2(\cdot, \cdot)$.

We have an a priori probability distribution $P_0(\cdot)$ on descriptions, \tilde{F}_2 , of possible $F_2(\cdot, \cdot)$ functions.

From $P_0(\cdot)$ and $G(\cdot)$ and S , we can use Lsearch to obtain a solution, $x_2 = F_2(\tilde{G}, S)$ such that $G(x_2) = S$. If it takes ~~more than~~ a total time t_0 .

to generate F_2 from \tilde{F}_2 , generate x_2 from $F_2(\tilde{G}, S)$ and test $G(x_2) = S$.

Then Lsearch will take about time, $t_0 / P_0(\tilde{F}_2)$.

At first, the probability distribution P_0 will be obtained from a universal algorithm (say a Turing machine or other universal computer) able to simulate the function F_2 , using its description \tilde{F}_2 . One easy way to get P_0 is to make

$$P_0(\tilde{F}_2) = 2^{-l(\tilde{F}_2)}$$

Here $l(\tilde{F}_2)$ is the length in bits, of the description, \tilde{F}_2 . So time for Lsearch is about $t_0 \cdot 2^{l(\tilde{F}_2)}$.

After we have solved the first problem, we will put the definition of F_2 in the "parameter list" of the language used to generate trials.

When we have solved several problems in this way, the language will have several useful functions in its parameter list. It will usually be possible to find common sub functions that can be defined, to increase the probabilities of solutions past solutions to problems, just as was done in the solutions to the induction problems of section 1. In general, all of the techniques for increasing the probabilities of solutions of induction problem in section 1, can also be used for the local prog to solve inversion problems.

Suppose we have ~~an~~ a new problem, (G_n, S_n) and we want the ~~best~~ probability that an arbitrary PST, $F_0(\cdot, \cdot)$ will solve this problem faster than any other PST considered. ~~We could~~ We could then use such a distribution ~~to~~ a very good oriented search.

How can we obtain ~~such~~ ~~a~~ distribution, $P(F_0)$? ~~We could~~ ~~use~~ ~~such~~ ~~a~~ ~~distribution~~ ~~as~~ ~~follows~~:

In section (1) we wanted solutions of maximum probability. For Inversion

00: 815.90

In section (1) problems, since a search time is ~~inversely~~ inversely proportional to probability of solution, ~~the~~ ~~higher~~ ~~probability~~ ~~solutions~~ ~~will~~ ~~lead~~ ~~to~~ ~~error~~ ~~us~~ ~~minimal~~ ~~search~~ ~~times~~ ~~and~~ ~~locking~~ ~~for~~.

Each of the $F_0(\cdot, \cdot)$ that has ~~successfully~~ ~~solved~~ ~~at~~ ~~least~~ ~~our~~ ~~problem~~, can be regarded as a ~~new~~ ~~problem~~ ~~solving~~ ~~technique~~ (PST).

In the foregoing the method of solving problems ~~just~~ just described, ~~it~~ ~~simply~~ ~~tries~~ ~~to~~ ~~imitate~~ ~~the~~ ~~sub~~ ~~functions~~ ~~of~~ ~~the~~ ~~PST's~~ ~~that~~ ~~have~~ ~~been~~ ~~successful~~ ~~in~~ ~~the~~ ~~past~~. The system does not "know" that it is supposed to

10

try to ~~find~~ ~~solutions~~ ~~that~~ ~~have~~ ~~been~~ ~~successful~~ ~~in~~ ~~the~~ ~~past~~. Though this method works ~~well~~ ~~to~~ ~~some~~ ~~extent~~, it can be much improved.

See if this works

Consider the set of quadruples: $(\tilde{G}_j, S_j, F_j(\cdot, \cdot), t_j)$

\tilde{G}_j is a string that describes the j th function to be inverted

S_j is the argument for F_j in case of \tilde{G}_j

$F_j(\cdot, \cdot)$ is the function of \tilde{G}_j and S_j that solved the problem

t_j is the time it took ~~to~~ ~~solve~~ ~~the~~ ~~problem~~ ~~with~~ ~~the~~ ~~function~~ ~~F_j~~ ~~for~~ ~~$F_j(\tilde{G}_j, S_j)$~~ ~~to~~ ~~generate~~ ~~the~~ ~~correct~~ ~~solution~~ ~~of~~ ~~the~~ ~~j th~~ ~~problem~~.

We have a quadruple like this for each problem solved ~~in~~ ~~the~~ ~~past~~.

L.C. And

See if this works

20

Suppose we have a new problem, (G_n, S_n) given any $F_k(\cdot, \cdot)$ ~~the~~ ~~previous~~ ~~section~~

We can use the technique of ~~the~~ ~~previous~~ ~~section~~ or induction, to get a probability

mean

density distribution of the time it will take for F_k to solve ~~the~~ ~~problem~~. That problem

Here we regard ~~the~~ ~~quadruple~~ ~~$(\tilde{G}_j, S_j, F_j(\cdot, \cdot), t_j)$~~ as Q_j and t_j as A_j .

Q_n is ~~the~~ ~~quadruple~~ ~~$(\tilde{G}_n, S_n, F_n(\cdot, \cdot), t_n)$~~ . As output A_n , we get the desired probability

density distributed on the solution time ~~of~~ ~~the~~ ~~problem~~. F_n may or may not be in the set of ~~the~~ ~~functions~~ ~~that~~ ~~have~~ ~~been~~ ~~tried~~.

From the foregoing, for arbitrary $F_0(\cdot, \cdot)$ and problem (G_n, S_n) we can obtain a probability distribution

$h_0(t)$. This is the probability that $F_0(\cdot, \cdot)$ will solve the problem in ~~less~~ ~~than~~ ~~or~~ ~~equal~~ ~~to~~ ~~time~~ ~~less~~ ~~than~~ ~~or~~ ~~equal~~ ~~to~~ ~~t~~ .

30

Thorp 3

Usually there will not be enough data to get a very detailed ~~of~~ ~~the~~ ~~distribution~~ $h_0(t)$.

A ~~parameter~~ ~~might~~ ~~be~~ ~~of~~ ~~the~~ ~~form~~:

34

$$h_0(t) = \int_0^t \dots N(x, \mu, \sigma) dx$$

N is a Normal Gaussian distribution on x of mean, μ and variance, σ^2 .

\int_0^t is the probability that F_0 will ever (in $t < \infty$) solve (G_n, S_n) .

error power

Often you will only have ~~partial~~ ~~data~~ ~~and~~ ~~you~~ ~~will~~ ~~often~~ ~~demand~~ ~~that~~ ~~we~~ ~~use~~ ~~less~~ ~~than~~ ~~3~~ ~~parameters~~.

on G_n, S_n

20: 817.40 Given a set $[F_e]$ of PST's and their associated $h_e(t)$ distributions for solving G_n, S_n we can obtain a probability that each F_e will get the fastest solution (smallest t) for the problem. This probability, P_e is

03
$$P_e = \int_0^\infty h'_e(t) \prod_{j \neq e} h_j(t) dt = \int_0^\infty \frac{h'_e(t)}{h_e(t)} \prod_{j \neq e} h_j(t) dt.$$

While this integral looks difficult (time consuming) to evaluate, not much accuracy is needed and it is relatively simple to ~~compute~~ calculate, using a Monte Carlo technique. We repeatedly sample from each h_j distribution and pick the ~~one~~ value that is largest. The frequency with which h_j is largest will be an estimate of

10 P_e . Since little accuracy is needed, we can use very simple approximations to model each h_e .

~~We can use the following~~ If the values of the P_e are initially expressed in Monte Carlo form, we can use that form to guide a Monte Carlo search for a solution to the problem. ~~Appendix~~ describes a Monte Carlo version of search.

17 Needs revision ~~The foregoing~~ ^{always} may not be a practical way to solve Inversion problems: The time needed to compute the P_e 's of the relevant F_e 's ^{can} be much larger than the time needed to ~~test~~ test each F_e empirically!

20 omit In this case, we should ~~use~~ use the foregoing analysis as a guide to approximating what we want. The clear superiority of this method is that the system ^{already} knows that it must try to find a F_e that obtains a fast solution. In this case ^{the present} model should be regarded as a kind of ideal that we try to approach as closely as we can, in ~~whatever~~ whatever time we have available.

10 See 819.13-15 for unobscured contents In this case, we ^{most find} will use different methods of solution — but the model described gives the ^{goal} goal of such methods — makes ~~possible~~ possible for us to improve them. A clear advantage of the system described, is that it is not a "blind" search — the system ^{already} knows that it must find a F_e that obtains ~~as fast~~ as fast a solution as possible.

~~to section~~ (maybe section 2.2;) They lost B was somewhat modified — sounds much better now. Also, in view of criticism 823.25-31. to ~~the~~ discussion of approx. method of 820.15-31.30 ~~was~~ was not included; 821.09-24 was detected from Report.

"new" ABCDefg / "use" ABCDefg

Subsection to section 2 on Inversion Problems.

00:18.40: Problem Solving Techniques. The first attempt to map problem descriptions into candidate solutions may be regarded as "Problem Solving Techniques" (PST's). Initially, they will be put into the system by the trainer. ~~some~~ number of such PST's will be inserted into the system by the trainer. obtained by the system using the "solutions" to suitably simple problems - obtained by Lsearch over PST space, using a universal distribution on functions - (each described in Appendix A). This will give us a set of PST's "known to be (at least occasionally) useful". We can then apply the technique of eq() or its approximations, to obtain a PC distribution for guiding Lsearch to any new problem - but our set of PST's would be limited to those that ~~are~~ → 921.00

10: Just as we do not expect an exact solution for equation (1), we don't expect an exact evaluation of eq (2)(2.05). We simply want to get as good an approximation as possible in the ~~time~~ available time.

12: Eq. (2) In such cases, eq (2)(2.03) ~~should~~ be regarded as a description of where we want to go rather than "How to get there". Approximation techniques should be guided by the goals and subgoals implied in this discussion.

15: **SN** to idea of "A well defined Induction problem", corpus is a prod program, it is clear what is to be optimized.

- 20: Some items I may want to include in either "T. report" or **Augmented Report** (My Version)
 - 1) How to do time series' 2) T. MXT corpus problem & f. Mad corpus Perm.
 - 3) Much detail on 2) T. 3 Goals. House Perm is Just why is to what extent "all info is in PD" enables Lsearch to be an optimum.
 - 4) Exactly what factors mean in "General Solns to problems" I.e. devoting 1/2 time to Lsearch, 1/2 time to updating GPD.
 - 5) The mechanics of optzn of continuous & discrete param of a system.
 - 6) Consideration of 1. set of all PST's: Obtaining a ~~single~~ ^{PC} best PST ~~each~~ each one will be (best / best) for a particular problem. I've recent written much on this recently, I may have a "Review".
 - 7) A "Protocol" Explains that this is a "Report", not a "paper": Also. ~~that~~ Abstract gives: General picture of an advanced system: a strong subgoal.
- 30: Better how can do this, we will implement a (limited kind of system that works QA, induction ⁱⁿ only). These kind of system that Algorithm it was only able to solve difficult problems, but of considerable power, it has the a certain deficiency its use of "BLIND" search algorithm it is a kind of "Blind search" ~~and~~ necessary weakness of "Blind search".

Tried:

- Abstract (General Learning - USA (Start - now) Knowledge for Gammer)
- Intro B is # Total in the.
- Section 1 | Induction Differences Sec 1.00X
- on prob 1/17/02
- 11/10/02
- Do better "Abstract" - Do it right.

Take faster pointer!

00: 8:19.40! For the more advanced system, ~~we overcome this weakness~~ uses a somewhat different search algorithm that overcomes this deficiency.

- 1) Simple Inductive Search: Looks for function \rightarrow 4:18:15 Mark ^{searches for} ~~to find~~ least order.
- 2) Next: Finds local at problem desc & produces soln. ~~we~~ search on each function.
- 3) For funcs of 2) we learn which funcs are likely to solve a given problem rapidly.
- 4) We apply 1. method of 3 to yr "library" ^{problem} ~~part of methods~~, (Bootstrap): Self-improvement.

8) Another item: Just $x^2 + 2bx + c$ or $\frac{b \pm \sqrt{b^2 - 4c}}{2} = \frac{b}{2} \pm \sqrt{\left(\frac{b}{2}\right)^2 - c}$

If it took ~~to~~ ^{loops} per trial, we have to ~~run~~ $n > 0$ trials ~~to~~ $n \times 10^6$ ~~to~~ n soln. in 1 second. Got pc of ~~it~~ ^{search}

This ~~results~~ ^{results} in a way of looking at ~~least~~ ^{search} ~~invariant~~ ^{show how additional complexity} ("a conf box \uparrow m yrs to wait")

9) Discuss how we get "logical reasoning" into System. Mark in section on TSO's & Scaling.



15 For end of first subsection of section 2: One ~~way~~ ^{way} to get approximate P_L distributions much more rapidly: After we have obtained P_L distributions ~~via~~ ^{via} $g(g, o_3)$ for several problems, over ~~several~~ ^{using} ~~PST's~~ ^(8:18.03) one or more PST sets, we can use the inductive system of section 1 to ~~verify~~ ^{verify} probabilistically & extrapolate these results.

A final result of ~~this~~ ^{this} induction: $W(\text{problem desc.})$ output is ~~satisfy~~ ^{satisfy} P_L pairs in $\approx P_L$ order (highest P_L first) with P_L being probability that PST's \rightarrow "Back" of first of P_L PST's for that problem.

what ~~is~~ ^{is} ~~gets~~ ^{gets} into report, it is partly an Essential part of System

I don't have time now to develop this adequately, but maybe one line in ~~the~~ ^{the} previous subsection of § 2:

29 30 One ~~way~~ ^{way} to get ~~the~~ ^{approximate} approximations of P_L , is to use ~~an~~ ^{an} inductive system ~~at~~ ^{at} § 1 to obtain a fast function that ~~relates~~ ^{assigns} P_L values to the $F_L(\cdot, \cdot)$ for a given problem description. ~~that~~ ^{that} ~~probability~~ ^{probability} ~~that~~ ^{that} ~~the~~ ^{the} P_L of $F_L(\cdot, \cdot)$ is to calculate the P_L 's ~~for~~ ^{for} a set of problems, using $g(g, o_3)$, then use the inductive technique of section 1, to extrapolate this relation to new problems ~~and~~ ^{and} possibly new $F_L(\cdot, \cdot)$ sets. The accuracy of the assignment will, as in all inductive problems, depend much on the size of the data set, and the amount of time spent searching for a solution.

[mentioned $O^i(A|Q)$ ~~is~~ ^{is} ~~one~~ ^{one} form of P_L . another: $W^i(Q) \rightarrow$ list of A, P pairs in P_L order.

GOOD!

SPEC
59

00: 820. ~~10~~: We have 23 problem solutions. We can expand this set of PST's in several ways:
One way is to have ~~a~~ a set of PST's inserted by the trainer — who thinks they should be good. Another way to extrapolate

change $f \rightarrow G.()$

for end of first part of section 2. u.c.

09: ~~10~~: One way to get approximate P_e distributions ~~more~~ rapidly:
After we have obtained the P_e distributions via eq (818.03) for several problems over several F_e functions, we can use the inductive system of section 1 to probabilistically extrapolate ~~these~~ results to new problems and possibly a new set of F_e functions.

yes

~~The~~ output of the inductive system
output distribution
in the present case

Can express ~~probabilistic~~ ~~functions~~ in several ways:

No omit

We regard the problem as a "Q",

yes

each and/or possible F_e functions as ~~a~~ a possible "A". The induction ~~function~~

takes the problem description as inputs and lists the possible F_e

No omit

functions as outputs, each with its associated P_e — in approximate P_e order

22: The system can express its output distribution in several different ways: ~~In the present case~~ a good way would be a listing of each of the F_e , paired with their probabilities with the pairs ~~in~~ roughly probability order.

yes

Actually: 22 is usually unnecessary unless there are an enormous no. of F_e 's,

We try all of them each $T \leftarrow 2T$ round, so the ordering which P_e for F_e is listed is unimpl.

If the no. of F_e is large ($\rightarrow \infty$) then we don't do all of the F_e each round: ~~we~~

In ~~the~~ we only do F_e 's that have $p_e > 2^{-n}$.

30

See 823.25 - .31 for Serious Criticism of .09-.24

(Section 3) Time Limited Optimization.

use PST concept: it was introduced in § 2.

20:

~~Time Limited Optimization problems are solved in a way very similar to that for~~

In the previous section, we described ~~two~~ ^{two} ways to solve induction problems.

The first way was in the first, the guiding probability distribution is directly modified by previous solutions to inversion problems. In the second, we use the induction technique of section 1 to obtain the probability that ~~each~~ each candidate function will give the fastest solution to the problem, and we use this probability distribution to guide ~~the~~ L search.

10

While the first ~~graphical~~ ^{analytical} method ~~is~~ ^{is} applicable to the solution of only time limited optimization problems, the second method ~~can~~ ^{can} be used for ~~more~~ ^{more} general time varying optimization problems.

We will first describe the solution of time limited optimization problems.

Then show how time varying optimization problems are solved.

In time limited optimization, we have a function from strings to reals $G(\cdot)$, and we have time limit t . We must find an X in time $\leq t$ such that $G(X)$ is as large as possible. The string X can represent ~~more~~ ^{more} generally, the string X , can represent a number.

These problems are solved in much like the inversion problems of section 2.

We have a set of PSTs, F_i , that look at the problem $(G(\cdot), t)$ and return after a requested time, t_i , a solution candidate, $X_i = F_i(G(\cdot), t_i)$, which is evaluated as $G(X_i)$.

19

20

... to test each F_i empirically. In such cases \downarrow (eg. B18.30) and the reasoning leading to it, may be regarded as a description of "where we want to go", rather than "how to get there".

We must find different ~~different~~ ^{different} methods of solution.

Solution methods

25

In a ~~to~~ ^{to} ~~empirical~~ ^{empirical} inductive procedure. (S2.19-24) ~~where~~ ^{where} it would be easier to use empirical optimization (in fact, that's why I suggested by Approximation). Doing so would ~~save~~ ^{save} ~~some~~ ^{some} as t . Old method of doing L search! ~~write~~ ^{write} ~~the~~ ^{the} ~~of~~ ^{of} ~~correctly~~ ^{correctly}! Another point: On basis of ϵ system's inductivity v. ϵ , it will be better to use "Method 1" for L search rather than ~~Method 2~~ (eg. B18.03) to guide L search. So B18.03 can be only used if ~~its~~ ^{its} ~~rather~~ ^{rather} ~~than~~ ^{than} ~~direct~~ ^{direct} ~~empirical~~ ^{empirical} ~~testing~~ ^{testing} ~~and~~ ^{and} ~~inductive~~ ^{inductive} system is already ~~very~~ ^{very} ~~good~~ ^{good}. "Method 1" may be regarded as one way to ~~try~~ ^{try} to approximate Method 2.

30

31

SN T. Soy effect is similar to + "FDA effect" in biasing results by selecting out an empirical "best": - may call soy \rightarrow FDA \odot . Also $\sigma \rightarrow \sigma^2 \frac{N}{N-k}$ effect. Also SMA also maybe "The 1024 10bit strings are seen"; 1 of 2¹⁰ 10bit strings is n't there any time for coin flips. Also + "Prayer Hospital" experiment(s).

Any way I MUST think out 22-31 carefully. I'm not sure its exact rate. - I really

need to know just what kinds of approx. Guiding f.d. ~~functions~~ ^{functions} are logic/usable.

823.19: Initially, we have a guiding probability distribution, $P_0(F_i)$, that is the same for all optimization problems. ← to given in Appendix A

$P(F_i) \approx 2^{-l(F_i)}$: $l(F_i)$ being the length of the shortest program we ~~could~~ ^{have found} for the function F_i .

~~For each problem G_i, t_i , we will~~ Many optimization programs such as F_i make many trials, and get better and better results for longer run times. In the present case, we will run the PST, F_i for ^{at most} time $t \cdot P(F_i)$, and ~~then~~ use $G(\cdot)$ to test the result. Checking all of the F_i will take time $\leq t \cdot \sum P(F_i) = t \leq P(F_i)$.

Using $2^{-l(F_i)}$ will give approximate P values so that $\sum 2^{-l(F_i)} < 1$.

~~so we have not yet used up all of our time:~~

We then double our testing times to $2t \cdot P_0(F_i)$. ~~This is repeated doubling~~

and retest all of the F_i . This doubling is repeated until we run out of time.

We then select the x_i that has given the highest $G(x_i)$ value.

823.19 → We have, as a function for Inversion problems, an ^{a priori} probability distribution $P_0(F_i)$ that is the same for all ^{Time limited} optimization problems. We run the generator and test any of the F_i and x_i for time $t \cdot P_0(F_i)$. Since $\sum P_0(F_i) \leq 1$,

the process takes $\leq t$ time. The x_i of maximum $G(x_i)$ is then selected as output.

If there is remaining time, we can ~~double~~ ^{repeat} the tests using times $2 \cdot t \cdot P_0(F_i)$ — doubling and redoubling until all of the time is used up.

As with Inversion problems After we solve the first problem, we put the solution $\$AZ\$$ function, F_i into the parameter list of the language ~~used~~ ^{used} to generate the trials ~~as~~ ^{as} we did with Inversion problems. ~~The remarks about~~ All of the

methods used to improve the ~~prob~~ ^{with} guiding probability distribution for Inversion problems can be used ^{with} ~~with~~ ^{minimal} modification in the present case.

In the discussion of the improved P_R distribution, we are trying to maximize G rather than minimize t , so eq (817.34) becomes

$$h_2(G) = \int_{-\infty}^{\infty} 2 N(G, \mu, \sigma) dG$$

Eq (818.03) becomes

$$P_2 = \int_{-\infty}^{\infty} h_2(G) \prod_{j \neq 2} h_j(G) dG = \int_{-\infty}^{\infty} \frac{h_2(G)}{h_2(G)} \prod_j h_j(G) dG$$

The circled class are t 's only difference from the eqs in section 2 mainly $t \rightarrow G$

30

TSQ & Lench: (More on GA) (i.p.c. of $\frac{b+(a-b)x}{2a}$)

Deleted from Report that refers to Sec in report - In Aug Report
1) Method of Gauss Ps (82018-821,30)
2) Time Varying Optm.

Discussion of Time Varying Optm.

say $h_1(G, t)$ is available from $F_2(G(\cdot, \cdot), t, \dots)$

$(G_1(\cdot), t_1)$ copies available: first consider given $G(x) \cdot H(t)$

So far now $G_2(\cdot)$ we have via ATM: $h_2(G)$ as a function of t .

So $h_2(G, t) \rightarrow h_2(G \cdot H(t))$

For each G, t value, what's prob that some F_2 will have $G \cdot H < G \cdot H(t)$?

$G_1(x) \cdot H(t)$

for given t , we have $h_2(G, t)$ probably F_2 will have $G < G$

Perhaps group some ideas from recursive trees zero dimension

My impression: that I originally had a soln for time varying G problem - not carefully worked out hrr.

But later I noticed my uncertainty about it's worth it out in more detail (with less hrr).

But I'm not sure of the recursion balance! \rightarrow For more recent soln, see NIPS 176.02-27

Best thing to do is leave time varying G out of report - it would be a nice

Branch Riv to report, but I would better spend time on other aspects of an already much overloaded report. I may have mention "Time var G" in introduction.

It was only mentioned in the brief brief dem of "Section 2"

So I cut it out.

Appx Note in Aug report NIPS 176.02-11
15 min idea

To define a common type of Recursion:

Consider dependency of factorial function:

$F(k+1) = (k+1) \cdot f(x); f(0) = 1. R(x) = x f(x); f(0) = 1.$

$F(x) = x f(x-1) \Rightarrow F(x) = h(x, f(g(x)))$ $h(\cdot, \cdot); g(\cdot); a, b.$

one may define this first by allowing "a" to be a form symbol.

For now end of Appx A:

A common class of The factorial function is a member of a class of recursively definable functions. recursively definable by Rec expression:

$F(x) = x F(x-1); F(0) = 1$

It is a member of a class of recursively defined functions, by $F(x) = h(x, F(g(x))); F(a) = b.$

A function of this sort is definable by the list: $h(\cdot, \cdot), g(\cdot), a, b.$

More general classes of recursive functions may be defined in similar ways.

So to define
Abstract
Intro
Section 1
Sec 2
Sec 3 & 4
Appx A Az
Appx B cont
under Sec 3 & 4
Sec 1 TSQ's
Sec 9 generalities
Appx C Lench.
Preface
References

29

30

New line

For n out of $A_n \times A_n$

Appendix A1

25: 25.40!

01 The factorial function is recursively defined by the expression!

$$F(x) = x F(x-1); F(0) = 1$$

It is a member of a class of functions defined by

$$F(x) = h(x, F(g(x))); F(a) = b$$

A function of this sort can be defined by the list:

$$h(., .), g(.), a, b$$

Formalisms of this kind can be devised to define more general recursive functions.



Alternatively, if we allow definitions of functs to refer to f. funct being defined,

we can simply include "a, b" after such a definition.

We introduce a set of interrelated interdependent recursive functs by first listing f. functs being defined, then give the defining eqns, then to "boundary roads".

It may be noisy to us "=" as a primitive.

Thus, a good method of expressing rec. (or other) funct. defns, should be closely related to one or more heuristics for discovering functs of that sort.

On "toaster learning"

Suppose we have two training sequences $[Q_1^1, A_1^1]$ and $[Q_2^2, A_2^2]$, each of length, n .

for extrapolation. If we regard the sequences as completely unrelated,

We will extrapolate the first sequence using

equation
$$\sum_{j=1}^{n+1} \prod_{i=1}^j O^{j,1} (A_i^1 | Q_i^1) \quad 1)$$

equation
$$\sum_{j=1}^{n+1} \prod_{i=1}^j O^{j,2} (A_i^2 | Q_i^2) \quad 2)$$

We will try to find formulas such as $O^{j,1} / \text{and } O^{j,2}$ so that

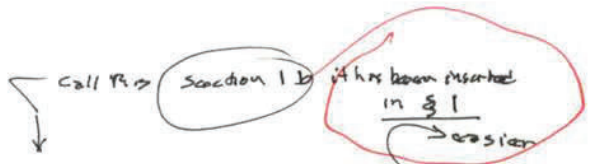
$$\sum_{i=1}^n \prod_{j=1}^i O^{j,1} (A_i^1 | Q_i^1)$$

and

$$\sum_{i=1}^n \prod_{j=1}^i O^{j,2} (A_i^2 | Q_i^2)$$

and individually as large as possible. In this case, each sequence extrapolates in its own way, and there is no possible transfer of learning between the two sequences.

two sequences



00: 827.40 : Continuation of end of § 1:

One value of P_{10} indicates "Transfer Learning". All of P_{10} data used in induction problems can be regarded as one large "Corpus" (body of data). At first it may be ~~too large~~ for the system to consider "sub-corpi" - data ~~sets~~ that have to some indication of their Q 's.

Each can be regarded as a separate training sequence, and they can be solved ~~individually~~ independently. It ~~is~~ is, however, better to regard all of P_{10} data as part of the same large corpus, so that the system can use ~~the same~~ common functions ~~in the inductions~~ on "sub-corpi" to ~~reduce the~~ reduce the ~~total~~ number of bits ~~to describe~~ to describe the data, thus materially increasing the value of ~~eq (4.1.18)~~ $eq(4.1.18)$.

This "sharing of functions" can be used ~~to link~~ to link any data in the system. What we have is an ~~effective~~ General condition Probability Distribution (GCPD) for all of the inductive data in the system. In sections 2 and 3 we will discuss ~~various~~ ~~probabilities~~ In various problems and Optimization problems. Though they are not induction problems, they have induction probability distributions that guide the searches for ~~the~~ solutions to these problems. These probability distributions will also be included in the GCPD.

to insertion § 2 (Inv probs):

After $P_2 \dots P_n$ all poss. $f_i(C, \dots)$ ~~of G, S~~

23 Suppose δ is the time to do ~~the~~ machine instructions: ~~for the P_i 's~~ we try all possible F_i 's, using a maximum time of δP_i to generate P_i , and x_1 and test x_2 . This takes a total ~~of~~ time $\sum \delta P_i$ - which is $\leq \delta$ since $\sum P_i \leq 1$.

If we don't find a solution we run all the trials again using 2δ as our time limit. We keep doubling the time limit and doing all trials until we find a solution.

10 The total time for all of the trials will be $< 2T_0/P_0$. P_0 is the probability assigned to the F_0 that solved the problem & T_0 is the time needed to generate and test the ~~correct~~ solution.

32

Section 3A (100-03)

Done → S3

Time limited

Since the induction problems of section 1 are optimization problems,

we can apply the techniques of the present section to their solutions,

In particular, we can use the ~~method~~ method that improves the P_2 distribution

to obtain ^{significantly} better performance.

Some things I need in to report:

1) How to merge and add PST's to system.

How TM can extract that set as a stack

factoring (partially factoring) the PST's.

in TSP Section 1
Grammar. How to merge can help by
The idea of a "Conditional Grammar".
in Section 1
TSP's

2) The idea of ~~Hyper~~ wrapper system (reg self replacement) in TSP

3) In §1 merge/division of L sub: refer to App C.

done!

4) §3 PL add a number of L sub

done

5) For §2 IIV probs:

with ~~828.23-32~~ replace 3rd P

replace 3rd P
time to generate P's

done

perhaps in TSP, discuss
GHT's.

All info in PD.

6) Discussion of GA - maybe in section

where it is first mentioned - ST → GATM

We could have section on GA.

7) Re update / search value fixed at 2. (by mentioned), - Discussion of exactly

what it "within a factor of 2" means. I could have a P or 2 on P's

- ~~Append~~ But why? "General Administration" - Gen remarks, Conclusions.

Section 57

8) The idea of SUMACS (Summary Machines). - Not done, but put in AR

9) Bibliography § Lev 72
4 Lit 93 (P)
5 Lit 97 (P) (Lisp-like lang) (NY 2001)
3 k02 95 Sec 1.

Sol 78

Sol 86

Sol 89

① Hute 01 - APXB

Get §1 in shape
Get rid of
Name author's
Get create
Put this on Entry
P's per.

10) ~~change~~ Appendix into w. names of new sections

00: 829: 40:

Stuba Appendix C: Levin's Universal Search Algorithm (Lsearch).

Levin's original paper (Lev 72) gave some properties of Lsearch, but didn't tell how to do it. The Li, Vitanyi book(s) (LiV93, LiV97) gave more details ~~and~~ ~~some~~ ~~ways~~ — including ^{some} ~~ways~~ ways to implement it. We will discuss three methods of implementation: Limit doubling, Parallel search and Random search.

~~13 of 251 vertices to L search method. For Random search read to Appendix C~~

Subsection 1 Limit Doubling & Lsearch

For Inversion problems: Given $G(\cdot), s$, to find x such that $G(x) = s$. We look for algorithms (PST's) $F_i(\cdot, \cdot)$ such that $F_i(G(\cdot), s) = X_i$. Then we test X_i to see if $G(X_i) = s$. ~~the time needed to generate F_i , generate X_i and test X_i~~ ^{in minimal time,}

The limit doubling method of Lsearch was described in section 2 (Inversion Problems) on Inversion Problems. ~~Section 2~~ described the application of ~~the~~ the limit doubling method to inversion problems. ~~It is a general method.~~ Section 3 described its application to time limited optimization.

Subsection 2 Parallel Lsearch

For Inversion problems: We work on all ^{PST trials} ~~trials~~ in parallel, time sharing our CPU cycles. The fraction of time ^{spent} on the construction of F_i , generation of X_i and testing X_i

is P_i , the ~~relative~~ a priori probability assigned to F_i . This method is at least twice as fast as the limit doubling method.

summary ~~with careful programming,~~ This method is about $-\log_2 P_i$ times as fast as the limit doubling method — P_i being the probability of the ~~best~~ P_i best found ~~the~~ successful X_i . However, it requires ~~much~~ much more memory and/or a clever disc swapping scheme ~~with~~

For Time Limited optimization, as with inversion problems, we work on all PST ~~trials~~ trials in parallel, using a ~~time share~~ time share of the fraction of P_i to PST, F_i . We quit when time has run out, and pick the PST that got the best output ~~at~~ at quitting time.

Random Lsearch is a variety of parallel search. ~~It is only twice as fast as the limit doubling method.~~ With probability P_i we randomly choose the PST, F_i , and we work on it for a fixed time T , ~~then~~ ^{ever} again ~~chosen~~ chosen, we ~~continue~~ continue where we left off, spending another T on it. T should be ~~small~~ small, but "large" with respect to time needed to switch from one PST to another.

10

20

30

00: 830.40: The properties of random Lsearch are ~~about~~ ^{about} the same as those of parallel Lsearch.

~~Random Lsearch is probably harder to program~~ It is particularly useful when the value of P_2 is available in ~~Markov~~ Markov form — e.g. F_i can be generated by a Monte Carlo procedure, say, by ^{randomly} selecting operators or nodes & with probabilities that correspond to their a priori probabilities

Maybe omitted.

Section 4 on Trailing Sequences discusses some of the properties of Lsearch and conditions under which they ~~are~~ ^{are} about the best way to solve problems.

10

SN | While Lots of info can be put into the PD that guides Lsearch, the GPD (which is supposed to contain all of this info) is not the same as this PD.

I really have to clarify this in my mind! Perhaps leave it out of this version of the Report. (I like to think of GPD as "containing all of this info", but it seems that the GPD does not "guide Lsearch"; Part of GPD "guides" Lsearch — what does rest do? The rest can be an arbitrary ~~stochastic~~ stochastic knowledge base.

Also "I. factor of 2" : ~~I~~ ^I factor of 2 // Lsearch "factor of 2" is in the update/search routine.
I really didn't mention rest in the report.

20

0

83290

00: 83290 :

~~XXXXXXXXXX~~ Presents a somewhat different direction that training can take

In section 2 on Inversion problems, the system had a number of PST's ~~that~~ that had been successful in past problems and it had to decide ~~which~~ ~~what~~ ~~to~~ ~~by~~ ~~them~~ what probabilities to assign to ~~them~~ them in working on a new problem.

The ~~trainer~~ ~~can~~ ~~greatly~~ ~~augment~~ this set of PST's by including techniques that the trainer ~~has~~ ~~has~~ found useful. The system would learn what kinds of problems to ~~analyze~~ which those new PST's could be applied. It would ~~be~~ ~~state~~ correspond to having a student memorize a set of PST's, and learn how to use them.

~~We would~~ We would like the student to understand the PST's ~~and why they~~ ~~work~~ — how they were constructed and why they work, so the student could

~~use~~ in want new PST's that were more appropriate to new problems. Correspondingly we would ask the system ~~to report the set of~~ ~~in~~ ~~the~~ ~~system,~~ ~~we~~ ~~could~~ ~~ask~~ ~~it~~ to report the set of

PST's as the "acceptable sentences" of a language ~~and~~ for which it must find a ~~stochastic~~ stochastic grammar — which it could use to extrapolate the set of PST's. This process would be much facilitated if the trainer first "factored" the PST's into ~~the~~ common useful subfunctions.

A very ~~experimental~~ ~~system~~ ~~might~~ ~~be~~ ~~able~~ ~~to~~ ~~do~~ ~~this~~ ~~without~~ ~~external~~ ~~help,~~ ~~but~~ ~~in~~ ~~the~~ ~~early~~ ~~development~~ ~~of~~ ~~such~~ ~~systems,~~ ~~external~~ ~~aid~~ ~~of~~ this sort would be ~~needed~~ ~~needed~~.

433.10:

(Section 15) ~~Candidate Marked~~ ~~Efficiency of the System~~ ^{fixed}

Subsection 1 ~~On Updating~~

~~After the system solves a problem, taking time, it must update the associated probability distribution. This involves looking for repeated ~~functions~~ functions or sub-function in the solution and/or ^{evaluating} ~~redundant~~ operations~~

317.34 and 318.03. ~~The latter are inductive~~ ~~Reactions are dependent for~~

P

Both of these tests are open-ended in the sense one can spend any amount of time on them. A good approximate solution is to spend as much time on updating as one spends on solving problems. This is within "A factor of two of optimum" in the following sense: ~~Suppose that by using a different ratio of problems.~~

Suppose that we have a system that uses a different ratio of ϕ update to problem solving time and is better than our system.

Then if our system has a clock that ticks twice as fast as that system and uses equal time for problems and updates, it will ^{results} ~~at least~~ ^{have (small) that are ~~at least~~ good} as good as the other system, since at any time, it will have spent at least as many machine cycles as the other system, ^{for} for both problem solving and updates.

Subsection 2 Ultimate limits of ~~the~~ system

Normally, in the course of its operation, the system looks for better PST's & since it uses a universal language to describe candidate PST's, there is really no limit on what PST it might consider. ^{It} ~~It~~ might find a PST that is much better than L search.

In which case, it would ~~substitute~~ ~~use~~ ~~as~~ ~~soon~~ nominally use L search, but since that PST is right close to one, so that all most all of its time would be spent applying this new PST to problems. It will have ~~replaced~~ ~~effectively~~ replaced L search by the new PST.

Subsection 3. Genetic Programming.

In section 1, on Inductive Inference, we mentioned the possible use of Genetic Programming (GP) for updating. ~~As Genetic Programming~~ It is a particularly good kind of ~~problem~~ problem for GP, since we are ~~not~~ always working ~~on~~ ~~the~~ "same problem". We can use the same population but ~~start changing the fitness~~ ~~each~~ each new problem given to the system changes its "fitness function" somewhat.

Current GP systems are much improved when the fitness function as well as the mutation and crossover ~~algorithms~~ are allowed to adapt to the problem being solved and to the ~~size~~ nature of the extant population. The techniques of section 1 ~~can~~ can be used to implement this "adaptive" ~~to~~ ~~the~~ "adaptive" ~~to~~ ~~the~~ ~~problem~~.

~~Mutation and Crossover are a very restricted set of techniques to obtain~~

In G.P.,

New populations are obtained from ~~old~~ ^{old} by ~~the~~ Mutation and/or Crossover — which ~~are~~ ^{were} originally attempts to simulate organic evolution. These operations can be usefully generalized. Mutation can be considered to be "induction ^{from} ~~with~~ sample of one"; and Crossover, "induction ~~from~~ from a sample of two". Classical statistics is not amenable ~~to~~ ^{of} small sample size. Bayesian statistics ~~(which also remains probability is a function)~~ ^{of} can deal with ~~the~~ ^{the} small samples, but the prediction depends much on ^{the} prior distribution. ~~Appendix 1~~ ^{Appendix 1} tells us how to ~~obtain~~ ^{obtain} a reasonable a priori distribution ~~and~~ ^{1, and 2} update it ~~with~~ ^{with} new data. ; Section ~~1~~ ^{1, and 2} tells how to update it.

It is clear that we don't have to use samples of 1 or 2 ~~in~~ ⁱⁿ any sample size

size ~~sample~~ up to the size of the current ~~population~~ ^{population} can be used.

The ~~induction~~ induction methods of sections 1 and 2 can be regarded as a kind of GA using the entire ^{population} of PST's ~~to~~ ^{to} induction to create ~~new~~ ^{new} PST's ~~whose~~ ^{whose} parameters to the old ~~RAST~~ ^{RAST} since the global fitness function changes as each new problem is encountered, the ~~number~~ ^{number} of population can be ~~parameters~~ ^{parameters}

for the Report "AR"

critical items

1) Fix into so it doesn't have ~~each~~ ^{each} one. ✓

2) Put in eqs. nos.

3) Bibliog.

4) Help: ID/IA, vis. prof.

417. → (6)
415.03 → (7)

equs: § 1 equs 1, 2

3, 4, 5

§ 2

6, 7

A P X A

(8)

ref. is seen in (8)

B

(9)

2. refs

6.11.

§ 3 refs to both 6 & 7 only 2 refs. by § 3 got them

Sec nos. men:
4 MV
5 TLOZ

§ 2 col (p) 2
3

2-updating) should be ~~sub~~ ^{sub} of section 1

20

30

cc: 835 [SN] In dealing w. correlas betw. cond. in Lsrch: One way was to take ~~the~~ equiv classes of cond. One member of eqv. class would be chosen is given all ~~the~~ Lsrch cc of all other cond. into eqv. class. This is a way of dealing w. b. fact that Pcost is often $\gg 2^{-kcost}$. But this eqv. class, idae puts even more wt. on t. single cond. This can be a way to really t. effective RC's $\approx 10^4$!

[SN] Sumacs can be a v.g. approach to long. ~~to~~ w. e. Azadi ~~log~~, I have ~~step~~ in this direction by including context dafns, - another stray step.

Imporant, ~~Sumacs~~ ^{Historically} successive Sumacs can all be universal, but t. "constant factor" by which each ~~one~~ deviates from t. needed soln. to a problem can \uparrow exponentially w. historical bias, this "constant factor" type in "Kolmogorov studies" really makes most of t. work irrelevant so Lrag: $\approx A.I.$

Notes on proofreading IBSIH Report:

1) The pros ~~to~~ are for LIV 1993 ~~to~~ 1997!

~~to~~ p2: ^{Abstract} transcribing: where is this discussed? ~~Actually~~ ~~last~~ ~~Ref~~ ~~of~~ ~~SI~~ (I omitted this! ~~to~~ original printed) P10, 10 - on context: ^{help} should be lead to discussion in Appendix 3 Scaling

P6 $\sigma_9(3)$ looks bad! And page: $(293-3)$: $P \quad R^k \rightarrow R^2$

P8 update ~~to~~ ~~the~~ ~~method~~ ~~5b~~: "Merging" P_2^i 's. We also have to "merge" b. assoc. R^i 's.

New R_2 "or" of old. However I ~~did~~ ~~already~~ say that we "merge" P_2^i forces and assoc R forces. So (at rfg).

P7 on updating: "Modify R^i " so New $R^i(Q_i) = R^i(Q)$ But: \dots \Rightarrow How possible to be done is not stated. Lsrch was mentioned earlier but not explained! Then p. 8 we modify R forces \Rightarrow p. finally, b. we don't say how we do it. On p. 9, Lsrch is introduced, but reader does not know its Lsrch R^i is being explained. \Rightarrow p. 9 is first born of Lsrch?

Tentatively: (P9) After ~~the~~ ~~first~~ $\frac{1}{2}$ time!

To start off we have only 1 ~~copy~~ ~~of~~ ~~the~~ ~~text~~. show how to ~~use~~ ~~the~~ ~~function~~ ~~of~~ ~~Lsrch~~ ~~to~~ ~~find~~ ~~a~~ ~~function~~ ~~sets~~

How to deal w. this? One way; (not unnatural): Do t. simplest case of R is discrete ~~then~~ first

~~to~~ we will first describe the operation of the system for ~~a~~ very constrained type early mainly of learning of the system, then deal with more complex cases, using simplified problems, then we will ~~fall~~ ~~how~~ ~~to~~ ~~update~~ ~~for~~ ~~more~~ ~~general~~ ~~kinds~~ ~~of~~ ~~problems~~. In particularly learning, we only allow one $R \dots$

(copy to end of p. 9) Then ~~we~~ ~~introduce~~ ~~more~~ ~~general~~ ~~case~~: start at top of p. 7

Some more remarks.

p. 8 near end!

20: 836.401 Intro of our sample case: When we are given a new Q_{new} , how can we ~~do~~ ^{for} more complex problems with several R^i and many P_i^k , especially

We ask, how is the system modified?
Start top of 7, a lim. first instance.

ff More Generally Suppose we have a more complex system with several R^i and several P_i^k .

When we are given a new:

The learning techniques must be ~~modified~~ considerably modified for systems with having several R^i and ~~associated~~ associated P_i^k . when we are given a new.

P17: Appx A: Introduce 'Agreement list' feature.

we start ... function definition. We will call this the "Agreement list". In the course of our calculation we will occasionally add to and delete items from the Agreement list. Some examples of primitive functions could be .X, Sin,

T. Appl. about "factor of 2" (make sure that we use || or random (such) in do regular update "Einsparung" Bayesian

Boys Boqas

CLRC, Royal Holloway, U of London.

Acknowl: ~~we~~ would like to thank IPSIA for providing productive month

our facilities that inspired this report. Particularly Marcus Hutter and

Juergen Schmidhuber for ~~their~~ comments and discussion.

~~Eliza~~!

Tips Workshop: Nikolo Case - Bianchi: (Ud. Milano, Italy)
Universal Lang., Anne Perry, P.H. reason

Eliza Newman man (d. Calh. Santa Barb).
Universal Lang.: A View of a Bayesian.

Representation OOPS (Jensen's paper)

Mention B.A.D. title.

P5 (Oswen) "new in vantage": This is identical to Lech!

3

He uses ~~Osma~~ ^{w.} a p.d. that changes adaptively: Is this ~~to~~ ^{to} distraction?

2) P8 on "Frozen Man": This is an attempt to get around the fact that in an old system of his, ~~any~~ any day in ~~man~~ man could be destroyed by a random "write" mechanism.

It seems like every A.G. soln. On the other hand, it would seem to be a good idea to store ~~all~~ ^{formally} ~~pool solns~~ ^{to problems}

3) P5: 3.3.4. So ~~perhaps~~ ^{is} saying that OOPS works w. proper ~~TSQ~~ ^{TSQ} (only)

4) It ~~is~~ is not clear that "Optimal" support to apply to ordered problem solver.

One interpretation is that it takes a bunch of problems and decides what order in which to solve them (P5 was my ~~original~~ ^{original} in the parallelish)

5) P4, because universal search ~~can~~ take much time → no reason to avoid universality

6) It's not sure but I suspect that what I mean by "Backtracking" is: ~~seq~~ ^{seq} ~~we~~ ^{we} ~~are~~ ^{are} ~~branching~~ ^{branching}

~~a~~ ^a ~~tree~~ ^{tree}; depths ~~test~~ ^{test}. We use ~~the~~ ^{the} ~~usual~~ ^{usual} ~~usual~~ ^{usual} ~~trick~~ ^{trick} to do this ~~trick~~ ^{trick} (Sussner "oss").

It works w. any radix - even ~~change~~ ^{change} radix.

ABCDEFGHIJKLmnopq
AB

7) Could I use his "Prefix string" (idea to implement Sussner?)

8) P32 t. Probly changes ("Braschifting") instructions: It seems that this is not a v.g. way to do this - but I must look into it. I is into "Credit assignment problem" (probly a "wrag problem")

9) P20 He uses 1.5 GHz pc: ~ 1000 clocks per "time step" seems very inefficient!

10) P20 ~~ff~~ ^{ff} discussion of expt results may be useful.

11) I want to find out how he changes bras (Also note 2.1.39 ff w. "branch" pc values; Also note PE P's 2,3,4)

12) First find out how it works w. constant bras, then → P32 perhaps

How it may work ~~trick~~ (w.o. changes of Bras): Machine ~~can~~ ^{can} request a

token (alphabet input) or work on a problem set. It checks a step state: Always goes to stop state when it has solved all problems - but then announces the fact. It may stop earlier - which is a "failure".

I think initially, all tokens have I pc. ~~ff~~ Also saves all solns pgs to problems in "Frozen Man"

To change pc's: new token can be derived as substrings of ~~existing~~ ^{existing} token set.

Also other ways to change pc's of tokens.

It may be that normal execution of trial functions, a revision of pc's is done as a unified task.

— Possible Good: ~~but~~ ^{but} elzn. can also be very useful. - see 2839.0d(±.02) ©.

So "between tokens" Machine can access frozen memory

- ① Input new tokens (\equiv strings of primitive data types)
- ② Assign PC's to tokens
- ③ Work out problem "Computer" PCs could be outputted or how to change PC's.
- ④ Output some of problem.

Since ~~input~~ one input is a token & a token can be $\frac{1}{2}$ byte, a lot of computing can be done on "one token"

Is the only requirement to system Guts is its correct solutions of problems? \Rightarrow insertion of these problems in frozen memory? — Also if P_n solves all probs upto $n-1$, P_n is used as a trial prefix for problem $n+1$.

I do use PC's of sols as additional f.f., info.

$\sum_{i=1}^n 1 \cdot P_i = 1 \cdot P_1 \cdot P_2 \prod_{i=1}^n (P_i + 1) = \text{converges!}$ its min is 2, its max is e.

for min, $P_i = 1$; for max $P_i = \frac{1}{n}$; $\lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n \rightarrow e$

$1 + P_i < e^{P_i}$ so $\prod (1 + P_i) < \prod e^{P_i} = e$. If $\frac{1}{2}$ bytes PC's is 8 then product is fraction $(1 + \frac{1}{8})^8$ and e

It would seem possl. to use .09-11 for $T \leq 3T$ but actually, ~~some~~ $(1 + \frac{1}{8})^8$ and e

T time to do a "round" is always $\leq e$ and if $\frac{1}{2}$ byte factor is over < 1 , then $\frac{1}{2}$ recursion is over. The $\frac{1}{2}$ occurs ~~whatever~~ whatever we have

to termination of a trial before time out — e.g. it does something illegal or it is earlier but to result can't be a solution to the problem.

What was the thinking about when he calculated $B \approx \frac{\text{time of entire search}}{\text{time of entire search}}$? P14.25 $\approx \frac{PC \text{ gain}}{CC \text{ soln.}}$

Did he want to modify all PC's so that he got B for PC of soln? —

i.e. modify PC's by something like numerator $\rightarrow num + k$, denom $\rightarrow all \rightarrow denom + k * (\text{no. of tokens})$.

This revision of the PC's would reflect the fraction of trials that aborted before "Time out" as in .14-19 i.e. $[B \text{ v.s. } \frac{PC \text{ gain}}{CC \text{ soln.}}]$ if none of these 2 would give "index of non-time-out reports".

Made discussion: ① title ② type on ~ ~~page~~ page 0. ③ content about $\frac{1}{2}$ such data, what cjs is out. ④ $\prod (1 + P_i)$

- ⑤ ~~series~~ discussion ⑥ Just what was his PC modifi. scheme? (I may know)
- ⑦ would PC's scheme work?
- ⑧ was his conclusion over "working to ~~1~~ 1^{n-2} problem helped solve "tower of Hanoi" logic?
- ⑨ I suspect that his search & finds are really grossly large.

Re: 839.30-31 : The idea of using info in old solns to help work new probs is help calculate pc's or next token.

In my reports, the system does have (about) all old solns "in the next list" — so it does do partly what .00-.01 does — but does use 'old info to completely recalculate pc's at next token? Is this a good idea? Normally, I think J uses "Laps rule" for pc's at tokens: The "boost" instructions is rather gross, as a method of how to recompute pc's!

I think his system is more general than "discovery of context" for new tokens to control pc's at new tokens — { the I did have this idea of "Generalized context" categorization of probs. — so maybe his system isn't any more general than mine! }

New Comments: I think his system is designed to solve INV problems only.

The w. my "contexts", I had a more limited idea as to how context could influence pc's of next tokens. The result would be that the system would learn these factors, but would not be as general. J's system is very general, but takes much longer to learn.

Also, it may be hard to train because of. training no good idea (e.g. model) of how to learn it.

As is, J's system can try any rules it likes for assigning pc's to tokens. we know a priori that only some such rules will not work. We have a way (realist) to generate rules that are very likely to work. (Unclear as to whether there are rules outside this class that work)

Another possible difference: J. uses strings to represent parts, such as: I use functors/frames.

Re: Optimization problems: Probably J's system could be easily adapted to do it if my section 1 does it: For more complex method of modifying the guiding pd ... I don't know if it would be possible.

[SN] In my original work on the QA machine, I "decided" that it could not really learn to work on non-optimal problems. That if I gave it examples, it's "best guess" as to how to generate these examples would be to do or "approximate" the program that created these examples: Hev. NOISE! If TM had proper vocabulary, debugging corpus would be relatively easy, a debugging program created to corpus rather difficult. (Law PC, by cc). So it would depend on TM's state of development, as to whether it would discover "the idea of optimization?" Whether it would discover it "completely" is unclear.

An analogous idea "Discovering Laws of Algebra" would be easier if TM took long time to do multiplication, division — maybe add & subtract also: So the laws of Alg. could save TM a lot of time. In very fact Arithmetic, the laws of alg. are harder to discover.

.00
.01

10

20

30

90.8
90.8
998.8

III

The General "pc of token" modifier seems to be a good general. of "Generalized context" to pc. of next token is a function of ~~previous~~ present "prefix" and probability.

Rec corpus for this induction is all (problem seen, previous ~~token~~ to learn) for plots. maybe pc known as part of corpus?

Could this induction create new tokens? ABCDEfghijklmnopq

Phi this result about "tau of Hanoi" being more Solvable after "2" was solved! I'm a bit incredulous of this; I want to think about it more.

I feel that his execution time ~ (10^9 - 10^8 ops) is excessive; for a problem size,

I would expect ~ 100; maybe 10000 at most.

I feel that his TSP is not V.G., in the sense that he is not really controlling things

1500 clocks
That it should take 1500 clocks to execute an instruction is rather wasteful - but at present time, you are not particularly concerned with speed - you just want to see it and how it works.

Probably ~~some~~ ops could simulate my method of ~~some~~ using pc's tokens by just giving it ~~some~~ using its program a "Token" - w. perhaps much initial "Booting".

Out "Success" of his plans: if "boost" inst seems critical to deciding which tokens will

Get reasonable pc's. Look into this!

The ~~description~~ way the system deals with induction has several phases of ~~some~~ sophistication. In the first phase "The early training" (P8 my version)

Induction is pretty much an inversion problem! We want to find a function FJ such that we want to satisfy eq 5 w. Pro stage $\epsilon \approx 2^i$ as possible. Since we do trials of FJ in approximate order of 2^i , ~~the~~ the first FJ we find that satisfies eq 5 is usually acceptable. The trainer may want to "overshoot" for a FJ w. ϵ larger 2^i . In this phase, we use the limited doubling technique of Lsearch -

we start out with small ϵ - τ_{00} ~~and~~ using 2^i as time limit, then double the time limit in next round, etc.

In the next phase, we do induction problems involving probabilities reduced on the A_i . We want to get a maximum value for eq 1 (P4).

In this case, using a τ_{00} $\tau_{00} = T_{max}$ as in section 4 - which is what you (Marcus) suggested.

(S4) Another drawback based on my & J's approach - the last part tried to do induction problems - but only using problems.

Note that .20 - .30 is not really "Lsearch" - which was designed to solve inversion problems only. Section 1 induction: find O_j such that eq 1 is satisfied; look for O_j in 2^i order.

For induction, we write O_j as a problem domain, A_i as an acceptable answer - as in INU problem. We could use same formalism for induction as well.

-00

because the machine is solving problems that it can't solve. (Obviously, the machine can solve a certain ~~kind~~ problems that it can't solve, but I'm not considering that kind of problem. There are lots of very difficult algebraic problems that I can't solve in a day.)

Major Q for it: where does he go next? what problems? - what would be good to do if he could be continued? (Experiments are expensive. possibly ↓ cjs by ↓ no. of units.)

(SN) An efficient way to do OOPS in 11! Have N machines: Each one does every ~~part~~ every point in the branching tree. ?? Not so clear!! Have first machine scan for "markers" at branches where another machine has been assigned. This ~~is not~~

Tree search is easily done using "stack".

20

30