

60 ID 843.40 Inverse of Conv. Perm:

- 01 1) Any distribution over strings that gives > 0 pc for every string, is Universal. $\rightarrow 2.12$
- 2) 1) \rightarrow seems easy for discrete over finite strings, but may work for infinite strings ("continuous" v. v. l. D. R.),

3) If conv. theorem is not true, Perm must be at least once μ_0 for which $\#$ $\mathbb{P} \leq \mathbb{P} \leq \mathbb{T} \leq \mathbb{B}$ is false: (e.g. $b = +\infty$, $P \in \mathbb{N}$) $\rightarrow 6$

4) in $\mathbb{P} \leq \mathbb{P} \leq \mathbb{T} \leq \mathbb{B}$ we can have $b \leq b(n) \leq n$ bits in seq. so if conv. Perm is false, $\#$ only need be false for $n > n_0$: $\#$ for $n < n_0$, Perm can be true but error/bit will be only bounded by $\frac{b(n)-1}{1}$ $\rightarrow 6$

5) if 1) is true for infinite strings, then for any positive δ , we can do $P \rightarrow \cdot P^\delta$. (normal constant) for large δ this means pc ratio of $\frac{P}{P^\delta}$ will $\rightarrow +\infty$ or 0. so system becomes

deterministic, rather than stochastic. Hier, ~~conv. Perm.~~ $2 \leq \text{error} < 1/\delta \rightarrow \delta \ln p_0$ so error is worse: we should do $\delta < 1/\text{diff. distribution converges}$.

22 - 6) If $\mu = \mu_0$ ~~Perm~~ $b = +\infty$ ~~Perm~~ must be zero for some n_0 \rightarrow approx proby

26 which gives KL distance $= \infty$ for strings of length $> n_0$, but that P could be true for $n_0 > n_0$! No! P will ϕ for all $n > n_0$: - P assigns prob to strings of length ∞ . If $P = 0$ for strings \leq then $P = 0$ for S_2 , where S_2 is any string, finite or infinite. \leftarrow This doesn't seem right! i.e. it should be poss. for \geq Perm to give ϕ pc to something that occurred, Perm gives pc's $\neq 0$ to subsequent bits in a usable way. $\rightarrow 2.00 \rightarrow$

30 7) $\# \mathbb{N}$ in this proof not eventually Perm will be an enormous ratio between normalized & unnormalized pc: My "sort of counter proof" was for $\sum p_i(v)^2$ bounded.

32 But this was not also for $2 \sum \text{KL distance} < \frac{1}{\delta} \ln p_0$. so ratios (sort of) have to converge

Bit tech at this carefully - see just what converges! It may be that my proof of $(pc)^2$ converging would only work for \downarrow ~~sq~~ $\#$ ~~KL~~ KL distance

20 : What to add to Report:

- 1) That stuff I wrote comparing to ~~last~~.
- 2) Design of OOPS
- 3) Tell how Sachin's App ^{and A} are
"Stand alone" — secs ~~2/3~~
on Int. s' ~~2~~. tell how to update
much better using Oprodution — a
ultimately how saci should update itself.

~~Compare to the~~
~~after work in this~~
Related work.
not need
(I have ~~used~~ ~~for~~ ~~last~~
work in Int. Lang
It's may be a problem.

Criticisms of OOPS!

- 1) S.2.2.1 He doesn't have lang in to put into computer now PC's for Tokens.
- 2) He has large CJS: probably will soon find it very difficult to write TSC's (T. "Autistic student" problem)
- 3) Scaling problem w. "Eosky" ~~is not~~ previously useful pems.
- 4) S.32

4) Perhaps pseudo pfm on how > a confederate
GCPP; Advanced updates work.

5) Begin work on my "Advanced report"
list of brief variations in previous, in Refs. in Notes.

6) Put talk itself on web, but add details
of a big slides — make slides ~~rather~~,
(when needed).

7) Send copy of report + OOPS to

<ul style="list-style-type: none"> • Ivan • Kurt? 	<ul style="list-style-type: none"> 2 Good students ^{at MIT} • Henry Krakerman • Salford • Jess Martin • Louis Martin 	<ul style="list-style-type: none"> • Phil Lasler ^{from Knabing} • Will Garsh. • Gary Wolf. • Eric Winter 	<ul style="list-style-type: none"> • Phil Apley • Guy in Berlin (ex. grad stud of Tuango..) • Guy in London UK. (Coke w. video (smars) ^{It's called hum}) • GAMBERMAN 	<ul style="list-style-type: none"> • Dick Harter • Steve Wilkin • Minsky • Carlos • Murray • Alex • Dowd • Wallace • Guy at Santa Point.
---	--	---	--	---

8) Get LaTeX 2/0 pdf version of OOPS to send out
w. my "Report" — Look at J's website.

CLUSTERING

ED
NIPS

This was in response to a proof (maybe at NIPS 2002) that clustering was impossible to do in an exactly rigorous way.

Grand clustering: Given a data set: each pt. of set has several parameters. \vec{x} is m dim vector

A "cluster" is way to ~~say~~ a) pick a pt. in n space, b) find a

Scalar function $f(\vec{r}_j, \vec{r})$ (\vec{r}_j is data pt. \vec{r} is "center of j th cluster")

$\rightarrow F(\vec{r}_j, \vec{r}) = \max$ overall second arg. values. Perhaps $\sum_j F(\vec{r}_j, \vec{r}) = 1$ & $F(\cdot, \cdot) \geq 0$

so $F(\vec{r}_j, \vec{r})$ is a ^{kind} probability that \vec{r} is in cluster j .

Anyway: The idea is that this is a way to decribe the set of points.

How it down goes: ~~Each pt. is assoc. w. any cluster j & all data~~

~~Each pt.~~ If there are j clusters & n pts., each pt. can be described j ways. We want f 's such that total info in f desc of f .

functions f plus f_j \vec{r}_j . ~~plus~~ f plus f_j info in f data desc, is ~~min~~ min.

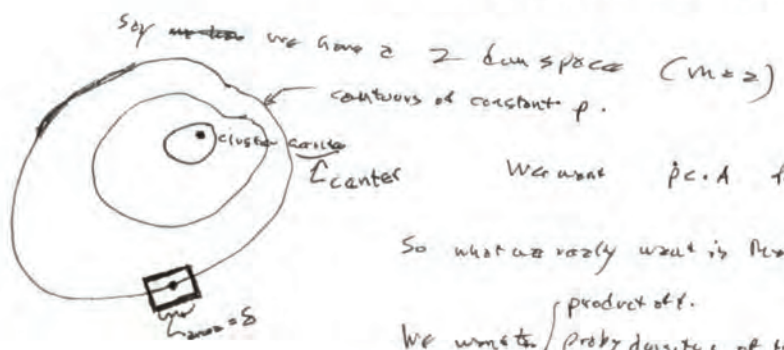
T. pc of each data pt & some of its pc's wrt. each of the j clusters:

We want product of ~~pc~~ pc of each pt's times (pc of j centers) \times pc of $F(\cdot, \cdot)$ to be max.

We can compare it w any other clustering of these pts.

T. forgone is a way to describe set. ? — Yes! .22 - .40 shows how.

.21 is unclear! In a cluster $\sum p_i = 1$, but how does this enable us to describe each of the pts. wrt. their center of k -cluster?



What we want is the product of desc by z data pt. So that it is within an area, A ,
[$(p_c \cdot A)$ will be $oc A$ for small A .]

$$\int_{\Omega} F(\vec{r}_j, \vec{r}) dV = \int$$

V is a volume of elements (Ω dim) space

So what we really want is that product of f .
We write f (prob density) of the data pts to be max, when multiplied by pc of f . $F(\vec{r}_j, \cdot)$ desc.

T. actual ^{pc of δ} desc of the data set to by accuracy can be obtained by multiplying above pc

by δ^n (for n data pts & we want to know each pt. to within area of δ).

but $f \cdot \delta^n$ factor is indep of the clustering method. So ~~we~~ choice of clustering method does not depend on δ .

101B ID

12.35-36 Seems to be: v.g. Genl. Soln. to this problem!

This ~~for~~

If system has no idea of what's being done - much less a concept

"Optimality" - If is, here, a ruff & dirty way to work on a problem. It extrapolates to set of solns, no. idea which solns were v.g. or very maximal - or which problem was related to present problem

J. would (perhaps) work on it in a different way! He would try various methods (PSMs) for assigning/medifying pc's of new relations. He would then return to ones that worked best. This is a simple but very slow method of learning.

A possl. way to improve this: T. data set being used in previous drawn was unordered set of soln. strings, $\{F_i\}_{i=1}^n$; want to extrapolate PLS set.

A slightly ^{nicer} way: Given $\{ \begin{matrix} \text{Prob}_i \\ \text{Soln}_i \end{matrix} \}_{i=1}^n$

Given Prob_i - to get good pd. over Soln_i : This is (ruff) operator induction.

It is better than my old concept of "Genz Context" - The Genz context ~~could~~ include probands ~~data~~ data. - so Genz context could be a "ruff & dirty" approx. to +12-13.

This may be a better way of picturing it!

N.B. There is int. system at 5.25 ~~ff~~ 6.10-12: 2 kind of "BUSINESS"

If is, to some extent, a "self confirming hypothesis" - The regularities observed in the past, are used for future trials & tend to be "reinforced": This is not the best way to do induction, but it is probly [much?] faster than the correct way.

T. Correct way (ALP) would consider all possl conditions, then do something \rightarrow 7.00

It is interesting that J. did not give all solns equal wt. for predicting next token. ^{the case} ~~Manually~~ ~~weights~~ of ~~the~~ ~~solns~~

$\frac{1}{2} \cdot \frac{1}{n-1}$ for solns 1 thru $n-1$ & wt. of $\frac{1}{2}$ for soln n .

I was thinking of very = wts for all solns. - so $n-1$ times as much wt. for soln n , as for the others.

His reasoning was that to soln. n did do all other problems. Also since he only had to ~~choose~~ "add on" the soln, ~~tests~~ is not checked ^{on} for other problems

That took much less time for trials. \rightarrow [I'm not quite sure how he did to "not run n soln continu." trials. ~~Check this.~~] Exactly what was done is certified in deciding whether system really did anything it was doing in solving (Tow) of Hangi after 2^{14} recursion.

0.19.02

IP
15/1/85

7

.00 : G.24: Like compare d.f. of ~~Garco~~ Garco at Continuum. At first glance, this seems difficult,
 .01 but maybe not so difficult! We have various problems, complete soln. desc., Gars;j triplets.
 The Gars can be different instances for different problems, & still be part of
 the same ^{date} corpus. Anyway, we use the data of .01 to induce the Gars distribution
 for all poss. taken combinations of the present prefix (wrt the present problem).
 From this Gars D.F. we can get the pc of any particular taken
 combination, giving max Gars.

.10 .00 It is a bit "fuzzy" & I want to write it out in more detail. - It
 looks (in present form) ~~is~~ easily forgettable"
 . O.k. Corpus is all known triplets of .01. Gars include time for solns
 of Inv. problems (for we have to make ~~factor negative~~ No!), as well as any other
 kinds of Gars problem Gars.

Looking at it from a slightly different viewpoint: ~~we~~ we have:
 (prod desc., positive, ^{plus for} continuation taken TK_n) want D.F. of ~~them~~ whatever
 Gars type is specified by problem. \rightarrow (as corpus for this induction, use
 Gars data of .01). \rightarrow (9.00)

.20 [SN] I'd like to be able to use negative info - Usually this
 takes form of knowledge of "No soln." for \geq time T_j for a given problem, soln. attempt.
 It is a INV problem. - but for Opten problems we always ~~know~~ (?)
 have a "Gars desc" for any Opten technique. This can be -oo, if
 no trials have been done yet.

[SN] For 1. primitive inst. set (like for FORTH) would best be obtained
 in a RISE cpu - but I think there aren't any any more! - Maybe
 Caruso? - Pick a fair sized inst set to start - Run the system
 cuts down on the set it uses - produce an optimally fast subset &
 giving it most of the pc wts. - we could obtain other insts at low
 pc's (but ~~was~~ possibly vary by pc's in certain "environmental contexts/conditions").
 - Actually, this is what J. did, but ~~hard~~ (i.e. selection of by
 pc wts ~~under~~ ~~car~~ run certain contexts) - but he didn't use cpu opcodes
 directly.

.37 Perhaps use inst. set of 8 bit cpu like G502. see 12.35-36
 \rightarrow (9.00)

T. ~~is~~ induction problem, w. "pc for best" calcn. of .50 - .10 will be
 to main (almost entirely) Make problem that system will work on. It is same as
 problem of ordering $BCFD_1 \Rightarrow GCFD_2$. Actually it is "T. same" in all ways! 8.10

ED
NIPS

"704TM"

in the early history of project.

10: 7.37: I may want/need, a hard wave clock, for time plugs. It may be easy to have a counter running at CPU clock rate, or $\frac{1}{8}$ clock rate. - just pickup clock signal somewhere & try to find fast counting. etc. as "off V. Shelf" item.

At first, system will be very slow, using many op codes per "in, medium", & I can use relatively slow counter.

09 The idea of 708. is a 704TM ~~with~~ idea of 5.00-07. (21)

0: 7.40! - Its + only meta problem, i. it includes all L such!

NB The present "Adventure Rush" of optimism, assumes t. cards are uncorrelated. T. result is a rather (limited/stupid) TM, I guess!

11: I had that very fast HW. idea for TM: Look it up now! It was frail frail.

2 saunders very parallel machine that moved addresses around. It may have also created various HW. functions that were physically duplicated if they were used much - then address "diverted" for different functions when a few functions ~~be~~ were more severely needed.

11:06 -> I think t. essential idea is that t. 2 level GPD update ~~include~~ can automatically include any recalculation of pc's of taking out

OOPS (in theory) does

25 -> Big trouble w. 5.00-07 = is discovery of needed leaves. T. system may be able to work under interesting/diff. packs, but it will not approach human skills. ~~text~~ ^{capt} by substituting ^(hardware) space for "heuristic insight"

perhaps study hours: list Ray try to ~~find~~ find common concs - try to water "Grammar" for extrapoln. : J. Pearl has book on Hours (but I may have: also most of rest of the community) - It's not sure how

This could be a main part (say of J.), because it will be able to "win on various benchmarks" U.S. other Mark Eng. systems. In general, days .31 will probably be most imp. part. The "Hours" are a "promise"

but nothing (perhaps) in way of provable superiority. Usually normal programming skill is buying of Biz Computers will be preferred, over development is understanding of Hours. perhaps they will try insubuy hours, whenever possible.

old sched.	lv	arr	Chi	Chi
	215	815	dep	895

My water is 30 minutes

PM	PM
9:28	5:45
2:28	2:45

enjoy of A.I. E. H. T. one is hard and E. now out

D.19.02
TD
NIPS

8:16 AM

ty -> p. is likely
to be after "8"

File UA 516:
(Lecture 8:458; Board 8:158)

8 1/2
21.15
9.28
9.15P

00: : So: perhaps impt. "Political" problems!

Re: Uvagon: Main guess is that he doesn't understand Operator induction, much less my formal Soln. (is more toward a practical soln. of it).

05 Re: T. Boost instructions its Applicn to enable ~~the~~ emphasis of tokens used by previously successful team: T. T sq. was certainly not key enuf for disc-usage of this "Heuristics": It would have to have several/many cases of its expected use before its "definition" were done. T. case in OOPS is not even Choc (hot Lugs!)
10 It is simply inserted by Kramer, not "L Fed".
11

18A
HARTER
W. Dip
Hitachi?

As is, there is very little (only Boost instruction) for usability Modifying pc's of Tokens. My guess is that J. would want to add more instructions - but that they would be to more or less A.H. implement heurs that he feels are useful. As his system is new - there seems to be no way for it to be usable learn new Heurs or any other useful Algms for modifying pc's of "Tokens":

22 Major Criticism of OOPS is that it really has no provision for its learning how to assign good pc's to Tokens. He says it has all these capabilities of modifying pc's of its tokens - but presents no good way to make his come about - other than by truman - pying in Heurs; invading new Tokens.

30 **[SN]** In Learning Text Book Algebra (5.00-07), it may be easy for it to know/learn Logical "reasoning".
To "identify" Heurs: Try to see why certain "answers" or whatever were "obvious" to me. - Also, Look at New-Simon on "Human Problem Solving".

32: : Compare Theoretical vs. actual cc of solns to Towk. of Hanci (as given in ods paper): See that it really finds anything - in its learning recursion" whether the pc of today is into "Proof" or is not as searched as it. aut. gained.

00: 7:21 : I want to work out details of 7.00 - 7.21 : see if its 2 hours or less
.12 - .14 seems correct!
Complete Soln. of TM (w.o.t. TSO, of course!).

02 Then look at various "Early" approxns that move toward that goal. 7.07
Also, look into possibl. use of Big Computers to reduce diffy of TSO does spn.
(Also enables larger CJS - which usually means! more "creativity")
07 02 → Actually, the original 2 part GPD was a very awful good ^{initial} ~~initial~~ approxn. to 7.00-20
T. only real diffrnce is that somehow, the PST is actually created
by GPD_{1,2} !?! Just how this is to be done, is quite unclear!

Perhaps a kind of modifn. of the old GPD_{1,2}, [PST] sort idea?
12 It would seem to be necy to consider an enormous no. of PST's, ~~is~~ do it.
standing old GPD_{1,2} on Mem. From the DEF on PST's, the d.t. out. must be taken
14 can be determined. While this does seem theoretically correct, one must
19 slowly ~~and~~ find reasonable approxns. - otherwise it ^{cc} sounds
Way too high. T. Goal, here, is clear. 10.00

10 12 ff looks reasonable, is not far from what I was doing before
(GPD_{1,2}), but considers all possibl. PST's: The obtaining of
PC's has to be taken (remember a "token" can be a large Macro)
is one possibly useful eliza of the problem.

Assoc w. infinite no. of PST's being optimized or what
is probly a Big SOY problem. - [see my second mullings
on SOY! - T. present (.12ff) problem, rather than an ^(regeneration) _(equivalent)
of the problem, may point toward a soln !!] (likely! → see 10.00)

30 SN Can Cross Validation (aka Faithful) get good estimates of
PC error in ALP?? ... Cross Validation may assume Stationary Data ~~Stream~~ Stream
Also, it is more for average error, not error of a particular prediction!

SN Perhaps mention talk in Report: that Lsreh is Adaptive Lsreh can be
used in sub system (standalone system) of S (2 Appendix A: Adaptive Lsreh is necessary
- Lsreh w.o. it is excessively diff/c.

SN In normal v.s. semi-recursive version of Lsreh Diagram: T. ratio of the 2 "B's"
Is the normal constant. Does this ~~const~~ is this constant a funcn of "N" f.
comp size? Can it → w. n? If so Normal const. much better than summed.

5 sec
00: 9.19

9.12 - 19 may be the long sought soln to a general SOY problem: I haven't yet figured out how to deal with in particular of PST's. — I don't know if it's a Big Problem, hvr. — Maybe trivial. I do various approxns. to 9.12-19, ~~hvr~~ — to degree of approxn, depending on how far TM has progressed.

So the idea is 7.00 - 21, then 9.12-19

10

Re: Soy: 2 aspects of the problem: (1) Is it really the "BEST" choice (2) What is the actually expected yield def? Have I'm only concerned (1) (I think!).

15

T. present problem (7.00-21, 9.12-19) differs in another way from classical SOY. — Am now only interested in the probly that a given PST will be "Best". → (19)

(SN) In section on "TSQ", explain how TSQ is to train trainer as well as TM. for trainer "Debugging".

19

(15) Perhaps a "Main Problem" is "How Logit. are the def's of G of the PST's.?"

20

21

(SN) In doing GP using "SPACE" for Gene evaln (2.1" pareto!) Koza should have spent as much time as selecting cards as he did on Gene cards! (and/or did approx Gene evaln, to start out).

25

As I remember, one set last by Bouts w SOY was a SM, w. large no. of Stock to bet on. STEIN would seem to be very relevant! — Seems to be a local soln, since it gives mean yields. — BUT I have considered this & there were serious diffys! I don't remember just what Ray were! Gaussian D.F.'s would seem not bad. I could use e^{-x^2} (σ < 1) for fatter tails, if necessary.

30

Anyway: If STEIN does compress & data better than other ways, Ray USE it!

(SN) I had this idea about how Evolution might be sped up, by having something like a computer try to extrapolate from the past. HVR. A big problem is how

to get info about ~~past~~ past success/failures into the system! It needs F.B. —

→ A wild kind of F.B. would be the set of ^{Genes} Macros found useful in the past.

That reorganizing of genes out. chromosomes so as to clump synergistic genes is away of speeding up evolv. Are there other ways to take advantage of available A.B. info, like that?

Could we find a better way to use past info?

(Alex suggests: T. older an individual is, the more likely that his genes kept him alive. — But evolution is interested in not merely age, but how many babies this for! Could be proportional hvr.

Bayesian result seemed to depend on first 4. Bayesian result seemed to depend on first 4. I did (apparently) run into some trouble, in first 4. exactly why on the spread of genes?

So: Say the PST d.f. is legit, & I have it for all poss. PST's. Inegration gives the p.d. for "Best" PST - to be used for Lsrch. At first Glance, how the p.d. is not in the best poss. form. for Lsrch. Maybe it is Ok! We need to pc's of the PST's in ruff pc. order - ~~we~~ we don't need the ^{"Takes"} ~~pc's~~ ~~for~~ For opten. problems ~~we~~ (≈ 0.2 , say) we have to CB = 0.06 , & spend ~~the~~ c_{co} p.s. on the ≈ 10 PST's. If we have time left, we double ~~threshold~~ $c_{co} \approx 2$ & add extra time on the ~~various~~ PST's (for all Lsrch) or, we just ^{run} ~~run~~ the whole thing at $\geq c_{co}$ limit. For Inv. probs, we do do doubling (or all Lsrch), but we really don't need pc's of "Takes"

So this looks linear adaptive TM soln, (except for correlas betw. pc's of PST's) If the d.f. on PST's is legit & the SAT effect doesn't invalidate the Analysis!



T. Jours will try to find PST's to analyze that (to compute d.f. of G) that have ~~the~~ large "expected" G's. After a few PST's are evaluated w.r.t. a given problem, it may become clear that we can ~~make a~~ assign a ruff single parameter to each PST giving it "expected G" - This is to tell where to look for "Good PST candidates".

We will look at the G d.f. of the best "PST Run for" and Order PST curves w.r.t. the probly of "Beating" ~~the~~ curve.

ID 779.37 (a preceding discussn) was pessimistic about the scales of PST opten problem: Just what was wrong? Rita ~~now~~ now it seems Ok! It was "T. problem" 746.00 & 771.00

24 T. complaint of 771.07-08 was that there were very complex, A.N. PST's that would work 25 Post corpus problems w/ly G but ~~not extrapolate~~ to present problem. Well, if it was clear that the present problem was not like the past ones - then there is no dilly. ~~the~~ Reading 746.00 ff ... its not easy to see just what the problem was!

10 Maybe a dilly of 746.00ff was to iden that a PST would "look at" a problem & decide on how to solve it: So the G (p.d. curve) of a PST was to be problem independent!

32 797.26 - 40 discuss 2. PST systems: In one, ¹ eval of PST is problem-independent in other, ² PST eval, does depend on problem.

33 Recently, I've been thinking about ²; But ¹ has advantage of the d.f. on PST's, being same for all problems, so presumably less frequent updates (?). In both cases, we update before each problem is "Lsrch" (= "solved").

I may have felt that ¹ & ² were somehow equivalent - so if there was trouble w/ ², there would be trouble w/ ¹. - I'm not so sure they're equiv, how!

GREAT BREAKTHRU! 1.35-1.36

But Note 18.33!

Perhaps 1. ① v.s. ② Question depends on whether I consider Q's to be One Grand QA seq. or a sub. QA seq. for each problem:

It's "Grand" QA seq. One can have "indexing" of Q's to help TM recognize them:

(Hvr, I make conjecture of 2 general problem "Inv, or OZ, or INP": I have TM decide ~~how~~ how to deal w. it. (≡ type ①: 1.32)

In type ② (1.33) T. problem is elementalized & divided into 2 parts @ getting P over PST's as a function of (problem density) Having PST work problem that's already somewhat appropriate to it. The PST's become specialists.

In type ① (1.32) T. PST's become generalists.

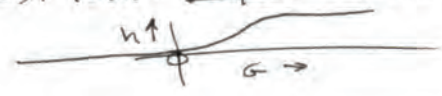
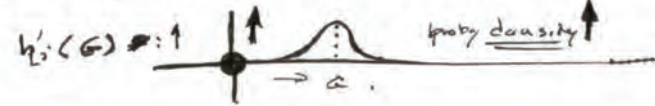
I'm finding it hard to get into the spirit of 77.00 (Sec 1.24) {I'm beginning to understand 1.24} we have PST's w. very low PE (very A.H.) that solve many problems. How so we make expert team So, if we pick one of them, it's a bit how good its likely to be on a new problem - answer would seem to be indep. of what that PST is of low PE.

Then, actually, low PE PST's don't extrapolate well - so fact that they did well on a certain set of problems doesn't make it so likely that they will do well on a new problem.

~~IF STEIN~~ IF STEIN of help here?

Well, Let us first Define the problem(s) clearly!

- 1) We have this infinite set of PST's - most (almost all) have very small errors.
- 2) For each P_{S_i} ($\equiv P_{S_i}$), we have $h_i(G)$. Thus to probly that P_{S_i} will have G is $G_{S_i} < G$ when given a problem of interest.



for OZ probs, the CE is fixed. we want to find for G for P_{S_i} for that problem (prob. den. includes CE).

For a few PST's we will have empirical data on G 's for various probs, T_k

$[P_{S_i}, Prob, T_j, G_n]$ Using this data set, from $\{Prob, T_0\}$ we want a P.D. over all poss $\{P_{S_i}, G_n\}$

Not Quite! May be want $\frac{Prob(P_{S_i}, G_n)}{Prob(P_{S_i})}$ i.e. $Prob(G_n | P_{S_i})$

I think we want $(G_n | Prob, T_0, P_{S_i})$ one or more O^j

So that's it! first find $O^j \Rightarrow \prod_{i,j} O^j(G_n | P_{S_i}, Prob, T_0)$ is a log-likelihood

then $h_n(G) \equiv \prod_{i,j} O^j(G | Prob, T_0, P_{S_i})$ where exactly, $h_n(G) = \prod_{i,j} O^j(G | Prob, T_0, P_{S_i})$

~~2^n~~ $2^n \equiv \prod_{i,j} O^j(G | P_{S_i}, Prob, T_0)$ with mass of all O^j 's.

see 18.33
to counter Argument

35
36

OOPS continuation .16

So how does hr of 12.36 stand up to criticism of 771.07-08 & 11.24-25
Maybe no relevance at all! — But then 12.36 may help solve other problems that were
discussed in 746.00 ^{→ 761.16} & 771.00 ff ← (Reviews: see 771.01 for list of reviews)

On 771.07-08 I was perhaps thinking about SM & Strat evaln.

769.11 was concerned w. "discovery" of good PST's. 771.06 seems to think this was
solved at 769.20-20 & 762.27-40. At present, I can see it as a discovery indeed!

After we have found a good O² in 12.35, finding PST's w. by ~~(2)~~ expected G, is a
Big discovery process!

Actually, finding ~~the~~ good PST's can be interpreted w. getting PC that each
PST is "Best Bet". Once one or more good ~~the~~ O²'s are found.
We begin to look for PST's w. good looking "h₂(G)" curves. — Using some sort of "goodness"
criterion. When we find a bunch of PST's that look good, we can do an approx assumption
of "prob of Best" to Planet. Looking for "Good Looking" PST's (after good O²'s have been found)
is a OZ problem — For it all fits together!

SN Criticism of OOPS: T. Ings mechanism for modifn. of probs of tokens is rather
rudimentary & first order Laplace rule (no how definitions). This criticism holds for
any other "Ings" in the system. In General, he has no way to doing

Combinations of concepts, unless they appear as sols to problems.
The hypothesis any kind of operators that could be added to the system as "Ings" or "tokens".

I think the only trick he now has is "Boost" instruction: So he can search w. an
excessively large no. of insts, then (slowly) solve some problems: T. insts used
& can then be "boosted" to be essentially a subset of insts used for new problems.

3. "Boosting" w/rt. several problem tokens can give ~~new~~ new, useful "concepts"
on the tokens.

Undoubtedly (I think) he only has 2 kinds of tokens: ① primitive insts ② "free" solns
to problems (But check on P13).

J. wants me to be more specific on how I do EPD₁ → GPD₂: Giving ages. 12.35-36
For both INV & OZ may be adequate. — I'd certainly want to do it for my own
reference. Right now 12.35-36 plus ideas on ~~curves~~ curves is how to do (w. Plan),
seems to "finish off" the non-esp. part of TDE! → TM

For TQR, I could adopt a "Don't look inside" pt. of view & try to have TM try to track a
"human" concept sequence — along w. harder ~~harder~~ harder concepts.
T. hard prob!

So: Perhaps Moha & study of hours used in "Early AI": Try to find ways to give a
TQR ~~for~~ to Ira Rean: — Just acquiring having Ph. hours isn't enough: TM must be
able to learn/invent new hours as it moves into new domains &/o ^{over} complexity of
old domains ↑. So far .33 oct. Go over 746-761.16 → 771.00 ff

long time problems

Note 3.00ff on "what to add to Report"

: Sections on OOPS; ID 8 36.07 - 843.12 5.12

SN **Autopce Principal in Statistics:** That ~~you are asking~~ More
are people asking. Q already eliminates many "conceivable"
options. That one should have no cut in t. a propd is ~~at~~
at probability zero is of questionable meaningfulness.

Also, it's such a d.f. would not be of practical interest. ~~there~~
There is no ^{problem} situation in which I could use such info.

→ 52.00

SN on OOPS: If Boost inst has rather by PC, the PC's of tokens will be \propto of f.
freq. w. which they occur in all ~~the~~ old "frozen" parts. Perhaps not a bad idea!
Start out w. an overabundance of ~~the~~ primitives & use this to track back to
pc's for the tokens: Q: This would include freq. of "Boost" - so this is a kind
of vocabulary: It's not clear how it works out! Perhaps if one starts w. a normal
(univariate) pc for "Boost" it will eventual converge to a ~~best~~ value. On the other hand,
if one starts w. a by pc for Boost, it may converge to a different value

(or, depending on the: it may oscillate - ~~or~~ drift in a direction!) ^{depending on the Q.}

Or, the pc of Boost may converge to some value, but very slowly!

Going over OOPS can take a potentially indefinitely long time. Perhaps Best!
Write out 3 or 4 similarities & differences to GATM. Some give some
positive & negative criticisms.

Later write ^{up} criticisms (like a reviewer) to J. (pg. 10-19)

One (apparent) criticism of OOPS: PC's of tokens are not all "Global":
He doesn't really have even context dependence! - He does have some
"long term average context dependence" - But this is not (I think) enough to
deal w. the scaling problem. J means a kind of context: the open paren ((),
you only have certain things following it - but I don't see how ~~OOPS~~ OOPS ~~works~~
is able to get the needed context dependence.

I'm not sure about how OOPS works!

Mention this to J., but not in report -
One possible way for it to get context is via "push pop" & "set pop" (V)
instructions. "pop" is an assignment of PC's to all tokens. In
so: we pop present pop onto stack & set pop # 7) say. ~~the~~ set pop (X)
~~can~~ X can be dependent on present "context".

In GATM, 12.35-36 may deal w. context in an optimum manner.

Re: OOPS division of search into prefix solving all n v.s. finding solving all x (i=1...n)

having = time: Closer to the way the Sci Community does Backtracking (= Recursion)
fraction of an ~~operator~~ solving problem; **fraction of 2^{n+1}** on prefix solving problem; ~~change~~ change may be $\frac{1}{2}$ or
some empirically determined constant. Look in the way OOPS does it: Unpleasant to
just how much memory allocated to various "prefixes".

>0 : **[SN]** I had 2 ^{pass} models for ~~problem~~ GPD \rightarrow [PST] In one, Φ PD was indep of Problem den. In ~~second~~ it was ~~fract~~ of prob den (in fact CB...). My "Gen'l soln" was for ~~case~~ only! See if it works for ~~case~~ as well.

Look at ~~com~~ ^mments I sent **J.** comparing 2 systems.

- 0.05 I think I had ~~3~~ ^{many} PTS:
- 1) OOPS has been ~~planned~~ ^{planned} OOPS had no concept of "optimum's" (links to ~~Q~~, ~~Q~~): Also It didn't work in such problems - but could ~~be~~ ^{the notes posted} if it had GPD \rightarrow GPD
- 2) I ~~Prob~~ ^{did} PC assignment of ~~Edwards~~ better than QATM (But ~~a~~ ^{probably} ~~probly~~ ~~not~~) Primarily true for early versions of QATM
- 3) My TSO's had too large CS - Mahay TSO construction hard
 discuss ~~numerical~~ ^{numerical} problem - ~~20~~ ²⁰ ~~of~~ ^{of} student
- 4) Also, ~~discuss~~ ^{compare} similarities (Frozen Man vs. ~~arg~~ ^{"arg list"})
- 5) use of ~~Stack~~ ^{complex} machine v.s. Function ~~machine~~ ^{Language (Function Trees)}
 $\{$ Very non-critical difference $\}$
 \rightarrow ~~Chances~~ ^{Chances} of not ~~inst~~ ^{inst} TSO $\}$ is of some ~~import~~ ^{import}.
- 6) J's use of "prefix p.g.m.s" $\{$ Is it poss. to have a ~~result~~ ^{result} that \rightarrow ~~get~~ ^{get} "pre" $\}$ \rightarrow 3 list
- 7) Differences in ~~set~~ ^{of initial} ~~initial~~ (Primitive Cons. ~~vs.~~ ^{vs.} Primitive Cons.)
 ~~It~~ ^{It} ~~whilst~~ ^{whilst} it is an ~~impt~~ ^{impt} difference, it is very apt to change.

8) **Re 2** J's choice of pc-modifying mbs is ~~critical~~ ^{critical}; T. Set because is ~~not~~ ^{not} ~~universal~~ ^{universal}.
 I imagine I ultimately use to calculate pc's is rather ~~general~~ ^{general} & certainly ~~universal~~ ^{universal}.
 9) OOPS ~~Do~~ ^{Do} ~~not~~ ^{not} ~~omit~~ ^{omit} (val from ~~it~~ ^{it} errors - it just lens to "amused & good" not "avoided")
 10) Scaling problem of Booting, when there are ~~many~~ ^{many} ~~from~~ ^{from} q's in fusion ~~planning~~ ^{planning} & ~~update~~ ^{update}
 [SN] How to do GPD, ~~when~~ ^{when} some of O's work w. ~~some~~ ^{some} of corpus \rightarrow ~~after~~ ^{after}
 \rightarrow ~~Sec 20.04 ff~~ ^{Sec 20.04 ff}

O's work w. most of corpus only. (Unclear as to what I was thinking about!)
 [SN] Did ~~it~~ ^{fractal of Levin's!} ~~effectively~~ ^{effectively} show that \rightarrow normalize ~~constant~~ ^{constant} for $\{$ semi-measure \rightarrow measure $\}$
 ~~Must~~ ^{Must} ~~be~~ ^{be} ~~un~~ ^{un} ~~bounded~~ ^{bounded}?

26 **[SN]** When I put solns of problems out "Arg list", this is ~~a~~ ^a way to implement "OSL"
 Its not ~~a~~ ^a ~~exactly~~ ^{exactly} correct, hvr. - ~~It~~ ^{It} ~~is~~ ^{is} ~~not~~ ^{not} ~~used~~ ^{used} in OSL! Also ~~the~~ ^{the} OSL ~~use~~ ^{use} may not be legit - depending on size of corpus, size of Ring ~~def~~ ^{def} ~~med~~ ^{med} ~~act~~ ^{act}. Hvr, in any case, ~~the~~ ^{the} way I computed ~~the~~ ^{the} pc's of Tolous is ~~legit~~ ^{legit}.
 So at least its ~~legit~~ ^{legit} (i.e. \exists $pc \leq 1$) - but I may not be getting ~~the~~ ^{the} ~~best~~ ^{best} pc's for ~~the~~ ^{the} best codes.
 On ~~the~~ ^{the} other hand, L_{norm} automatically ~~normalize~~ ^{normalize}, so ~~the~~ ^{the} way I ~~get~~ ^{get} pc's ~~is~~ ^{is} ~~legit~~ ^{legit}.
 Hvr. on ~~the~~ ^{the} ~~22.22~~ ^{22.22} OSL is sort of ~~irrelevant~~ ^{irrelevant}: This doesn't seem to ~~work~~ ^{work} ~~for~~ ^{for} jibo w. ~~26~~ ²⁶!

Initially, the two systems are similar in this respect. The ^{updated} probability of a token is obtained using "Laplace's rule" — roughly proportional to its total frequency of use.
 Both systems augment the set of tokens ~~with a program to solving all current problems~~ by including any programs that successfully solves all problems "up to now".

Alpha has, in addition, facilities for cloning new functions. OOPS has facilities for "boosting" — increasing probabilities of tokens used ^{by or more} in the successful past. In OOPS, the token modification instructions are designed by the trainer, (and are ultimately "reinforced" by the boost instruction.)

After Alpha has had some success in ^{induction} prediction (using a special ^{method} "baker" (6.35L) form of Lsearch), this ~~induction~~ induction facility is used to update the probability distribution for Lsearch. ~~It~~ ^{It} attempts to do this in the best possible way — in view of the available data and available ~~search time~~ search time.

This is probably the most important difference between ~~Alpha and OOPS~~ Alpha and OOPS.

In OOPS, any augmentation of the ^{update} system beyond Laplace's rule, depends upon ~~the~~ ^{token} special instructions added by the trainer. In Alpha,

modification of the ~~probability~~ ^{probability} distribution is obtained ~~through~~ through the universal induction algorithm. ~~Its~~ ^{Its} optimality

is limited ^{only} by ~~the~~ ^{number} of initial instructions ~~for~~ ^{for} the sequence of problems.

The amount of time allowed for ~~the~~ ^{the} induction process. Both Alpha and OOPS are ^{for} conditioned on initial

choice of instruction set and the sequence of problems used to train the systems.

Another Big Difference is "Scaling" (5.19) Imp for early Alpha OOPS & Early Alpha. Alpha has "scaled context" to help deal with it. OOPS will have to add special insts to deal with scaling — may be able to use generalized context.

In Advanced Alpha, the induction-based update algm. can be regarded as a form of Generalized Context. — No expert that it will handle scaling problems. OOPS can ~~deal~~ ^{be able} to deal with scaling problems about as well as is possible.

ommit The more mature Alpha does not directly assign probabilities to tokens! It directly assigns probabilities to PST's that are to be selected ~~via~~ ^{via} Lsearch.

The mature Alpha assigns probabilities to PST's — it does not assign probabilities to tokens directly. — Rough such probabilities can be inferred from the PST probabilities.

In d.ign. of α : mention ^(Section 2, APPA) Stage I & ^(Stage II section on Env, 02) Stage II ^(Phase) ^(Phase)
 Remember we 2 dates: ① α . ② Stage I

GA: .25 : Looks Imp!

Discuss Difference betw A2 long. & "forth". → 20.33

[SN] Th. Self-Confirming hypothesis problem. If a certain PST is fairly successful, TM will give it much wt. on new problems & collect much data on it. — Mechanisms data on other PST's. — So it might tend to not realize that another PST was really better & eventually switch to it. → .16

From short samples of behavior of clear PST's: if they look promising (by induction) TM will give them longer trials. Hrs, a better TM could really "logically understand" ^{i.e. do logical reasoning} PST operation, it probably would not be able to see (recognize) that a certain PST that did poorly w. low CB, would, indeed, do very well w. large CB! → .16

[SN] What would it be useful analogy: Ram job: Uses conventional rockets to get up to critical velocity — then → Ramjet: For QATM!
Phase I → Phase II: Phase III | under "logical reasoning" or "Correlations betw Cands"

Some ways to Elm. Self-Confirming hypothesis (1) random choice of TSP cands.

[2] On PC of PST cands, use $\delta < 1$. $(PC)^\delta$: If we use random search, is it poss. for N to diverge? — what happens if $\sum (PC_i)^\delta \rightarrow \infty$.

If $\sum (PC_i)^\delta \rightarrow \infty$, it will be impossible to create Mt. Carlo trials!

Is there a way to "flatten" the PST d.f. so that normally, least likely cands are more occasionally tried? ^{stochastic}

If we have S CFG; can we multiply them all PC 's by $p \rightarrow p^\delta$ w. $\delta < 1$. This

Will "flatten" the d.f. — Make all branches of closer to $\approx PC$ for each sub-branch? — But may be very close to flattening to entire PD!

[SN] GA for GPD₁ & PD₂ induction using usual copy & 4.14.18 — (i.e. for α)
Also Mut as induction w. ss of $\binom{z}{r}$ — implement N is very GA is 4.14.18. is "usual copy" usual copy is.

Is it true that most any algorithm system that can get "reasonably O.K. induction",

(enough to do GPD₁ → GPD₂) can be used to start off the system, so it becomes eventually (soon?) indep of the original induction system? — no longer useful.

for such a system. Well, it finds a good O^j to generate $N(C)$ pd's (GPD₁)

Then → BPD₂ to choose/generate good PST's. → 19.07

[SN] Re: + "Soln" to a soy problem of 12.35-36: T. counter exp. to "best looking

cand" is the best choice is: we have 10 random vars, w. zero mean & $\sigma = 1$.

we know 11 random variable w mean $\mu = \sigma = .01$. If we pick the best off. "1", first 10 will almost always be the best — yet Var₁₁ is best on average.

I'm not really sure that 33-36 is really relevant to 12.35-36. — There were

00 actually interested in which ~~method~~. t_i prob of a given stand, having ~~max~~ G (or ~~what~~ whatever)
 In 18.28-36, If t_i distributing over of "G", we would be interested in what GPD_2 gives us: t_i variables will indeed be very unlikely to be "Max", even tho it has a higher N than t_i rest of t_i variables.

06 In SM applic of SOY, we want t_i variable w. best reason, \hat{P}_{13} is quite direct from 12.35-36.

07:14.32: T. "Good PST's" would (initially) include P_{13} "OK induction" of 18.28 that got t_i system "Up to speed". Eventually, other PST's would ~~be~~ found that would be better t_i pc of t_i original "O.N. induction" would fade into \emptyset .

I may have a bit of trouble explaining just how t_i updating process automatically, simultaneously ~~re-evaluates~~ re-evaluates t_i of PST's & "creates" new ones. The O^i for GPD , give z pd on all past PST's. $\in T_i$ "creation of new PST's" is an additional pm. That ~~is~~ Given t_i O^j that gives z pd on G of a PST $_i$, to find a PST $_j$ that is "better Good" — This is t_i process of "discovery" of good PST's. It amounts to an OP problem.

Sections of report on $GPD_1 \rightarrow GPD_2$:

17 P12 §2: INV prob. elim linear 2,3,4 of P13

we took top of p 13 (2nd line):

Atan " $(\tilde{G}_i, s_j, F_j(\cdot, \cdot))$ as Q_i and t_i as A_j ."

Using search, ~~we~~ ^{probabilistic} took for functions, O^i , ~~such that~~

21 $z_i = \prod_{j=1}^n \tilde{G}_i^{A_j} F_j(\cdot, \cdot)^{A_j} O^i(t^j | (\tilde{G}_i, s_j, F_j(\cdot, \cdot)))$

24 is as large as possible, z_i ^{being} the a priori probability of $O^i(\cdot | \cdot)$.
 The j summation is over all of the known $Q_i = (\tilde{G}_i, s_j, F_j) - A_j = t^j$ pairs.

"From t_i paying t_i probability density $h'_i(t)$ that $F_i(\cdot, \cdot)$ will solve the problem at time t_i ."

29 $h'_i(t) = O^i(t | (\tilde{G}_i, s_i, F_i(\cdot, \cdot)))$

30 $h_i(t)$ is the probability that F_i will solve \tilde{G}_i, s_i in time greater than t .

31 $h_i(t) = \int_t^\infty h'_i(\tau) d\tau$

32 "Usualy"

Q: What's better, a soln. w. prob p_1 at time t_1 or " " " p_2 " at t_2 ?

There is a partial ordering if $p_1 > p_2$ & $t_1 < t_2$.

Is $\frac{p_i}{t_i}$ reasonable: it's prob per unit time. (Looks like Least criterion!)

for a curve, $h(t)$ maybe w'd like max $\int_0^\infty \frac{h(t)}{t} dt$ or min $\int_0^\infty \frac{1}{h(t)} dt$

$$P = \int_0^{\infty} h'(t) dt$$

$$T = \frac{\int_0^{\infty} t h'(t) dt}{\int_0^{\infty} h'(t) dt}$$

$$\frac{T}{P} = \frac{\int_0^{\infty} t h'(t) dt}{\int_0^{\infty} h'(t) dt}$$

HA! If it's 19.21 I only include data on $F(1.1)$'s that succeeded, I will get a very inaccurate result!
Including unsuccessful data is difficult! How to do it: we only know for each case, that $t_{soln.} > t_{giving}$. - we know t_{put} in each case, but t_{soln} can be ∞ .

One poss. way to deal w. this: instead of wanting large $h'(t_i)$ at the time of soln, for a 2nd problem, for each factor in the product involving unsuccessful failure to find soln, opt till time t_{ol} we want max $\int_{t_{ol}}^{\infty} h'(t) dt = h(t_{ol})$. for unsolved problems.] i.e. this is f./probability "Failure to t"

Does this problem occur w. OZ probs? I think not in general (usually!) Pro for less than caution CC's, a PST could give no G at all. **No G** could be "or G = -∞" one of the possibilities that the prodn. function considers. - say for small CC or for PST's that are not much good for the kind of OZ problem being solved.

So OZ is close analogous to IND in this respect. - "No G" corresponds to (for G = -∞)
"No soln" to a given cc.

O.K. - so we have this $O'(1.0)$ funct. from 19.19:

As an OZ problem, we want a funct that looks at $O'(1)$ and to problem $[G, s]$ and finds $F_0(1.1)$'s w. $h'(1)$'s $\rightarrow 1.00 \frac{T}{P}$ is as small as possible.

Actually, what we want is a batch of F_0 's w. large $\frac{T}{P}$'s & we then compare their $h'(1)$'s to get probab of each being "Best".

As a result ²¹⁻²³ is not a "normal" OZ problem. Many kinds of PST's solving OZ problems give a seq. of excs of by G. PST's of this type are appropriate for 21-23. But there are other PST's that give a single output after f. CB has been attended Such PST's are less appropriate here.

necessarily
This is a standard OZ problem. Not an induction problem.

Next: try to write this into report in intelligible form.

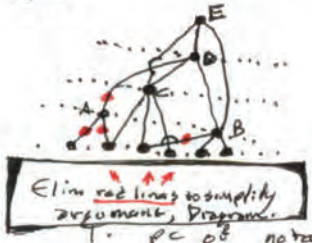
SN One reason I use "AZ" rather than "font", say! AZ is a "Fractional font", so subt over meaning. What this in fact would be "Memorable"? - well Progs in Direction! This is an input Q. A language is good for TM to content that copy in its forms are easily detected. Actually it's "most common copy" that need to be easily detected. Any copy, no matter how obscure & difficult to detect, in any computer lang, is a copy. - but it's hard to detect, it is accorded low attn.

21.00 is about 1.00 - 2.00; "uncompleteness" of ALP
22.00
Spec

0 (SPACE 20.40)

Re: "Growth Dictionary". In both forms: AZ, we can write a function & use it directly, or define it & use it. Normally, Definitions should not be made unnecessary unless
1. Thing defined has accord + better.

What has to be done: In constructing new trial (ex: we (somehow) have to randomly try subfunctions (of the "function thuster"). This is equiv. to OSL. ^{one shot try} → But see (2.2) !



Suppose D was an initial problem soln. form. primitives: B, C, A, B, C were used, they were not defined (not in dictionary or in "Argv list" or in "free zone")

We would like E, to be a function of D (defined) & C & B.

By "retro defining" C & B. ^{any part of node D}

One trouble in this is that pc calculations after C & B were first constructed (but

not "defined" will be incorrect if we "retro define" C & B. (2.8)

Actually, I may be rarely interested in OSL. This is because in the update phase,

TM knows to write answer, so SSZ is at least 2. In the prediction phase of TM operation,

if it does not use OSL, it will be less good at prediction, but this could be fixed up

with A.H. for periods in which we have lots of time for produ. phase. Normally, produ. takes

little time, since no updating is done at that time.

But see 15.26 for what looks like Counter Example

The point is, here, that w.o. OSL, TM's updating will not be defunct from with OSL.

Tho, for actual prediction, & this is often important, OSL is important.

2 → For Long term Growth of TM, OSL is ABSOLUTELY irrelevant. I am mainly interested in long term Growth. After TM has grown up and, ~~we~~ can give it the problem

of "What's the Best way to do OSL". For long term Growth, TM needs to do prediction.

(The pc's assigned to corpus on update) ... This is all to checking & needs on TM's progress

28 .12 → An approximate folly: when E refers to C, its pc is twice what it would be

otherwise: T. Reasoning: In deriving E: we mention "D" which has been defined!

C and B have not been defined, but if they had been, we have to pc of deriving them which

is already covered by pc of "D". D's pc for C & B is w/ by $\frac{1}{\text{total number of bits in } D} (\approx \frac{1}{2}) = \frac{1}{2}$.

So we get $\frac{1}{2}$ factor. E's pc \rightarrow ~~is~~ (pc of D) w/ by $(\frac{1}{2})^2$ since C & B have accord twice. So ... (This probably needs to be done but I don't have time now!).

0: (Saves) 20:32 : Starts w 19.17 : From 19.32 Usually to Deriv of data 3 parts. - They:

01 Given ~~F~~ a good O^2 function, there are an infinite number of $F(\cdot, \cdot)$ functions ~~that~~ it could be ~~for~~ which ~~could~~ be ~~the~~ $h'(t)$ distributions.

We would like to focus our attention on F 's that are likely to be "good".

A rough ~~figure of merit~~ ^{figure of merit} is $\delta = \frac{\mu}{\sigma}$. It is something like the ^{expected cost} probability of solution ~~of solution~~ ^{of solution}.

~~of solution~~ and is similar to the ~~conceptual~~ ^{conceptual} ~~jump size~~ ^{jump size} of ~~the~~ solutions ~~to a problem~~.

~~z~~ $\approx h(0)$ and ~~z~~ $\mu \approx \int_0^{\infty} \text{~~h'(t)~~ + h'(t) dt} / 2$

so $\frac{\mu - h(0)}{\sigma} \approx \frac{\int_0^{\infty} h'(t) dt}{(\int_0^{\infty} h'(t) dt)^2}$

We will call this ~~the~~ ^{the} figure of merit of ~~the~~ F .

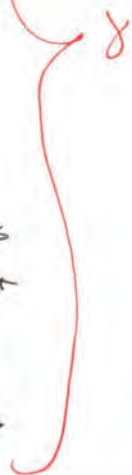
Given ~~the~~ ^{the} function O^2 , the problem of finding a set of F 's with ~~the~~ ^{small} δ values is a time limited optimization problem of the type that ~~is~~ ^{is} solvable by the methods of ~~sketched~~ ^{sketched} in the next section.

Most time limited optimization techniques give a sequence of trial solutions of increasing merit ~~that~~ ^{that} can be used in the present context.

18

0 \rightarrow ~~Given a set of~~ ^{Given a set of} good F 's and their associated $h'_g(t)$'s, we can obtain the probability that any particular F will solve the problem faster than any other in the set. This

R_{13} is text of β below eq. (7)



Ray Nil add. tax

To be added to

nips02.Tax

- date 1.6.03. 5.30 PM

$\alpha \in 19.20-24$

Using L search, we look for probabilistic functions O^i such that

$$z_0^i \prod_j O_j^i(t^j | (\tilde{G}_j, \tilde{s}_j, F_j(\cdot, \cdot)))$$

is as large as possible.

z_0^i is the a priori probability of O^i and the \prod_j summation is over the known

set $\mathcal{Q}_j = (\tilde{G}_j, \tilde{s}_j, F_j)$, $A_j = t^j$ pairs.

$\beta \in 19.29-31$

$$h'_R(t) = O^i(t | (\tilde{G}_n, \tilde{s}_n, F_R(\cdot, \cdot)))$$

$$h_R(t) = \int_T^\infty h'_R(\tau) d\tau$$

\Rightarrow $h_R(t)$ is the probability that $F_R(\cdot, \cdot)$ will solve \tilde{G}_n, \tilde{s}_n in time greater than t .

$\gamma \in 23.01-12$

Given a good O^i function, there are an infinite number of $F(\cdot, \cdot)$ functions for which it can obtain associated $h'(t)$ distributions. We want to focus our attention

on F 's that are likely to be "good". A rough "figure of merit" is $\phi = \mu/\alpha$.

$\rightarrow \mu$ is something like the expected cost of solution and is similar to the "conceivable jump

size" of solutions to problems. We want ϕ to be as small as possible.

$$\alpha \approx h(\phi) \text{ and } \mu \approx \int_0^\infty t h'(t) dt / \alpha$$

$$\text{So } \phi = \frac{\mu}{\alpha} \approx \frac{\int_0^\infty t h'(t) dt}{(\int_0^\infty h'(t) dt)^2}$$

Given the function O^i , the problem of finding a set of F 's with small ϕ values is a Time Limited Optimization problem of a kind that is solvable by the methods

of the next section. Most time limited optimization techniques give a sequence of solutions of increasing merit that can be used in the present context.

The behavior of Alpha is analogous to the operation of a ramjet engine. A slower, less efficient propulsion system is used to get to a critical velocity - at which point the ramjet begins to operate and the initial propulsion system becomes unnecessary.

Phase 1 is needed only to get to Phase 2 - at which point the updating techniques of Phase 1 ~~become less and less~~ are used less and less, ^{would be} i.e.

A variety of induction systems ~~is a departure for Phase 1~~ ~~in addition to the systems described in Section 1~~ ~~could use Genetic Algorithms~~ ~~are some possibilities~~ ~~of a modified form of OOPS~~ ~~system for~~ ~~learning in OOPS~~ ~~is not oriented toward induction~~, but ~~it could be~~ ~~fairly modified to do~~ but ~~it could be very easily modified to work~~ induction problems ~~of a nature~~ ^{in various problems only,}

Wraith
1x200

Though OOPS ~~is not~~ ^{is} ~~designed for~~ ^{in various problems only,} it could be very easily modified to work induction problems - as a possible implementation of Phase 1.

OOPS and Phase 1 of Alpha are very similar. They both use universal distributions to solve problems using Lsearch.

Footnote: Schmidhuber's terminology is incorrect: what he terms "Osearch" is really Lsearch; what he terms "Lsearch" is really "SIMPLE", a ^{minimally complex} program devised by Li and Vitányi; to illustrate an important characteristic of Lsearch - i.e. its ability to solve all solvable inversion problems within a constant factor of the speed of an optimum solution. The "constant factor" to SIMPLE is ^{however,} much larger than that for Lsearch.)

The Language used by OOPS is a ^{kind of} ~~stack~~ stack language, related to FORTH.

The Language used by Alpha is AZ, similar to LISP. Both of these languages ~~can~~ can use distributions to compress code - ~~which is an essential feature in~~ incremental learning.

~~AZ was designed so that~~ For AZ is a functional language: functions are represented by trees, and all subtrees are legal functions.

In AZ finding common sub-trees in a large function, enables compression. In the language used by OOPS ~~we~~ ^{guess that} ~~we~~ ~~think~~ it is less likely that common sub-trees (that had not already been defined) would occur with usable frequency, an analog of this technique would often be useful.

Both systems use incremental learning to update the ~~probabilistic~~ probability distribution that guides search. ~~to most parts, the update systems are similar,~~ while the update systems are similar. Alpha uses context of various kinds to deal with scaling effects. ~~There are uncertainties as to~~ how OOPS ~~addresses~~ a proaches scaling — though apparently, it is a serious problem. OOPS was able to solve the "Towers of Hanoi" about 1000 times faster, by using ~~search~~ a search pattern (a kind of probability distribution over ~~candidate~~ problem trials) from ~~the~~ the solution to an earlier problem. Since there weren't many "earlier problems", it was ~~relatively~~ relatively easy to find the appropriate pattern. If the system had solved 1000 problems before the "Towers of Hanoi", finding the correct pattern ~~would~~ may have taken 1000 times as long — ~~thus~~ thus canceling out the gain in search time.

The most serious difference between OOPS and Alpha is in OOPS not having any concept of optimization ~~experimentation~~. It has nothing corresponding to Phase 2 of Alpha. Phase 2 has many useful features. ~~It is able to use its~~ ~~updating~~ ~~scheme~~ ~~to~~ ~~improve~~ ~~its~~ ~~induction~~, which further improves the updating scheme, and ~~it~~ ~~is~~ ~~able~~ ~~to~~ ~~use~~ ~~information~~ ~~about~~ ~~failed~~ ~~trials~~, as well as information about success.

It is able to use information about failed trials, as well as information about success. ~~It used data from previous problems to assign probabilities~~ ~~It is able to~~ ~~invent~~ ~~new~~ ~~PST's~~ ~~and~~ ~~assign~~ ~~probabilities~~ ~~to~~ ~~them~~ ~~with~~ ~~respect~~ ~~to~~ ~~an~~ ~~attempt~~ ~~to~~ ~~solve~~ ~~a~~ ~~particular~~ ~~problem~~.

Another difference is that Alpha is largely a procedural system. The behavior of ~~Alpha~~ ~~is~~ ~~a~~ ~~variant~~ ~~of~~ ~~Phase~~ ~~1~~, using a stack-based language ^(much like OOPS) has been analysed for the ~~sequence~~ ~~of~~ ~~problems~~ ~~involving~~ ~~learning~~ ~~of~~ ~~evaluation~~ ~~of~~ ~~algebraic~~ ~~expressions~~ (Pau 1994). ~~The~~ ~~scaling~~ ~~effects~~ ~~were~~ ~~observed~~ ~~when~~ ~~we~~ ~~computed~~ ~~the~~ ~~CJS's~~ ~~of~~ ~~problems~~ ~~of~~ ~~increasing~~ ~~difficulty~~. While some analysis was made of more difficult problems no CJS values were computed.

OOPS, on the other hand, has been realized as a computer program — demonstrating the importance of early learning in facilitating solutions of difficult problems.

It really tries to find the optimum solution to a problem, rather than an improvement over the ~~best~~ best possible trial.

little balls?

10

10

20.04-19 is an important simplification of the system; Consider 14 v. probs (20.04-19)
 "h'(t)" can be divided into 2 parts.
 Can I somehow integrate this into the idea that

^{0.2}
 A correction for P.13 of PS document

Suppose $O^i \in (0,1)$ is a probability density function of the type discussed in section 1.

~~From definition~~
 ~~$h'_{i,j,k}(t) = O^i(t | G_j, S_j, F_k(\cdot, \cdot))$~~

define $h'_{i,j,k}(t) = O^i(t | G_j, S_j, F_k(\cdot, \cdot))$

This is the probability density (according to O^i) that F_k will solve G_j, S_j at time t .
 i.c. $\neq t$.

$h_{i,j,k}(t) = \int_t^{\infty} h'_{i,j,k}(t') dt'$

is the probability (according to O^i) that F_k will take longer than t to solve G_j, S_j .
 than t to solve G_j, S_j .

Then ~~using~~ O^i 's such that

$z^i \prod_j h'_{i,j,j}(t^j) \prod_l h_{i,l,l}(t^l)$

is as large as possible.

z^i is the a priori probability of O^i

The j product ranges over those cases in which F_j has solved G_j, S_j at time t^j i.c.

The l product " " " " " " " " has spent time t^l but has not yet solved G_l, S_l .
 i.c.

It is notable that this update scheme uses data on past failures as well as successes.

usually there is not enough data. ... (omit subscript on h, h')

$h_e \rightarrow h$, etc. (output altered) it says: "Partial data will often demand $W_{new} < 3$ params"

also note <http://cyc.com/>
 D. Lohat "Cyc: A Large Scale Investment in Knowledge Infrastructure" Communications of ACM 38, no. 11, 1995

Communications of ACM 33(8):30-49, August 1990

refs
 M. L. Cramer
 J. Schum 2002
 D. Lohat et al Cyc: Toward Programs with Common Sense? <http://www.sover.net/~michael/hlc-publications/cyc285/index.html>
<http://www.ldsia.ch/~juergen/oops.html>
 Optimal Ordered Problem Solver TR IPSIA-12-02, 31 July 2002 J. Schmidhuber.

W. Paul; R. Sol. 1994 Paul, Wolfgang, Solomonoff, R. Autonomous Theory Building Systems, Annals of Operations Research 1994.


D. Lohat. D. Lohat, R. Gupta Building Large Knowledge Based Systems Reading, MA, Addison Wesley (1990).

0 } "But on p. 13: We have a quadruple list for each problem solved.
2 } We also have an identical looking quadruple for problems not solved. However, it is not true that
before the trial was ~~discriminated~~ discontinued.

0 Revisions: 26 Nov 02 original report
9 Jan 02 : Section 2.1 on "Improved Update system" augmented
Section 6: "Related work" on Lemarié OAPS added,
additions to Bibli. & ref. material, for Crowder added.

Mailing list for 1.9.03 report

- 1) Kurt Alan Stein Kraus kurtas@ai.mit.edu
- 2) Ivan Chardin chardin@media.mit.edu
- 3) Alex Gammerman alex@cs.tul.ac.uk
- 4) Natalia Krasnogor NataliaKrasnogor@nottingham.ac.uk
- 5) Gerry Wolff gerry@informatics.bangor.ac.uk
- 6) Doug Campbell ~~██████████~~ dugie^{→ "sto" "u"}@yahoo.com tmq.
- 7) Tom English tenglish@jam.rr.com
- 8) Henry Lieberman lieber@media.mit.edu
- 9) Oliver Selfridge ogs@media.mit.edu
- 10) Phil Apley PGA@ALIENLANDING.ORG
- 11) David Dove dd@mail.csse.monash.edu.au
- 12) Sapp Hochreiter hochreit@cs.tu-berlin.de
- 13) Eric Winfree winfree@hope.caltech.edu
- 14) Steve W. Yam SW@tiac.net
SW@tiac.
- 15) Marvin Minsky minsky@media.mit.edu
- 16) Carlos cazhv@nyc.rr.com 1.16.03
- 17) Jess Morton jmorton@igc.org
- 18) Lewis Morton lm@lmc.demon.co.uk
- 19) Will Gersch ~~Gersch@Hawaii.edu~~ gersch@ics.hawaii.edu
- 20) Murray Denofsky mur252
mdenofsky@yahoo.com 1.16.03
- 21) Alex Sol als
- 22) Cosma Rohilla Schalizi Schalizi@santafe.edu
- 23) Kirby Smith kirbysmith@mindspring.com | try ksmith@ieee.org.
- 24) Irv Wieder WiederSci@aol.com
- 25) Marcus Hutter macus@idsia.ch
- 26) Juergen Schmidhuber juergen@idsia.ch
- 27) Richard Gardner rdl@rtuh.com
- 28) Gerald Sussman gjs@mit.edu 617-253 5874
- 29)

00 : Call to system of 1. Nov 9, 03 report  ALPHA). (ALP+H).

Consider Alpha as Phase 1, Phase 2. At first glance, we have much clearer how Phase 1 works: It does have to be able to do induction, to extent of finding

O^2 \rightarrow eq 5.5 (the first eq n ≤ 2.1) is "by and"

05 **[SN]** Hur, even Phase 2 is not the "final + H", we have to deal w. "corollaries" betw.
06 pc's of PST's, I'm not sure that even that is the "final system". But not (.11-.12)

But anyway, in Phase 2, the system is able to ~~the~~ derive such things like:

09 **[PST]** ~~set~~ ^{set} - which (I think) make explicit consideration of "Context" / predecessor,

11 - Also (I think) Phase 2 is able to replace itself - (the ultimate "wrapper").

12 If .11 is true, then the "context" of .09 becomes no longer necessary, - also "corollaries" of (.05-.06) would be taken care of.

14 **Re OOPS:** Use of prefix codes for causes: The way it supposed to work:

A particular prefix code q^i is able to work all problems ~~in~~ thru n^{th} . After working n problems, it halts. So, even the problem in put, then, the q^i 's form a prefix set.

For a new problem not in n (thru n set), the q^i 's may no longer be a prefix set.

20 This seems like a nice way of being sure that all new trial q^i 's will ~~not~~

work ^{through} $(1+n)$; if we make all new trial q^i 's; old q^i 's w. post fixes on them.

Hur, I should think that "backtracking" would work as 14.35: (spend time ~~on~~

$\propto \alpha^{j-n+1}$ on q that work prob; thru; ~~on~~ $\propto \frac{1}{\alpha}$)

What α should be is unclear. Perhaps find out empirically?

Actually, the function & mechanics of this "Backtracking" are unclear! Say you never Backtrack.

Would it be possible to somehow "go on", w.o. revising to "General program this far"?

30 **[SN]** Perhaps it would be well to work on how the parts function, since

Part of P_2 is Goal of Phase 1. A Bern seq. can be defined by

31 ② ~~to~~ n steps (i.e. overlap rules) ③ non-deterministic ④ homomorphism.

.31 is not so clear! Perhaps Bern seq. for a radix k alphabet is clearer. List symbols in alphabets, followed by concatenations.

[SN] Jess, M. \rightarrow this can it split into / recurrent, repetitive loop of "improvement of GPD", then improvement of induction \rightarrow system \rightarrow improvement of GPD \rightarrow and ? E.g. for improvement "read" "modify."

40 **[SN]** Also ~~PS, FP3~~ \rightarrow prop $\rightarrow 2 - 2(0^2)$ & is length... was unclear to Jess a "naive reader"

empt. Argmt.

or more exactly
A kind of BLP: A kind of Comp. Prob
A way of trying models before CB occurs

— e.g. A method other than RLP

00

On: Q of 1.00 - 1.40: Any computable prob. distrib. can be regarded as an approximation ALP - which is to "True" p.d. If there is a code for the data that is much shorter than that given by this computational, then the error will be very large.

.05

Whether such a much shorter code exists is unknown: \Rightarrow whether is how large the error in the c.p.d. is, is unknown.

.06

Counter argt. "I don't agree that ALP is really the correct probability".

f

v.s. (06): (.00-.05) can be rephrased by saying: Whether there exists a "much better model" for the data is unknown. More exactly: we can never get to the point where we can say "There are no significantly better models"

0

Thermodynamic entropy difference between state A and state B $\equiv \int_A^B \frac{dH}{T}$

dH is incremented heat added in path from A to B. In a solid box,

$dH = dt \cdot c$ (since no. of molecules) so $\int_A^B \frac{dH}{T} = \ln H_B - \ln H_A = \ln \frac{H_B}{H_A} = \ln \frac{T_B}{T_A}$

How ^{Thermodynamic} Entropy relates to "Informational" entropy is unclear.

In Cover's book, he defines Thermodynamic entropy as \ln of no. of poss. ^{micro} states of the system.

However, he expects probabilistic transitions from one state to another, so

But anyway; If one knows Temp, press, vol of gas, one can compute its entropy change to new Macro state.

Actually in Thermodyn., If T, P & V are constant,

On S functions: "many" different forms? (1) 3 input vars (2) Bernoulli params: numerator, denominator

(3) explicit functions $P(A|Q)$, for ~~normalized~~ A, Q , ~~or~~ $P(y|x) = F(x)$ (this is a normalized pd for unnormalized) but $\int f(x) dx < \infty$. $f(x) \geq 0$, usually $f(x) > 0$ strictly.

(32) a common form $F(x) \rightarrow f_1(x), f_2(x)$: i.e. we express 1 or more parameters of the d.o.f., f , as functions of x . More exactly: Q is of form S, x ; S is a string that refers to the problem, S is the continuous set of params in the problem. "S" may say the params are the coeffs of a ~~linear~~ Taylor expansion, or any other parametrized approx.

Criticism of OOPS: Its method of modifying f.p.d. doesn't look universal to me.

(My "Phase 1" isn't universal either - which is one reason I have Phase 2)

OOPS gets its f.p.d. on functions by modifying f.c.'s of tokens and defining new tokens
Hvs. actual limitations/defn of what kinds of f.c. module is allowed, is unclear.

At one ^(time) I was uncertain about whether a system I had, had this limitation - I think I finally got a system that was universal (in the usual code way).

ON SUMACS: A Sumac ^{has} 2 imp (Characteristics) functions and prediction and updating:

The base of each is not inputs & f. accuracy of growth. (which interacts w. t. "accuracy (per unit time) of updating.)

So Study various Sumacs wrt. these characteristics. 2 or 3 come to mind:

- ① True ~~sumac~~ based CB or Sumac
- ② ① with finite CB: Problem Backtracking is possl.
- ③ AZ-type Lang: Bernoulli f.c.'s of "Tokens"; New tokens desirable: Backtracking f.c.c'l.
- ④ OOPS language: Stack based.

Remark: In General, Backtracking is an attempt to deal w. fact that CB for each sumac is finite - so we don't know all codes for "Complex Problem".

In ③ or ④ if we had CB = ∞ , would we have a universal d.f.?

In ③ I admit "oversearching", so any code, is eventually findable! A rich lang is indeed universal. T. answer to .19 is "Yes"

In ④, I guess f. lang is universal, but he doesn't allow oversearch, so "No" for that specific case. If OOPS allows "oversearch", then .19 would be "Yes".

In both AZ & OOPS, we are considering D-induction only.

Since OOPS uses lists, given enough time, it will find any short code possible.

Having found that short code, does it extrapolate well? I think that OOPS is pretty much to same as Phase I Alpha: It finds short function pieces to solve the problem.

This is not a conv. program so perhaps T. Conv. Program of Appendix B applies to OOPS - It has over solved enough!

Hvs. his problems are all different types: So, for the $2^4 b^n$ problems & examples have to be indeed; Similarly for f. sources of hazard. - I think we may do quality QA problems & conv. program holds!

The OOPS doesn't "oversearch", it does "backtrack", so T. Conv. Program may apply - (no maybe is "less exact way") ; Q: How does Conv. Program work with "D-induction" - BTW spec

00 : Since "Backtracking" is impl, its impl that we do it via: My impression of a ^{possibility} ~~good~~ ^{way} ~~is~~ if P_1, \dots, P_n are k successful perms, then it spends fraction α on another Q_n .
 fraction α of remaining time on Q_{n-1} , fraction α^2 of remaining time on Q_{n-2} ... etc. Its not clear what α should be, or what constitutes "A problem" (so fraction α should be α)
 Is + success α^k set a "single problem"? Is P is something for "Trainer" should decide?
 (i.e. is the time distribution a part of the TSO definition? — it could be.)

→ COPS doesn't give
 into pc's to
 successful perms
 Q_2 .

10 : W.O. special "pc modifying insts" (other than derivations & Bernoulli pc assignment)
 We should get good pc's for functions. Yet his "Bart" inst. was critical in solving the "Towers of Hanoi" problem. So clearly certain insts can be very implicit in working certain kinds of TSO's. These insts can be "Heuristics"

12 : Hvr, for a TM to discover heuristics: It has to have "Requisite Variety" of probs with TSO.

12 : 33.40 For D-induction, $\text{error}(pc)^2$ is either 0 or 1, — so we count no. of errors. T. expected no. of errors is bounded — indep of corpus size, if there is a "short code" corpus that doesn't ~~grow~~ grow w. corpus! (corpora \neq (word))
 This D-induction may be related (or identical) to Chaitin's ideas about how much math one can get from "how much axioms". The "errors" are the axioms (?).
 Anyway, if one has a large corpus w. no errors, i.e. deriv (= axioms) have n bits, one could have n errors in the program's applics, (ways) — its "Expected no. of errors!"

20 : Hvr, "Expected values" when $p=0$ or 1 are quite different from statistical "Expected values!"

I should end up w. a Math Program (not a statistical Program!).
 Maybe put Conv. Program in RL form — (is stronger) than $p \rightarrow 0$ or 1 .

Actually, unnecessary to go at it from the Conv. Form — I think its poss. to do it directly.
 Consider sequential predn: (D-predn). we have a long sequence ~~described~~ ^{described} by k bits. Deriv. is "u". We have another ^{deriv} k bits long that describes a corpus of length n , created by u .
 (u can create a corpus of arby length $(=\infty)$). Say k has correctly predicted n bits of u . Say i at $0 \leq i < n$, there is some "productive" gamma.

The applicability of the Conv. Form to D-predn. — ~~both~~ ^{both} (sequential, bay induction & A induction)

30 : is interesting, I don't want to spend time on it now. Maybe I will for paper on "Convergence Theorems"

Anyway: what I was working on, was the Q of just how good COPS was.
 W.O. special pc modif. instructions, it seemed like it would give a regular Universal D.F.
 The pc modif. insts could give a ^{highly} ~~highly~~ biased Univ. D.F. — T. Q is — would it be a good Bias?

Re: SOMACS. By limiting Backtrack depth, we make it impossible deal w. certain "errors" (\equiv find certain Regys). What kinds of Regys are these, and

20: a new problem: It may give output (usually similar); But extending to codes of Progs will usually modify some responses to $\frac{\text{older}}{\text{problems}}$ in S_0 .

Essentially, there are 2 parts to understanding OOPS: 1) The Mechanics of "TRT".

2) how the Language works: how it recalculates new pc's for tokens, etc.

One imp. idea is use of prefix codes, to reduce time spent on decoding new codes.

It may be that by suitably doing // search & ~~avoid~~ duplicating/saving states as needed one should be able to get better than bias optimal of 1; perhaps better by factor of length (in tokens) of solution (?).

Is the OOPS formulation more amenable to expressing S functions?

On S functions: for continuous PD's: say I have a $\begin{pmatrix} \text{arbitrary} \\ \text{uniform} \end{pmatrix}$ pd on $(0,1)$ interval.

Any monotonic function $w: (0,1) \rightarrow \text{Domain}$, will give us a pd on its Range. What I (have) use \rightarrow that normalized pd's \rightarrow normalized y.d's.

I can probably translate this to discrete d.f.'s.

In Continuous case: for each x , we have a (defmt) function. $\text{sort}(X, y)$ a difmt. mapping from x to a common (?) domain to a common (?) Range. It is y 's domain.

What we want for discrete case: each $x \in Q$ gives a difmt (usually) d.f. on discrete A space.

For each Q we have a discrete P.D. over A space. In general, A space covers all possl. Discrete objects. So it is a p.d. on all sings (finite alphabet) - It could

be a grammar. It could be a Universal D.f. - perhaps it would be desirable for it to always be a univ. D.F.

NB. Just because a D.F. is Universal doesn't fix it down very much! All ∞ grammars are Universal, but they evolve a preorder as $\#TSC$ grows!

A Universal grammar can start w. a small number of primitive functions & a discrete P.D. on Param . In AZ I have a (I hope!) univ. d.f. on functions. How can I modify AZ to get a d.f. on strings? Well, OOPS does that! OOPS also asks for a new token occasionally.

Perhaps use OOPS formalism to create general script \rightarrow Universal P.D. - i.e. to generate all possl. P.D.'s.

Perhaps have a machine that's a mix of AZ & OOPS! $Q, R \Rightarrow A$

Q input is AZ ; R input is OOPS.

Note, w/ 3 input one; one need only have 1 input tape: T. machine itself decides w/ boundaries betw. Q & R .

There has to be some convention on when Q terminates: maybe a "stop" symbol.

10:36:40: T. Machine has no output until ~~it starts~~ ~~Q starts~~ Q starts. It can have some output w/o knowing the entire value of Q - but usually (I should think) it waits for the end of Q before starting to "speak". I could have a special output, (that can be ^{erasable}) that tells when T. machine wants another input bit.

35 BUT I think my idea about "forms of S. functions" was that I would express them pretty much the way that scribbles normally express S-functions - from possibly genz. to methods to make them Univl. (if necessary). ^{If they are not already Univl.}

0 So introduce problems w. S-funct solns. that one would expect a person to be able to discover. Then find ways to express these S-functions in a common, general way

I should wait until I begin getting problems of that sort - then write S-functs that seem "Natural" to a human:

-15 Initially, use Bern P.P. for Discrete probn. $\hat{=}$ (radn. of conv) + other continuous functs.

Re: My Understanding of α is cap!

1) I think I understand α , (phase 1) w. one R-funct. > 1 R as I've written it may not be so good, using R's w. pc's that overlap might be better.

α_2 is O.K. to some extent "In Procy", but I don't have good genz. (form for S-functs (the seq above .05-.15) also for "correlation bands" is a problem - later: "Clustering" may help. - Cross corr between some bands may help.

2) Cops Understanding Cops has several parts:
 a) General lang. used. How PC of tokens are changed; ^{When does} ~~how~~ system requests now to learn? etc. : As a sub-Q: Just what are good sets of primatives? (Q's last is also a Q in α_n)

b) How does Cops solve a bunch of problems? or more exactly, how solve a bunch, how does it solve a new one? My tentative understanding:
 2 methods, ^{used} ~~two~~ time share!

(I) try extensions of q_n (~~that~~ q that solve all probs)

(II) start w. a bunch of $(n+1)$ a bunch $\hat{=}$ try to solve them in Π . T. only lang used

is the pc's of tokens plus t_i frozen solns, q_1, \dots, q_n . Methods of "editing"

(III) "Modifying, comparing" f . q 's must be derived! J. expects COPS to

learn how to do this - but he doesn't suggest I try to implement it. He may want to insert them "By Hand" - but in many instances ~~to have~~ system "learn" which are erasable

T. feasibility of this work (II) would seem to depend on n . 33-34

Otherwise, the use of prefix code q_i having time share $\propto b^{-i}$ (for some $b > 0$)

looks better than (I): extensions to q_i would try to solve probs $\hat{=}$ $n+1$ in parallel q_i (time share).

MEM:...

6 v.s. 42
1 v.s. 7
Gen Schwartz
Koff

- 1) Unkn about UN
- 2) Don't mix Sci, Politics
- 3) Don't want to appease Eragu.

MAX ENTROPY METHOD



It seems to me that MEM \equiv Max Likelihood!

Say P_i are the ^{unknown} prob of various symbols in a corpus.
 N_i is case count of symbols.
 We have ^{known} constraints on P_i .

Then we want $\prod P_i^{N_i}$ to find a set of P_i , \rightarrow

$\prod P_i^{N_i} = \max$ subject to known constraints. This is Max Likelihood.

It is also (so MEM!)

More exactly, the def of the P_i (w. finite size) is simply

$\prod P_i^{N_i}$ subject to P_i constraints.

P_i gives a precision of the \vec{P} vector.

Maybe not! According to R. Christenson!

MEM says choose $\vec{P} \ni \vec{P}$ s.t. P_i is max, (subject to constraints)
 — so the N_i of .07 would not be "known" — (or used!).

In this case, $\prod P_i^{P_i \cdot N} = \max$ subject to constraints on \vec{P} .

N may be unknown, but "large".

Max cross entropy (see R. Christenson, 1986) is
 $\min \sum p_i \ln(p_i/q_i)$ where q_i are priors of the P_i .

20

21:38:40 More on OOPS! This method of solving Problem (net) r_{net} : Divide into 2 ways: \pm from each.

- 1) Try to extend q^+ to solve r_{net}
- 2) Try to extend Δ (null) to solve r_1, \dots, r_{net} .

Actually, this is problem 21(L), & it is not obvious that these are the best 2 ways — or that either is particularly good! It amounts to a WDN problem: "Add" net, but w. much aux info — i.e. that we know that q^+ solves r_1, \dots, r_2 .

30

\rightarrow J. says or implies in paper that 21 is only off optimum by a factor of 2!

I see no reason to suspect this is true: It should be paid on it if it has 2 reasons:

[I have, in α_1 , proposed a way of solving this problem: I make no claims of optimality, but it's not clear that (21, 2) is any better or worse. (The J uses a trickier way of a kind different from those used. — (if α_1 is a \pm sign and it stays same place for all steps).]

NIPS

LV book 1993 p 237 Lemmas 4.6 & 4.7 (perhaps)

Does L-V have a proof or discuss?

Also P220 Lemma 4.1 + Discrete P223 introducing semi-measures.

John Levin / Zvonkin perhaps.

It's known result in § 75, this problem is L/Z paper.

ON Normalized Universal D.F. For a continuous D.F.

Normalize then, then there is no one that is $>$ all others (written constant factor).

If μ normalized constants were bounded, then $P_{n,2}$ couldn't be true, so Normal constants must be unbounded (may be not for ≥ 1 corpi, but for "Most"?).

That μ normalized constants infinite has strong effect on why μ Normal form of univ. d.f. is better than unnormalized form: i.e. its error is much less.

look imposs! The expected score bound is $-\ln P(\mu)$ where μ is parameter of f -data.

If $\ln A(n) \rightarrow \infty$ { A is norm constant (≥ 1), n is no of bits in data } so $E \sum \epsilon_i^2 = -\ln P(\mu) + \ln A(n)$.

$P(\mu) \rightarrow K(\mu) \cdot A(n)$

say $p = 2^{-10}$; $A = 2^{+2}$

$-\log_2 p = 10$ $-\log_2(p \cdot A) = -\log_2 p - \log_2 A = 10 - 2 = 8$

$$E \sum \epsilon_i^2 \rightarrow \underbrace{-\ln P(\mu)}_{\text{Bounded}} + \underbrace{\ln A(n)}_{\text{unbounded}}$$

which $\rightarrow -\infty$ as $n \rightarrow \infty$!

Maybe $A(n) \rightarrow \infty$ only if $P(\mu(n)) \rightarrow 0$

So, for

Perhaps Levin should (0.00-0.01) if we considered all methods of Normal, not just μ type I used.

Another Posey: But my proof of Sol 78 T3 is incorrect! That μ conv. form is not true for

μ Normalized Universal D.F.

$\prod (1 + \epsilon_i) \geq \epsilon_i^2 < \infty$

First study ratio of normalized Universal semi measure to any semi-measure (including μ).

$\sum p < \infty$
 $\sum \sqrt{p}$ can be ∞

Univ. semi-measure $\leq K \times$ any (semi) measure w. finite den. + positive value n
Norm. Univ. measure \leq Norm const $\times K \times$ any (semi) measure w. finite den. + positive...

The many semi-measures have lower finite den., f -lengths of f -den. can be very long.

N.B. all normal methods are not equal: some poor ones: ① always chose 1, ② always chose 0, ③ chose 0, 1 w. = probly ④ chose f ~~or~~ if $P(f) \geq \frac{1}{2}$ choose f .

This set is obtained. ⑤ if $P(f) \geq \frac{1}{2}$ chose 1; otherwise chose 0.

Not all methods of normal have normal constants! Only μ one ϵ used by μ normal constant! Anyway: Look at μ Li-Vit refs at top of this page to get idea of why Norm. Univ. d.f. have no "best".



FIRST DATA MERCHANT SERVICES
20 MALL ROAD
SUITE 350
BURLINGTON, MA 01803
www.firstdata.com

00

Things to do in Revision of report.

- 1) Contents; w. some comments.
- 2) Long discussion of differences ~~in~~ Between ops in TSQs Actual differences ~~the~~ approach to the problem.
- 3) Discussion of how to use intro failures.
- 4) Go thru ~~the~~ analytical paper carefully - fix typos, etc.
- 5) Do put in section that evaluates what has been done - what is programmable - what isn't (yet). (see 2 Feb 7, 03 letter to J. Firing our report report - i.e. how to implement s-funct is 2 ~~TSQ~~ TSQ to get Ram to work phase 2 ~~TSQ~~ SPD options.

10

20

30



FIRST DATA MERCHANT SERVICES
20 MALL ROAD
SUITE 350
BURLINGTON, MA 01803
www.firstdata.com

00:00 follow centuries for the scientific community to acquire, and I'd rather not have Per system take that long!

F

G - See my letter of 27 Nov 02 for a discussion of just this point.

H - I have written this up, but have to find good way to insert it into the report.

10

I - Am now working on 2 talks: One "popular" lecture for ^{the} London, Kolmogorov Celebration; Another more technical ~~talk~~ ^{seminar} to be given at Royal Holloway Computer Learning Research Center.

Will try to finish up the report as soon as I can - perhaps after ~~the~~ ^{Per} Kolmogorov Lecture.
Much thanks ~~for~~ ^{for} comments!

19



20

Have been looking at OOPS. You propose several ways to get a selected single solution to ~~the~~ $r_1, r_2, \dots, r_{(n)}$, given q_1, q_2, \dots, q_n .

It is not clear that any of them is in any sense optimum.

In general, L search is a ~~bad~~ ^{good} idea, if the only information you have is in the coding probability distribution. ~~When you have~~ ^{While} auxiliary information, like q_1, \dots, q_n , is in your probability distribution, information on how to use this data is not ~~under~~ ^{under} those conditions, L search is not near optimum. ~~just what is optimum in this case~~

I don't think you should claim that your search technique is optimum within a factor of 2. We really don't know how close ~~to~~ ^{to} optimum ~~the~~ ^{any} of these methods are.

30

At first I thought that you had made a "Great Breakthrough" when you decided to use prefix codes for your programs. A bit of thought made it clear that prefix codes were no panacea. The codes I used in "Phase I" were also prefix codes - ^{but} each legal code ended with a "stop" symbol. ~~Your prefix codes differ from mine, in that - My~~

prefix codes differ from yours in that the prefix set is independent of the agreement. In your codes, there is ~~a~~ ^a (possibly) different ~~prefix~~ ^{prefix} code for each agreement.

0 minutes - I'm uncertain about how to split



FIRST DATA MERCHANT SERVICES
20 MALL ROAD
SUITE 350
BURLINGTON, MA 01803
www.firstdata.com

Answers to J's letter of 13 Jan 03!

A Actually, I think that's eq 6 and associated discussion, is about the only part of the paper where I give examples. I will, btw, make it clear that the definition of E_{α} being proposed is independent of just what α does. Also, I really should have more examples!

B "It has enough skill", when it is able to do what needs to be done in Phase 2 — i.e. in section 2.2, ^{first equation (unnumbered):} it must be able to find good O^i functions. Usually the trainer will judge when ~~it is ready~~: it can usefully try to ~~maximize~~ maximize this expression — just as the trainer has to decide when the system is ready for ~~any~~ ^{it} ANY new problem in its training sequence. If the trainer gives Alpha a problem that is too hard for ~~Alpha~~ to solve, it will not solve it — which was his timer.

C This is certainly a good idea and I will put in section 2.2 discuss these matters.

D The idea of self improvement: ~~is~~ ^{spends half of its time solving} The system normally ~~solves~~ ^{The over} problems using L search, guided by the GCPD. Half of its time is spent "improving the GCPD" by trying to get better O^i values for the first equation in section 2.1. ^{By "self-improvement" I mean} "Improving the GCPD" ~~is~~ ^{is} self-improvement.

The reason you can't implement "self-improvement" is that I ~~didn't~~ ^{didn't} explain ~~how to~~ ^{how to} realize stochastic functions like O^i . In ~~the~~ ^{the} section on ~~how to~~ ^{How} ~~to~~ ^{to} ~~realize~~ ^{realize} stochastic functions and what needs to be done, I will discuss this.

E True, ~~also~~ ^{Also note that} Phase 1/using only one R can do the same thing, ~~but~~ ^{by} building various hidden functions into the single R. ~~The~~ ^{However} in both OOPS and ~~Alpha~~ ^{Alpha} Phase 1 of Alpha, it ~~will~~ ^{will} take a lot of training before it does this well. In the "Multiple R" system, I am telling the system about ~~a~~ ^a very important heuristic device used by search lists, and saving a lot of ~~the~~ ^{the} computation time, ~~which~~ ^{which} I have no idea ~~how~~ ^{is} as to what kind of training sequence I'd need to get the system to discover this trick. ~~Also~~ ^{Several} heuristic tricks have

10

20

30



FIRST DATA MERCHANT SERVICES
20 MALL ROAD
SUITE 350
BURLINGTON, MA 01803
www.firstdata.com

But a solution to r_n is not always obtainable as a consequence of q_n .

The great feature of your prefix codes is that if q_n solves r_n , then any extension of q_n will continue to solve r_n, \dots, r_n . However, extensions of q_n will be required to solve r_{n+1} . A possible advantage of your prefix codes: If you use r_{n+1} as argument for q_n , the program could stop with useful output (or do something illegal) or print some thing clearly wrong before it asks for a new token. This will save time, since we then know that no extension of q_n will solve r_{n+1} .

In Prolog prefix codes that I use, the prefixes can never be extended, since Prolog and λ has not made make one. The advantage (if any) of functional languages is that it may be possible to devise a system that combines both extendable prefix codes and functional languages.

Does 3 input code do that?

It subclasses in non-functional languages are dangerous to have, then Prolog is no guarantee of meaningfulness in a sequence of tokens.

Looking at pp 20-21 of OOPS: After OOPS solves r_n it has this solution. But when it try to solve r_{n+1} , it remembers PC's of tokens for r_n but it doesn't have to derive a code to solve both problems!

In λ , I would have to put indices in the λ 's so TM could tell which problem was which!

In both of J's problems: $1^n 2^n$: Pro extension solves them w/ arbitrary used.

But didn't the PC's they inspired help obtain a final soln? In a TOFF, extensions were never found to be solns. In $1^n 2^n$, the final soln did not benefit from earlier extension "solns" as re it, took 20 times as long as it would if it had to be imposed by the in-line solns.

So it looks like "extension" did not help solve any of Phase 2 problems!

My solns of PC's as $(70)^n$ are not nice: J broad'd initial token PC's to Prolog across not all =. see P17 of OOPS: However, he doesn't say just what his initial bias was! He may have simply passed custom insns to Prolog. 3 insns for $1^n 2^n$ in Prolog, unbounded PCs & best insns for TOFF.

So it looks like: for Phase 2 problems, the extension property of his 5 different insns for TOFF

languages did not help solve the problems.

His choosing an initial set of insns to "bootstrap" amounts to selecting which insns will be (manually) used in the soln of the problem.

Assigning custom insns low numerical "names": This should be done more carefully: It may indicate a fault in the addressing system. The groups of 6 insns can be controlled either initially by trainer (which J does do by selective "booting" - but doing a different spread on tokens for each problem - doesn't sound so good!



FIRST DATA MERCHANT SERVICES
20 MALL ROAD
SUITE 350
BURLINGTON, MA 01803
www.firstdata.com

10

5 → 3
P → A
10 hrs.

SAY Learning recursion is usually done by writing that $f(n) = \text{some}$ func of $f(n-1)$ or $f(n-2)$ — Not by trying various recursive functions.

Also, addition, mult, expn... could be used as recursive functions.

These could be in unary notation ($n \equiv 1^n$). Then A could seem to xlt. unary to its internal operators on integers

Ex: J's assignment of (low nos. to certain insts) so they would be addressed / more likelihood!

He should arrange so that instructions don't have pc's of being put on stack that are + normal for pc's — or some variant of that if he likes — but putting it to 4. top array of how numbers are

2000
v.s.
 $6^4 = 1296$

constructed, seems unnacy & Bad!

pc's last ff! (346): No parts on stack they don't boost. — so he boosts q^2 .

This should not be done that way! q^1 But q^4 should have = likelihood of being boosted.

Unless, there is a special reason to reboost a particular q^2 .

Is "BOOST" somewhat like a "Recency" context?

On p22 J. says that by2, dec, boosty have been boosted! I don't see where or when!

Also the frozen code to recursive solution to $1^2 2^2$ problem was supposed to help

2 above "2" why "2"? There were 4 lines to the $1^2 2^2$ problem. — See 4.2 was

f. recursive soln.

See Page 17: But why dec is not inc? (it has both insts). i.e. $c3 \rightarrow c2$ via "dec".

if he used inc instead of dec, it would be difficult (unlikely) to get "c2" — which is what he needed

Superficially + ambiguous here. Reconstructing — better way would be to have both.

In general, the idea of starting with an "overadequate" set of insts, then have it

solve many simple problems to see which insts were most used. seems good.

But, for more difficult problems, recursion, say, is needed. So solve easy recursive problems to get needed insts by and pc's.

I had this idea (long ago) that TM would have to learn special techniques for

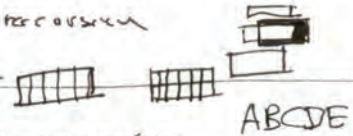
solving very diff (long search) problems.

That diff probs could be "essentially different" from easy problems.

Note: If we have an overadequate set of insts, then we solve a bunch of problems:

If we arrange the pc's diff. solutions so that the pc at 4. entire set of solns. is max, we will

get the "straight rule" to compute pc's! So Lap's rule is probably fairly good.



00

10

20

30

2.7.03

NIPS

Letter to Juergen.

20:42.19 : I noticed on your website a discussion of the demise of refereed journals. While I pretty much agree with your ideas on this, I think that an additional problem has to be addressed. As is, there ~~is~~^{is} a very large number of papers ~~is~~ accessible - some of interest to me - some not so interesting. Finding the potentially interesting ones might be done by ~~somebody~~ a generalization of Google, but finding papers that I actually want to read, is a much more difficult problem. ~~Refereed~~ Refereed journals cut down the search space a lot (admittedly, in a procrustean manner!).

0

At the Dartmouth workshop, we considered the idea of "editors" - ~~the editors~~. Anyone could set themselves up as an "editor", and list papers that they thought were good or bad, and why. A person could then ask a searcher for papers "recommended by A, B and C and not ~~recommended~~ by D" or various other search constraints.

10

I'm not sure of the mechanics of this, or whether it is workable. To some extent, when we write a paper with ~~our~~ citations, we become an "editor" of this sort - but often it is difficult to find useful references in a paper, to the papers cited. - i.e. I may know that ~~John Doe~~ paper x was cited in paper y, but finding the actual citations in the text of y may be difficult - ~~and which books, that I have read, that~~^{it could be} and very often a paper will not tell much about a paper it cites. ~~a person~~

30

Anyway, the general problem of finding papers that ~~is~~ really wants to read, remains an important problem - even with a lot of refereed journals around. We must find a way to deal with this problem ~~which~~ which becomes even more serious without ~~refereed~~ refereed journals (I!).

7 PM 290
snow 2-4" 320
NY

00

Some random Comments on OOPS:

- You proposed several possible ways to ~~solve~~ get a single solution to $r_1 \dots r_n$ given $q_1 \dots q_n$. You then picked 2 of them and spend equal time on each, saying that the resultant search was within a factor of 2 of optimum.

I'm sure the list of possibly good ways to use L-search to solve problems is quite long, but I don't know which is "the best" (if any). — ~~still~~ I'm still working on this problem. That either of the 2 you picked ^{should be} optimum — is very unclear!

01

When I first read about your use of prefix codes for programs, I thought it was a great idea! More thought convinced me that it was not at all ~~obvious~~ certain!

The apparent advantage is that if q_n solves $r_1 \dots r_n$ then any extension of q_n will continue to solve $r_1 \dots r_n$, thus saving lots of time in checking new trials.

The question is: how frequently are extensions of q_n able to solve r_{n+1} ?

02

I have no good ideas on this. In the 2 examples you gave, $1^n 2^n$ and Towers of Hanoi, you would have done better without considering extensions at all.

It may be that with a properly designed prefix language that it is, indeed, more likely that an extension of q_n will solve r_{n+1} .

Prefix ~~code~~ code languages seem to be useful in describing universal probabilistic functions. I have ~~found~~ ^{developed} a functional language with ~~prefix~~ prefix properties (inspired by an idea of Chaitin's) — **but I** don't yet know if it's really practical, Needs work!

03

There is some possibility that you may be able to double the speed of OOPS by changing a few symbols in the instructions! — This would change it from a limited doubling L-search to parallel L-search. — However I'm not yet sure I understand your program I may be completely wrong about this.

omit
looks
on library!
may need but
change!

~~Preparatory to solving of Hanoi, you boosted~~

At a certain point you boosted by 3 , dec and $boostq$. Why not "boostinc" also?

It would seem that inc and dec should have about equal probability. "dec" was specifically useful in creating $(C3, dec, boostq)$, which enabled the particular solution of TOH that it found.

Without it, this particular solution would not have occurred. If you had used "inc" instead of "dec", ~~wasn't~~ it may have found some other solution — though

00 it seems less likely. Probably it would be best to boost both "dec" and "inc" —
This would reduce the probability of solution somewhat, but not much.

I don't understand just which solutions to problems you "freeze". I ~~was~~ initially
thought you freeze all successful solutions. In this case, there would be 4 of them for
the $1^n 2^n$ problem. So why is the code for the 4th solution, "2"? —

It would seem that "4" (or perhaps 3, if one starts with ~~exp~~) would be appropriate.
~~the~~ {later say: perhaps was }

7	6	5	4	3	2	1	0	
1	2	3	4	5	6	1	2	3

 humming of them! with
1324 1234 1234
Not unnumbered

10 I am still trying to understand some details of OOPS, so that I may be able to
11 usefully criticize, and possibly improve — and possibly get good ideas on
how to do Alpha 1. → 48.00 new TP

~~Will write more on this later~~

I want to make a list of my comments (+ or -) on OOPS! Also points that I
don't really understand (like which parts get frozen, when),
My main immediate goal now, is to get alpha "off to end".

Def: α_1 is α , ~~can~~ can work 5 induction problems. ABCDefghij ABCDE abcdefg.

Perhaps a main Q about OOPS: Is the "Boost" heuristic any good? —
If good, could it be improved? What's the best way to do a boost? (if it's a good idea at all!).

As is, \checkmark . decides to freeze certain codes. Which codes to freeze are unclear.
Having frozen them, they are in ~~memory~~ frozen memory & accessible for & boosting
& /o for editing & use in programs. (perhaps "boost" is a kind of "editing"??).

Presumably, all code ~~chunks~~ chunks in ~~the~~ frozen memory have = pc (as in SAARD ANH-PCM)

30 Boosting to ~~single~~ single insts byz, dec, boostg — how is this done as part of it.

System! Sounds a bit A.H.! He may have just decided that ~~was~~ an initially
uniform applied over all tokens was unreasonable, & that this would be the way
to fix that. (Hvr. "Boost" is unextreme change of: p.d. on tokens ~~at start of~~
for an induct system! As system matures, $\cdot 73$ (hypothetical) occurrences of
a token will be less important — In fact, in general, just how to best do a "boost"
is unclear.

10 : ~~_____~~

01:47.11 The main immediate concern is the ~~boost~~ ^{boost} operation: Using ~~you~~ ^{you} were ~~able~~ ^{apparently able} to multiply the speed of solution of the Towers of Hanoi ^{problem} by a factor of 1000.

Questions:
~~Was~~ ^{Was} the use of boost legitimate? (i.e. will it continue to help solve difficult problems? What is the best way (if any) to use boost? If it is imperfect, how can we fix it ... etc.

I use something like boost in ~~Alpha~~ Phase 1 of Alpha, but in a very "muted form".

Should I modify it?

Will write more on this later

— Ray.

Perhaps ~~start with~~ introduce univ. dist. via section in Sol 86

00: Kol talk: Start w. section of ^{Amended} ~~Amended~~ Kol lecture ^{notes} about discovery of Universal D.F.

He was surprised to learn of my earlier discovery of the universal distribution and its application to inductive inference and prediction of all kinds. He told ~~me~~ his colleagues about this and so for a long time, my work was discovered was much better known in Russia & the Soviet Union than it was in the rest of the world.

~~He was surprised to learn of my earlier discovery of the universal distribution and its application to inductive inference and prediction of all kinds.~~

~~He was surprised to learn of my earlier discovery of the universal distribution and its application to inductive inference and prediction of all kinds.~~

At first, The Universal Distribution seemed to have two serious flaws.

First it was incomputable and second, it was dependent

on $\frac{1}{2}$ hr. for part to end of page.

Do an outline of talk, first.

The discovery

First explain how Universal D.F. is obtained from any Universal Hypothesis.

Then how Kol discovered it comp.

how he discovered found out about my discovery.

Univ. D.F. had 2 Appeal flaws:

1) Depends on just what unc. used.

1) It was "incomputable".

Discuss incomputability: i.e. test probabilities is incomputable.

Test any computable eval. of prob. ~~is~~ ^{will} have errors in it of unknown size, unknown size & unknowable size.

These are two > kinds of errors in prob. & statistics!

1) $\epsilon \leq \epsilon$ ~~error~~ \rightarrow to account data uncertainty ($\epsilon \leq \epsilon$ - we don't know rest of ϵ)

2) Model error \rightarrow Model uncertainty - we haven't tried out all "reasonable" models.

Study of Univ. D.F. of ALP ~~has~~ have revealed much more about it than that

PC is incomputable but all ϵ computable approx. have errors in form of unknown size,

which can be arbitrarily big.

[Give (0.1) as example of Model uncertainty.]

Discuss Dependence on Unc.

- 1) 2 kinds of Approx
 - 2) Philosophers (zero info to start)
 - 3) Statisticians (Entire Data is known).

Express w. examples.

Perhaps discuss Statistical Axiomatic Principles.

What to do if you have no info?

What to do if you have no food, water or air: ~~just~~ Do nothing - you Die.

That ALP tells how to use past experience in terms of info in hypothesis to construct reasonable approx.

So dependence on Unc. is necessary & desirable ... our usual way to create approx.

20

30

2 aspects to Goal: Priority of in class
(2) Heuristics: how to find by pc codes.

00:49.40

Second Part: on Alpha.

1) Discuss General Goals: "Very capable scientific assistant".

2) 3 parts α . phase 1, α phase 1 advanced, phase 2
d produ s produ Env, Oz, v.g. INP

Not clear what this is good way to develop project.

(3) Give equ. for α , i.e. (4.14.18)
Also "second version": and explain.

4) Tell about convergence param.

5) Discuss Heuristic prog. - a fast way to solve (3) ... find short (3 hy pc) paths.

6) Discuss Phase 2 ::::: :::::

Discuss Lsearch!
experimental cc
"cheat" by size of psm.

Computer Chess v.s.
Human Chess ^{has a near} 1-3 ^{moves/sec}
2-4 ^{moves/sec}
Deep Blue ^{3 M} million Bds/sec

Deep Blue may ^{20 M}
Bds/sec.

perhaps: Early!
Discuss Identity of
a) PC determination, Estimation
b) Learning.

00: 1909: **epistemology Anthropocentric Principle**: A certain initial principle + suitable TSCQ \rightarrow Apex of Modern Man.

W.o. proper state of TSCQ, we would not be here to ask Qs about principle!

Situation similar to that in physics: If laws of physics were not close to what we observe, we would not be here to observe them. Hvr. in physics we have to determine what basic laws & constants are.

In a way in studying Learning, we have to determine/what man's principle is: perhaps

guess at some initial principle: T. TSCQ that led him to his present principle ("state").

for talk: Explain Anthro-principle in physics: Gen Anthro-principle in Logic.

10
11 TSCQ's in OOPS: prob1 = 1921; prob2 = TOTT:

It looks like t. early examples in both probs were unnecessary (minimal)

13 { In both cases, t. early problems were used only as "tailors" to help find TM, \Rightarrow ~~TM~~
14 { prob_{1,2} (k) had a recursive solution for TM (k).

But by way: Just what was t. "21" problem? Wasn't that TM find a path w. n as input is 1"2" as output? The TOTT was, I think, given initial set of piles of disks, to legally transfer them to new piles.

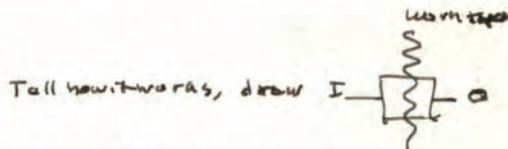
20 It looks like .13 - .14 is true. we would have saved some time by giving prob1 (3) than prob2 (3) (I think \Rightarrow was solved recursively in both cases, but check) prob1 (4) wouldn't have taken much longer than prob2 (3); but prob2 (4) could take 2 times as long as prob2 (3) so t. TSCQ prob1 (4) prob2 (4) would be much slower.

24 In human logic of recursive functions (I wrote about this recently) t. mode of logic is, say prob1, prob2, prob3 ... T. logic of t. recursive function occurs
26 when TM notices relation of soln. of prob₂ (k+1) to soln. of prob₂ (k).

In OOPS TSCQ, t. only ~~thing~~ thing that OOPS learned from prob1 (6) that enabled it to solve prob2 (3) was that t. known initial states used in recursion were useful. — Just which initial states? (Notice \Rightarrow CFG (k) \equiv 1921 (k) \equiv prob1 (k) or TOTT (k) = prob2 (k).) Using TM have better factor of 1000 improvement \leftarrow but looking this more carefully using t. work of .24 - .26 t. factor of improvement should be much larger. (much smaller CJS).

30 While OOPS got "factor of 1000", CJS was still very large (I would say "too large")
N.B. t. initial states for CFG (k) was correct
• 3rd " " TOTT " "

00 : **Vol. talk: Part 1: (see 49.00-40)**
 Derb's discovery of ALP, Univ. D.F.



unresolvable output type -
 => **main constraint.**

What is PC that machine will write "x" before stopping (it reads a character code).

U: 20 machine may be easier easier to derb.

G No

$$P(x) = \sum 2^{-2^i}$$

$$P(x) \approx 2^{-k(x)}$$



Tells how many confs
 use by many of you
 have think for!

Discuss k complexity in this approx. Tell about k 's interaction Randomness

is in properties of complexity... was surprised to learn of ~~any~~ earlier work

on induction Accuracy: Discuss Conv. Problem.

Also: w. ALP: Overfitting as on index fitting
 never occurs: it fits just riter
 No need to divide data into fitting set and testing set.
 just 1 set is needed

Diffs: 1 Incomputability: ~~is~~ have π or $\sqrt{2}$ not computable exactly but

approximable to any precision. ~~ALP~~ ALP also approximates to any precision.

But one non-learn useful assumes not cover. It is provably impossible to know what

This might be regarded as a Bug in ALP! Actually its not a Bug but a feature, that means it better than other proxy evaln methods:

Let me explain: Ordinary ^{computable} methods of proxy evaln. ~~show errors~~ ~~inherent~~ ~~under~~

you compute it out your finished. If there are better models of the data than you have used, your estimates will be in error. In all cases, it is impossible to know if there are better models available.

Perhaps discuss ^{uncertainty} ~~error~~ v.g. model ^{uncertainty} ~~error~~

$$\frac{1}{2} (p^2 + (1-p)^2) = p^2 = (2p^2 + 1 - 2p) \pm p^2 = p^2 - p + \frac{1}{2} + p^2 = \dots = \frac{1}{2} (1-p)$$

Means sq. - sq. of mean means $(p+1-p)/2 = \frac{1}{2}$

0 < 1; $p \neq 1/2$ so $\frac{pn}{n} = \text{small norm of sq. } p = \text{symmetric. ; means } p \text{ so } p - p^2 = p(1-p)$

so $\text{SSZ error} = \frac{1}{2} \cdot \frac{1}{\sqrt{n}}$ stands down.

Compare w. Model uncertainty error $\approx \frac{1}{2}$.

$$\sqrt{\frac{p(1-p)}{n}} = \text{approx error} = \frac{1}{2\sqrt{n}} \text{ for } p = \frac{1}{2}$$

ALP is honest about uncertainty: other computable methods are not

ALP suggest how to \downarrow error: suggests non-better models.

other ~~in~~ computable methods usually don't.

Summary & Conclusions

Next Difficulty: ALP depends on choice of unc.

Co variance from $\text{error} \leq \text{area} \leq \frac{1}{2} \ln 2 \cdot b$ $(\text{error})^2 \leq \frac{1}{4n} \ln 2 \cdot b$ v.g. $\frac{1}{2\sqrt{n}}$ for usual statistics

Here "b" is very much UNC dependent: so ALP is very "sensitive".

200
1998
2000

Meaning of A Priori

1) Classical philosophy: case of No a priori ^{Not a priori - what to do? can occur}
No a priori - what to do. - very unlikely to occur.

2) Biostat A Priori: Info before data is seen. - change of ~~the~~ a priori during (6) as new updates
after each problem is solved, ~~the~~ data predicted.

What we have w. is largely a priori, but not well known.

What we have as Adults is to solve extrem known: we can use language of our science to construct instruction sets for universal computers. That useful vocabulary of science has very
into in it - - - to be used for a priori.

In ~~the~~ Physics we have Anthropic Principle. ^{52.00} This laws i constants of physics happened to be such that Intelligent life can evolve, and thrive. This seems to indicate only a narrow set of values for those laws i constants are possible. - a random initial state would not occur.

In a similar way, initial life in universe has to have certain a priori built into it a "reasonable" TSg to bring it to present day, in the big bang.

In Physics, the laws i constants have to be discovered i determined.

In ~~the~~ statistic in ling occure, the initial a priori have to be guessed at in the beginning.
[Actually, this evolution to Man may be part of Physics "Anthropic Principle".]

In the case of ~~the~~ laws of physics, small deviations from present laws would do, by now -
But I suspect that modulus of a priori ^{initial} TSg could be rather brandy still intelligent
human would evolve.

Forst. is Bayesian View: Non Bayesians do consider a priori info, but they doubt
position form of a PD of S. Forst I'm under stand.

One common way to put in a priori info is to limit the model set considered - give
all other models a priori = 0. There will always be models that have not been
considered but might give much better prob values. We are not entirely certain that
we are not entirely certain that there are no such models.

If no models seem to work so far factually - the statistician should be happy
that will always be in the background living
the possibility of discovery that he may be in a god's paradise!

Using the universal distribution, we are not surprised when radically new
models are discovered!

60:5240: **OOPS** (cont.): Ad. front. part: Say OOPS calculates pc of each token it takes: One interesting point is that OOPS instructions can change state of system much (much more). Boost does this. Simply presenting a token changes system a little (it's not/don't of various token PC's) — but this is a small change compared to ~~Boost~~ "Boost."

In what sense ~~is~~ does it above give a Universal diff? If it does, does it have any advantage over a lang (like AZ) that doesn't have such large jumps in PC's?

[Pro, if I properly implement OSL, AZ might have "sudden" large PC changes for tokens]

Also if AZ discards a v.f. conc. that it had a known buffer, then PC's could change much.

— If AZ really allows many big discrepancies (since ~~the~~ PC's of all concs must \neq to be produced (or run OSL)) — then \rightarrow less likely in AZ.

SN on "Negative PC's" As TSO continues, TM can discover new concs, which could be a way of implementing PC's "correction"?

"correct/modify" PC's calculator in front. Could "negative PC's" be a way of implementing PC's "correction"?

A better way to interpret "boost": Boost looks at stack, if K is large enough for Boost, N mod K .

But in general, Boost should give about equal priority to all frozen ppm's, (unless we want to be able to style out certain ppm's for re-priority)

T: first time Boost is used, it really narrows down: insts that are at all likely!

Since this is usually a small set of insts, we can quickly do all of them. If one of them is a Boost, we can augment it: instructions. In solving TOH,

for $1^2 \cdot 4$ we start w. board of $c1, c2, dac, boost, byz$. 5 insts

for TOH " " " " " $c3, c4, c5, dac, boost, byz$. 6 insts

Initially, the solns are usually limited to 1 boosted insts: w. 1 or 2 others used. 0 PC's insts have

PC of max 3 "boosted" insts. $\frac{1}{7}$ boosted $\frac{1}{7.73}$ for unboosted $\approx 25 \approx 7 \times 7$

Essentially what OOPS did in solving TOH: It found a $1^2 \cdot 2^n$ ppm that had insts needed for TOH to solve.

For a recursive soln, def, call and end are impl. insts: That they should be useful for both $1^2 \cdot 2^n$ & TOH is not surprising.

It took longer time to solve (recursively) TOH w. $1^2 \cdot 2^n$ (TOH w. $1^2 \cdot 2^n$ & PC's via "Boost").

SO GENL Conclusions: T. OOPS mechanism enables use of tokens of "successful" ppm's, but in a rather exclusive way: So if a ppm is picked pretty much only one token for that ppm will be given reasonable PC's: Other tokens could be in PC general used or not used at all!

This is a rather sharp opposed to "recursion" way to select out tokens:

A "softer" way would not \uparrow work as equiv to $Q_N (=73)$ case count, but something better — perhaps the denominator could be \uparrow by a certain β , if numerator A by some (absolute) \uparrow as denom. Also note that in more mature TM the $\Delta = Q_N$ will be less imp. (But amount to Q_N ratio)

There is some attraction in the Boost idea: That a token should consist (nearly) of 1. tokens in a previously successful ppm.

5 3
6 2
1 1
2 0
3 ...

$$\frac{9.3}{4.5 \cdot 10^{-3}} = 2k$$

00: 5pac
01 Another view: that Boost keeps together tokens that are "synergistic" - that work well together (like defnp, callp, endp as in TOH in 1st 2nd recursive solns).

04 I noted that if one wants to get tokens pc's so as to (post hoc) maximize pc of all solns then, pc's will be for straight rule (or LoP/cyclic paths) - which is what AZ is doing (w.a boost) gets. → How can we insert the synergy effect of .00701? ← .09

07 I suspect that the success of the 3 (.01) insts in a single recursive path is not enough evidence to justify using them in a "Boost" only one case. → 57.20

09 One may to perhaps get more wt. for the recursive 1st 2nd soln; Use this pfm on many values of N as ~~TSQ~~ TSQ

09 .04 In .04 if the "boosted" set were more probabilistically less wt., it would seem like a more reasonable approach: As is, one is much more constrained to the "boosted" set(s) - i.e. better "boost" of 55.34 seems more reasonable.

On the initial boost of by2, dec, boostp. + C1, C2 or C3, C4, C5.

The initial 5 or 6 insts can do better by themselves, so they must boost some previous qⁱ to obtain a useful inst. set. So we have this qⁱ boost + by2, dec & C1, C2 or C3, C4, C5.

But means we can generate a forest to select various best qⁱ's boosted. - Also have those integers (i by 2) available w. the boosted set of insts. We also have 6 insts, but

They have a pc of (maybe) 3 boosted insts.

14 Q: Can the system afford to boost > 1 qⁱ set of tokens? It amounts to
+ pc cost amounts to boost; is some with inst. My impression is that in addressing of qⁱ to "boost" is okay: that all should have about same pc's. If illegal addresses for i (in boost qⁱ) are quickly recognized as illegal & ~~ignored~~ ignored, then they cause no trouble. - would like to be sure that all legal i have about same pc's here. Probly to lang. could be modified to do this: (see 55.12-15) One way, after insts boost is put a stack, any integer up to base i of domain qⁱ is legal token. (P13 is in line w. OOPS computing of pc of next token) - We may want those integers to be tokens in their own (denominator) domain for pc.

10 This would resolve 19 a bit: > 1 boost could be good. 2 or maybe 3.
- Pro probly 2 ^{would be} and "diversity" in instructions - Please see something like "crossovers" in GA.
A 3 way crossover may be unvary.

On "weakening" boost: As is boost amounts to few additional case counts for boosted tokens. One kind of weaker boost would be a few case counts w. a < 1.

looking at page 22 of OOPS: Trade of pc of soln. to TOH with boost vs w/o boost!

$$\frac{(1+73)^3}{7 \cdot 73} \approx \frac{(1+73)^7}{12 \cdot 73} \approx \left(\frac{1}{12}\right)^7 = \frac{(7.73)^4}{12^7} = \frac{2^4}{12^7} \cdot 73^4 = \frac{73^4}{15 \cdot 16} = 2.8 \cdot 10^4$$
$$\frac{(1+73)^3}{7 \cdot 73} \approx \left(\frac{1}{7 \cdot 73}\right)^4 = \left(\frac{1}{7 \cdot 73}\right)^4 = \frac{1}{(7 \cdot 73)^4} = 1.9 \cdot 10^{-3}$$
$$\left(\frac{1}{7 \cdot 73}\right)^4 = 1.9 \cdot 10^{-3}$$
$$\left(\frac{1}{7 \cdot 73}\right)^4 = 1.9 \cdot 10^{-3}$$
$$\left(\frac{1}{7 \cdot 73}\right)^4 = 1.9 \cdot 10^{-3}$$

2.8M inst
15N
73⁴
1516
2.8M
2.8M
= 2.8M
2.8M
= 2.8M

00:56:40 So ratio = $\frac{73^4}{15K}$ for bracketed. i.e. pc + some. $(5K)^4 = 11$ i.e. $73 \rightarrow 11$.

$$\left(\frac{1+73}{73}\right)^7 = \left(1+\frac{1}{73}\right)^7 \approx \left(1+\frac{1}{73}\right)^{73} \cdot \frac{7}{73}$$

for $73 \rightarrow 11$ $() = 2^7 = 128$

if $73 \rightarrow < 11$ we have much loss: if $73 \rightarrow 1$ we have loss

factor of $\frac{128000}{128} \approx 1000$ factor of loss! I don't understand why we lose / lose any thing!

OH! we lose because of the initial boost of

WOPPS! for boost < 73 my figures are wrong! I have wrong domain factor.

$\left(\frac{1+73}{73+11 \cdot 73}\right)^7$ vs $\left(\frac{1+x}{73+(1+x)}\right)^7$ v.s. $\left(\frac{1+x}{73+6 \cdot x}\right)^5 \left(\frac{1}{73+6 \cdot x}\right)^4$ for $x=0$ to ratio is 1.

in **F: /SM/ OOPS!**

for $x=0$ we got ratio = $2.57 \times 10^{-6} = \left(\frac{1+x}{73+11 \cdot x}\right)^7 / \left(\frac{1}{73+6 \cdot x}\right)^4$

$$2 \cdot \left(\frac{1}{73}\right)^3 = 2.57 \times 10^{-6}$$

T. reason $x=0$ has ratio 73^{-3} is it involves the pc + 3 dec boost, which $pc = \left(\frac{1}{73}\right)^3$ total loss since 73^{-3} gives nothing at $x=0$.

We notice cost of (c3, def Boost) \rightarrow 7 symbol soln at TOH

7 " " " " above.

so ~~ratio~~, c3 dec boost has a "drop" value & if $x > 0$ it has pc of next 7 symbols.

If $x=0$, we only have to drop cost & no time pc of ≥ 7 symbol soln.

20 : 56:07: On "synergy": A very basic idea in TM is that good new ideas are formed by

combining old good ideas & parts of them. If we run after ≥ 2 boosts (50055:22)

We can regard them as tentative recombinations of old "Good ideas" & "Boost", but

is very vague about how to do the recombination, and it doesn't preserve order or spacing of

taken in the "boosted" pm. For 1st 2nd his soln was $\boxed{\text{defnp}, c1, \text{calltp}, c2, \text{endnp}}$

for TOH soln was $\boxed{\text{defnp} | c4 | \text{calltp} | c3 | c5 | \text{calltp} | \text{endnp}}$

The TOH soln had same order as ("2" job), but it also had an extra "calltp".

If DOPS did a lot of recomb steps, it might realize that defnp was a good start, and np was a good end of a recursive pm. But one or two "calls" were usually used.

I really don't understand either pm, so I can't really tell how one takes advantage of tags in the other.

If we assume 4. end symbols defnp, endnp, we only need $\left\{ c4, \text{calltp}, c3, c5, \text{calltp} \right\}$

If we cross formulae to have formulae $c1, c2, c3, c4, c5, \text{calltp}$. (6 symbols)

its pc is $\left(\frac{1}{6}\right)^6 \approx (46656)^{-1} = 2.14 \times 10^{-5}$ rather large!

SN OOPS

Doesn't really learn much from previous prob. solns: only smallish biases in Tokaupcs (via Laps rule) plus gross pc changes in disturbing via Boost. Humans certainly learn transfer more into from previously solved problems. My own approach to try to put any hour I can think of into a form that TM can acquire, use. - Don't let a fort tell you that there ... possibly w. "Hints".

code?

Kop. talk

20: 54.40:

This Lecture is going to be about three things:

- 1) The Universal distribution and Algorithmic Probability.
- 2) Its relationship to Kolmogorov Complexity
- 3) The application of these ideas to a system for Machine Learning (which has always been my main interest.)

One way to introduce the Universal distribution is by considering the first "No Free Lunch" theorem. A very natural way to do prediction of sequences of symbols, is to assume an a priori probability distribution on all possible sequences.

The obvious first candidate is the Uniform distribution: All sequences equally likely. Turns out that this is ~~really~~ very bad. And the solution is independent of the part: No matter what happened in the past, the next symbol ~~has~~ ^{has some} probability as any other. **Back to the drawing board!**

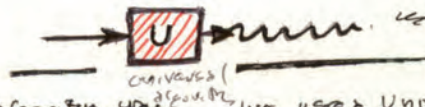
If all strings can't have some probability, what about descriptions of strings all having some probability? This works out quite well.

A better way: Suppose we have a long sequence of symbols: for example, binary symbols — and we want to extrapolate it. One way to do this is to hypothesize that the sequence was created by some algorithm — then use that algorithm to generate the sequence. There is a nice way to do this.

We start out with a universal algorithm, that maps input strings into output strings. "Universal" means that the algorithm can simulate any other algorithm using a program input to a certain length longer than the algorithm being simulated.

omit

(draw pictures)



Since we don't know just what the algorithm was, we use a Universal algorithm one that can simulate any other algorithm. Since we don't know just what the input to the algorithm was, we assume all inputs to be equal (likelyhood). This amounts to saying that the sequence we want to extrapolate was created as the output of a universal device with random input.

Aside: There are many possible forms of universal algorithms! One of the best known is the Universal Turing machine. And most any computer is a good approximation.

This is the Universal distribution — the output of a universal device with random input.

The ~~distributed~~ probability distribution on output strings can be used to extrapolate any finite string of symbols.

Superficially, this looks like the long-sought "philosopher's stone" — the device that will answer all questions.

In more concrete terms, if x is a finite string and we want to assign probability to it.

Find a set of strings $\{s_i\}$ such that $f(s_i) = x \dots$. Then $p = \sum 2^{-l(s_i)}$



personnel had an application to Sci Method: Choice betw. Computing Resources that both adequately explain data.

It can be also modified to extrapolate an ordered set of finite strings — (the legal sentences in some unknown language).

If we use the technique for prediction, the errors in probability obtained decrease rapidly with the amount of data used.

We can use this model to determine just how many coefficients to use in linear or non-linear regression. — Thus solving ~~one of the most serious~~ a very serious problem in modern statistics.

When we use the method for prediction, there is never any problem of "over fitting" or "under fitting": The fit is always just right.

When we do prediction with data we don't have to divide the data into a "training set" and a "test set". We use all of the data for training and the expected errors are always just right.

It would seem from all of this that the universal distribution is a kind of "philosopher's stone" ~~that solves~~ that solves all scientific problems.

There do, however, appear to be two Bugs with First Post-Proc ~~universal~~ probability in the universal d.f. ~~are~~ incompatible

Second, the P_{ij} values obtained, depend on just which Universal machine was used to generate it. It is my thesis that neither of these are Bugs — they are ~~rather~~ very desirable features.

Let me explain: The real difficulty is that probability itself is incomputable. Any ~~theory~~ ~~claiming~~ ~~to~~ ~~be~~ ~~a~~ ~~model~~ ~~of~~ ~~true~~ probability must be also incomputable.

A bit more detail: The value of π or $\sqrt{2}$ are "incomputable" in the sense that we can never find exact decimal values for them. ~~The universal distrib~~ but we can make approximations that we know will eventually have arbitrarily small error.

In the case of probability, we can also make approximations that approach zero error. In the case of π , however, we can always know a limit on how large the approximation error is. In the case of probability, we can never know a useful limit on the error size.

In ^{probability} statistical Estimation, there are ^{two} ~~two~~ kinds of uncertainty:

- 1) sample size uncertainty
- 2) Model uncertainty.

To illustrate sample size uncertainty: Suppose we are given a binary sequence and we want the probability that the next bit will be 1. We could just take the relative frequency of 1 in the sequence and use that as a probability estimate. However, the expected error in such estimates is proportional to $n^{-1/2}$; n being the sequence length. — ^{When} we have a longer number of symbols, $n^{-1/2}$ becomes very small.

Things that I will not discuss in this lecture:

~~We can approximate the universal Turing machine~~ ^{is} ~~algorithm~~ can take many forms: One form is an ordinary computer language. In this case we can express our a priori information to some extent, by making ~~the~~ ^{the} instructions set so that it expresses very compactly, regularities that we have observed in the ~~data~~ ^{data} - that we feel are important. The language APL is able to express ~~our~~ ^{our} commonly used mathematical concepts in extremely compact form. The dependence of a program on the unit. algorithm is something ~~we~~ ^{we} can't express to express our personal grasp by choosing the characteristics of ~~the~~ ^{the} algorithm.

So ~~to~~ ^{to} summarize, incompleteness and dependence on resources are not really bugs in the system. They are necessary parts of ~~the~~ ^{the} probability concept that we want to learn to ~~use~~ ^{use} to take advantage of.

~~interact~~ I have tried anyway to use the ~~design~~ ^{design} of unit. algorithm to express needed sp. info. We can take advantage of incompleteness by using ~~the~~ ^{the} formalism to suggest new models/induction models. Incompleteness implies that there are always models that one has not yet tried, and it gives a somewhat different suggestion on ~~where~~ ^{where} to find new models: e.g. by considering ~~more~~ ^{more} ~~simple~~ ^{simple} models that are rather simple, but may be very time consuming to evaluate.

SN I think that I should outline just what I want to cover. As I write (~~on~~ ^{on} the "so to speak") I tend to wander off into branches that I can't complete. Just list topics, indicating just how far I want to go into it. Then try to find ~~the~~ ^{the} ~~best~~ ^{best} ~~way~~ ^{way} to do each section.

- Kol. intro
- S&S
- Barc
- NIPS talk
- NIPS Abstract

SW Discussion of ~~the~~ ^{the} system: Inductive int. ~~as~~ ^{as} "well defined problem" Use of ~~the~~ ^{the} classical ~~heuristic~~ ^{search} to solve it - using ~~the~~ ^{the} ~~ATSP~~.

Discuss Lenat's AM: Artificial Mathematics: Goal was to find interesting concepts, deductions ~~from~~ ^{from} forms. It had set of rules for how "interesting" something was, also rules to ~~generate~~ ^{generate}. After rules to find ~~new~~ ^{new} "interesting things" found

It is a bit out of goals ~~more~~ ^{sharply} defined: we want to ~~make~~ ^{make} ~~it~~ ^{it} ~~work~~ ^{work} ~~as~~ ^{as} ~~well~~ ^{well} ~~as~~ ^{as} possible. ~~Lenat~~ ^{Lenat} had very rudimentary understanding of induction - we have a much better model & a much better search procedure.

It has to solve ~~specific~~ ^{specific} problems that were given to it. At first it solves them ~~using~~ ^{using} "standard" search procedure, using a GCD to guide its search. After it solves a problem, it tries to improve GCD. It does this using heuristic tricks that have been either programmed into it or that it has learned from its experience in problem solving.

Perhaps shorten it by ~~using~~ ^{using} ~~the~~ ^{the} ~~idea~~ ^{idea} of Lenat, but ~~do~~ ^{do} ~~have~~ ^{have} ~~used~~ ^{used} to ~~heuristic~~ ^{heuristic} ~~things~~ ^{things}. For the most part, heuristics are implemented as ~~modifications~~ ^{modifications} of GCD.

So outline: Goetz: ~~approx~~
Dorb. Univ. Df. : Perm. Qiro has 2 apparent dfts.

1) Incompleteness \Rightarrow subjectivity.

show ~~uncomputability of fact~~ 1) is property of property itself

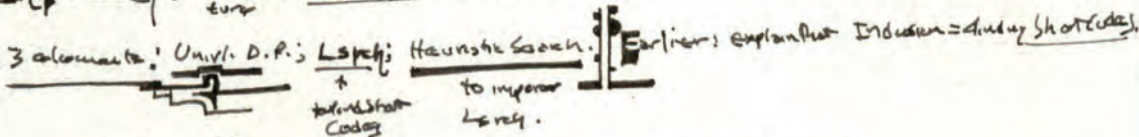
2) is natural to Bayesian statistics: Non Bayesian have different ways to include a priori info

but a priori info is essential. Attempts to make non-subjectivity of necessity result in poorer induction

Much more limited options

I will discuss a particular application ~~of~~ ^{in which} ~~of~~ ^{classical} ~~of~~ ^{Universal D.F. approaches} Induction - Machine Learning.

Introduce as ^{classical} ~~of~~ ^{Universal D.F. approaches} Induction.



Q: Is Heuristic Search a good approach to explanation?

Application of Univ. D.f. to Machine Learning:

1) Dorb. QA problem: say a great variety of induction problems can be put into this form. Q can be string of nos of instances Σ smart for A.

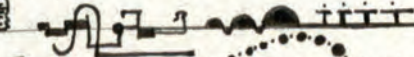
Formalization $\Sigma \Pi O'(A/Q) \dots$

$$\Sigma \Pi O'(A/Q_{uni})$$

$$\Sigma \Pi O'(A_i/Q_i) \text{ -- wt of } O'$$

To find good O'

We use search in simpler way guided by ~~univ. d.f.~~ ^{univ. d.f.} O' functions. We're able to make finite trials in appropriate probabilities.



One way: Start w. trial lecture notes

Discuss Kol's not using network use for induction. I used Levin's Koloid pattern induction wasn't math problem. Another pass: at first time there was no v.g. reason to believe that we use for induction had any accuracy at all - it had some accurate bits to support fact for long strings of bits, it ~~was~~ ^{could} converge of vito value.

In 1968 I ~~found~~ ^{found} that it would converge very rapidly to vito value. Simply stated: if there was any regularity in a data set, the universal distribution would ~~converge~~ ^{eventually find} it using a relatively small ^{amount} ~~sample~~ ^{sample} of that data.

"Eventually" refers to computer time, not ~~length of data~~ ^{amount of data used}.

I call this property of our universal distribution "the computational" property.

From here go into discussion of ~~the~~ ^{my} NIPS talk: Perhaps discuss extension of conv. theorem.

I'm going to ~~start~~ ^{start} out with a discussion of ~~the~~ ^{the} universal probability distribution and some of its properties and ^{how} it is related to Kolmogorov complexity.

The ~~idea~~ ^{birth} of this idea ~~was~~ ^{was} most ~~probably~~ ^{was} ~~connected~~ ^{connected} with

two events: First, my ~~publishing~~ ^{publishing} a long report ~~on~~ ^{on} Machine Learning.

I ~~first~~ ^{first} discovered the universal distribution in 1960 - as a direct result of some ~~earlier~~ ^{earlier} work I'd done on Machine Learning and Ginzburg's early - 64.00

1) General UTM!
random input:
At time T:
What is probab. that
output will have
property X?
What is probab. that
a stop output will have
property X.?

KolTalk: Start: Many years ago — in 1960, I discovered ^{what we now call} ~~the~~ The Universal Probability distribution. At the time, it seemed to solve all of the problems of Bayesian Statistics. It gave what looked like a Universal a priori probability distribution for all possible strings of symbols.

- At first, there seemed to be three main difficulties:
- 1) I wasn't really certain that it gave the correct probability values, or that it would work with a reasonably sized data set.
 - 2) The values of probability it gave were incomputable.
 - 3) The distribution itself depended on choice of a certain Universal reference machine ~~for reference machine function.~~

~~As for a reference machine, I did have a theorem that gave an upper bound on error and depended on reference function, but the bound was far too large to be of practical utility.~~

Five years later in 1965, Kolmogorov's idea covered ~~what we now call~~ Kolmogorov complexity of strings of symbols X is the length of the shortest program for ~~reference computer~~ a reference computer that produces X as output program as output.

Kolmogorov complexity of X is length of s : $M_R(s) = X$ such that ~~reference computer~~ $M_R(s) = X$ ^{is} length of s : $M_R(s) = X$ ^{is} length of s : $M_R(s) = X$

~~Kolmogorov complexity was~~ Kolmogorov was ~~interested~~ was interested in this complexity as a way to define randomness and as an interesting mathematical concept in ~~the~~ ^{its} use of ~~the~~ ideas of bits sort for induction ~~of these ideas to inductive inference.~~ How was surprised to learn of my ~~work~~ ^{work} ~~in~~ ⁱⁿ a ~~paper~~ ^{paper} ~~of~~ ^{of} ~~these~~ ^{these} ideas to inductive inference.

to me that Kolmogorov, who had ~~invented~~ ^{invented} the basis of modern created the (axiomatic) basis for modern probability theory — who had contributed so much to physics (to Higgs — ~~who~~ would not notice that his complexity could be used to approximate a priori probability. is.

$$P(x) \approx 2^{-K(x)}$$

Note that these probabilities are always integral powers of 2, so Bayes is usually in error by a factor of about $\sqrt{2}$ — an extremely large error. This may have steered him away from considering it as a candidate for approximate probability.

He did, however, publicize his earlier work, and so for many years ~~it was~~ ^{it was} ~~known~~ ^{known} in the Soviet Union ~~then~~ ^{then} in the ~~rest~~ ^{rest} of the world.

2.15.03
Nils



10:62.40 worth on formal Grammars: - both were invented about 1956

My main interest at that time and subsequently, has been in understanding induction in reference well enough to get a machine to do it very well... Essentially Machine Learning at a very high level - "Strong A.I."

5 years later, in 1965, Kolmogorov Complexity was discovered ^{what was called Kolmogorov Complexity} ~~and~~ thought it was very closely related to the universal distribution. Kolmogorov was a great thinker, and he read my earlier work. His interest, at first, was ~~primarily~~ in this kind of complexity ~~was~~ at first was in defining randomness and understanding complexity itself.

In complexity as ~~was~~ an extremely interesting mathematical concept. He was surprised to learn of my earlier work on application of these ideas to inductive inference. He did, however, publicize my earlier work and so for many years it was much better known in the Soviet Union than in the United States.

I'm going to start out by ^{talking about} ~~discussing~~ the Universal probability distribution, and how it's related to Kolmogorov Complexity. I will discuss some important

I'm going to start out by discussing the Universal probability distribution, and some of its properties - and how it's related to Kolmogorov complexity. Then I will talk about some work I've been doing in machine learning that employs these ideas. Then I will talk about describe the application of these ideas to Machine Learning.

Many years ago - in 1960, - I discovered what we now call the Universal probability distribution. It is the probability distribution on output strings of a universal computer with random input. ~~At that time, it seemed to solve all of the theoretical problems of Bayesian Statistics.~~ It gave a probability distribution over all possible strings of symbols. It seemed to solve ~~at that time~~ ~~some~~ very serious problems in the foundations of Bayesian Statistics.

Five years later, in 1965, Kolmogorov ^{independently} discovered "Kolmogorov Complexity". ~~The~~ The Kolmogorov complexity of ~~a~~ string of symbols is the length of the shortest program for a reference computer that produces x as output. Kolmogorov complexity is closely related to the universal distribution. If ~~k~~ is the Kolmogorov complexity of x then 2^{-k} is an approximation to the probability ~~universal probability~~ obtained by the Universal distribution.

Initially, Kolmogorov was interested in this complexity as a way to define randomness and as an interesting mathematical concept. He was surprised to learn of my earlier work on inductive inference. He did, however, publicize my work and so ~~it~~ ~~because~~ for many years it was much better known in the Soviet Union than in the United States.

~~For several years after my discovery, I wasn't really~~
certain that the ~~idea~~ ^{distribution} would give good probability values. I had lots of heuristic
agreements, but nothing very certain.

This all changed in 1988, when I ~~discovered~~ ^{discovered} the convergence theorem.

Nip take page 1 P 2.

The time needed to solve a problem will be, T_1 / P_1 . Here P_1 is the ~~probability~~
probability that GCPD assigns to the successful PST; and T_2 is the time
taken for PST to solve the problem. For hard problems, P_1 is ~~too~~ ^{too} small and T_2 is too large
solution time is ~~excessive~~ large. — ~~to solve hard problems by transforming them into~~
easy problems. This is done by "Adaptive L search." The idea is that we use solutions

of easy problems to modify the GCPD.
L search cannot solve hard problems directly. — But it can solve hard them
indirectly by first making them small.

Problem for hard problems the CJS will be impossibly large.
 T_2 will be too large and/or P_1 will be too small. L search cannot solve hard problems
directly.

~~There are two ways to~~ ... "It's too hard" ...
In the early education of the

... Solvable by L search
In the early education of the system, updating is done by looking at successful solutions

~~updating the GCPD so that previous solutions~~

In the early phase of education of the system (which I will call "Phase 1")
the system looks for regularities in ~~previous~~ ^{past} solutions and ~~uses~~ ^{uses} ~~them~~ ^{them} ~~to~~ ^{to} ~~guide~~ ^{guide} ~~future~~ ^{future}
searches in view of these regularities. It tries to ~~find~~ ^{copy} "what worked ~~in~~ ⁱⁿ the past". It has no concept of ~~optimal~~ ^{improvement} "improvement" or "optimization".

In a more advanced state of In Phase 2 — ~~update~~ ^{update} ~~more~~ ^{more} ~~advanced~~ ^{advanced}
stage of education, the system ~~tries~~ ^{tries} to get the best GCPD possible in
the available time. — it really "understands" the goal of optimization.

Phase 2 updating is a problem for operator production.

There is a tendency to avoid ~~this~~ ^{issue} ~~issue~~ ^{issue} by pointing to the incompatibility
of the Universal distribution. ~~but this is not~~ ^{the incompatibility of the Universal}
distribution is not the problem — the problem is that empirical probability
itself is incompatible — it always has some sort of uncertainty.

2.18.03
Naps

PP 67-100
don't exist.

00

The way it works: Suppose M is a universal computer and S is a binary string. Then $M(S)$ represents the output of the machine after it stops. If $M(S) = X$, X being a binary string, then S is a program for X .

If we give as input to M , a binary string that starts out with S — the machine would read S , output X and stop. — So if S was only 10 bits long, only the first 10 bits of S that particular input are relevant.

10

If we give a random input to M , the probability that the response will begin with S , is just 2^{-10} . If there are other programs for X , and they have lengths

L_1, L_2, \dots then the probability that M will output X under any given

random input is just $\sum 2^{-L_i}$ — the sum of all of the ways X could be produced.

If L_1 is the shortest program for X then 2^{-L_1} is an approximation to the universal probability of X , with respect to M .

IN 1965, 5 years later, Kolmogorov discovered independently by Kolmogorov complexity — it was the length of a string — it was the shortest length of the shortest program that could produce that string.

20

through "Infinite Kolmogorov" . . .

The way it works, say X is a binary string and M is a universal computer, and S is a binary string, such that $M(S) = X$ — i.e. S is a program for M that produces X as output. If L_S is the number of bits in S , then the probability that S will occur as input is just 2^{-L_S} .

problem

It seemed to solve all kinds of problems in prediction and seemed to be the deep with some difficulties in Bayesian statistics.

2/20/03

colth

not $\frac{1}{2} L/hr$ $\frac{2}{3} L/hr$ $\frac{1}{3} L/hr \rightarrow 800/hr$

20:
201
202

TC of production q 's \rightarrow would be done for PST in work $\frac{P}{x}$ ^{assuming irreversible}

But I had a bunch of PSTs w. some $\frac{z}{u}$ \rightarrow Ray ended up at 0. Some of it "win".

If $\frac{z_i}{u_i} \rightarrow$ constant, will $\frac{z_i(1-z_i)}{z_i^2}$ be constant? No reason for it to be \rightarrow so far results w. z_i were not equal - they were free in this one case only!

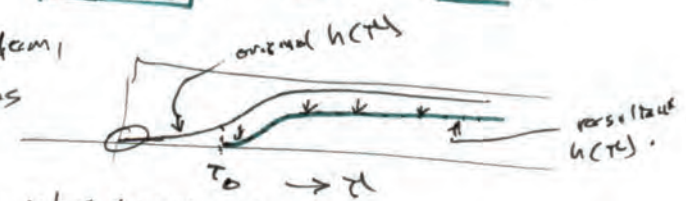
While $\frac{z}{u}$ ordering may not be true, it does have some law-like w. z I decided to use it for approx ordering of PST's. \rightarrow look at my reasoning for Ray's to have was partly the "max pc of win (per unit time)" - but I think Ray was another "proof" ("reason").

0: 3.17!

If PST₁ w. μ_1, σ_1 has failed for $CBS_1 = 3 \times \mu_1$, then PST₂ (say) which is dependent to ϕ by highly correlated w. PST₁ (wrt to same problem), will have its μ_2 moved out to $\approx 3 \times \mu_1$ at least.

T. reasoning here is unconvincing: we are sure now that PST₁ will fail for time $< 3\mu_1$. (probly = 1).

Say $h_1(t)$ curve for TST_1 was



If we work prob for time T_0 , what does it do? curve then look like?

It would seem that it would be ϕ for $T < T_0$.

If we look at $h_1'(t)$ maybe its clearer.

Looks like $h_1(t)$ comes just toward so far $(T \geq T_0)$ W.A.W.A.L

$h_1(t) \rightarrow h_1(T) - K \neq 0$ for $T = T_0$.

This becomes true also, for any other ST highly correlated w. PST, unc prob, T. above results (-10ff) and for $\sigma^2 > 0$, if $\sigma^2 = 0$, then it

\rightarrow PST₁ fails for $T > \mu_1$, then it has failed completely.

N.B. T. long. discuss. looks related to WON!

If we approximate $h(t)$ by 2 S func's:



If we fail in $T_1 \leq T < T_2$ then $h = \frac{z}{2}$;
 " " " " $T \geq T_2$ $h = 0$.

we can approximate $h(t)$ by any no. of S func's.

30

10 : More detailed inserts for Kol Talk

slide 15

Insert 1. ① ~~data function of strings to strings X~~ $X = M(S_i)$

$L(S_i) = \text{length number of } S_i \text{ (no. of Bits)}$

$2^{-L(S_i)}$ is probability of input S_i to $M()$

$P_M(X) = \sum_i 2^{-L(S_i)}$ is total probability of all programs for X
? still "invariant" it includes partial recursive funcs, & ~~partial recursive funcs~~ ~~swelles recursive func.~~? No.

10 ~~Slide~~ slide 2: add: Converges faster than $\frac{1}{n}$.

20

30

109 9

Perhaps in A.I. $\frac{1}{2}$ of programs CS
 in SW, $\frac{1}{2}$ in HW!
 Maybe for Seminar Talk!

20:

Model Uncertainty:

List of items I want to include, that I should "write up":
 1) On ^{Insert at 3.09} "Subjectivity" accuracy: 8.17-.70: meaning of SSZ v.s. Model uncertainty
 Give ~~example~~ example.

06

2) "Subjectivity" 7.10-8.01 2) Defn. "Info before data is ^{known} seen"

b) Q: philosopher: Zero data --- What what? (NFL Thoms) ^{buffer lunch}
 analogy for "no Air" → no Reg.

0:

How massive form w/ exp. info. Talking, walking, synthesizing of v. visual/auditory scenes
 Inserting into into U.B.F. ← it's hard to describe into detail.

Maybe another C++ (1.25, Fourier, Lisp, Maple, Mathematica libraries of functions,
 ALP as compact way to represent data.

In Living ~~creatures~~ ^{creatures}: Where does a principle come from?
 Idea of evolution! ^{process in response to environmental stimuli} A (to use suitable "CS" of events, challenges, by intelligent
 would not ~~have~~ ^{have} developed.

16

"An Epistemic Principle" in Epistemology ← perhaps rapid "Throw Away" time.

0:

3) OOPS: First discuss Phase 1 v.s. Phase 2. OOPS P is debts TOH
 proposed ~~subject~~ for Kol lecture; ^{Mano detail for Seminar} "1" "2" "3".

possibly explain of ~~Boost~~ "Boost"

Make list of insts used in 1st 2nd; TOH: How Boost was effective.

Inherently Boosted for both: ~~insts~~ ^{insts} for (1st 2nd; ~~insts~~ ^{insts}) (3rd 4th 5th for TOH)

Insts boost the both: byz, dec, boostq

+ Boost for 1st 2nd only c1, c2 Boost for TOH only c3, c4, c5

0:

Sets of 1st 2nd: [dec, byz, boostq, dec, byz, boostq] 5 insts.

Soln. of TOH: [c3, dec, boostq, dec, byz, boostq, c4, colltp, c3, c5, colltp, endup] 10 insts

Only use boosted insts! So from original PC = $\frac{1}{23}$ per inst → $\frac{1}{7}$ or $\frac{1}{12}$.

6¹⁰ = 6 = 2.3
 12¹⁰ = 6.7
 17 $\frac{17}{2.3}$
 36⁵
 36. 384
 = 36 x 10⁶
 3 x 10⁷

OSL .00, but .02-.06 (looks much more promising!) One Stat Lang

0: : **[SN]** Could I list various rtues or functions in "sig list", but w. relatively low pc's.
These would be functions that have occurred only once, say, & are somehow to be used in OSL.

02 Tho, normally, if may OSL is implemented: A part of a function occurs. Then we look
in corpus, to see if that part has ever occurred in fi. past. — if so, we can try OSL, to see
if it's ok. quantitatively.

04 **[NB]** .02 looks like L-2 compression routine: so we might be able to devise a
06 very fast implementation! May apply to routine to function trees, subtrees, hvr.
10 would Polish Notation catch much of it?

12 **[What happened in OOPS]**: (How T.O.H. ~~was~~ obtained factor of 1000 in speed, from
14 Sarnot 1724.)

T. ~~initial~~ boosted by 2, dec, boost: ~~made this 3 insts for E and P into 3 insts~~
24 must all wk. Then boost of 3, c4, c5 included except in the "effective set of insts"
So we have pred'n "c3, dec, boost" fairly likely: The new ^{additional} boosted insts
26 are now **[def np, c1, collp, c2, end np]** ← successful seq. comp. gain to 1724.
from this, the sequence **[def np, c4, collp, c3, c5, collp, end np]** has an acceptably
30 hy pc. ~~successful seq. soln to~~ successful seq. soln to T.O.H.
could be used to construct fairly complex recursive functs.

Note that "def np" & "end np" always come in pairs, at the start & end of a funct defn.
— One should be able to build this fact into. Lang. used by ~~the~~ OOPS.

Problems for diagn. of model uncertainty,

26 → Let me explain: ~~But first a brief~~ discussion of uncertainty in prediction.
30 into 2.27 There are two kinds of uncertainty in statistical results: The best known is
uncertainty in probability values due to finite sample size. If you have
40 ~~one~~ bits and half of them are zero, then the probability of a 1 being
the next bit is about $\frac{1}{2} \pm \frac{1}{\sqrt{n}}$. The larger ~~the~~ the sample size n, is,
the less error in our probability estimate.

The kind of uncertainty I'm trying to talk about is not ~~sample~~ due
to sample size but due to "Model uncertainty." ~~When~~ When analyzing empirical
data, there ^{are} normally an infinite number of models that can be used to
analyze the data — some will give good predictions, others will give
poor predictions. ~~It is~~ In any finite amount of time, one can only evaluate
a finite number of models. From the ~~unique~~ ^{finite} ~~strip of the~~
universal distribution tells us that whenever we have ~~one~~ ^{find} a model $\gg 1200$

Things to be typed

20

for P. 1.06: Suppose we have a finite string, x and we want to know ~~the probability that~~ ^{its probability in universal probability} ~~it will be produced~~ ^{it will occur} with respect to machine, M . There will be many inputs to Machine ~~that~~ ^{that} x is output. Say S_i is the ~~input~~ ^{input} such ~~that~~ ^{that} x is output. If S_i is $L(S_i)$ bits ~~long~~ ^{bits} the probability that ~~x would be produced from a random binary input is just~~ ^{the probability of} $2^{-L(S_i)}$. To get the probability that x will be produced by any of the programs S_i , we sum ^{all of} them up.

07

down to $\sum_i P(S_i) \leq 2^{-L(x)}$. \rightarrow 13.00
 (H. Shannon S. 1.5) \rightarrow To do production with $P(x)$
 distribution is fairly simple.
 If x is a binary string

law of evolution
 application of this
 led to industrial
 revolution.
 Had did, however,
 public it for my
 discovery
 many years it was
 better known in
 E. S. U. 1911

0

for P. 2.27 after "feature": 10.26 - 1.40; 12.00 - .03
 (expresses Model uncertainty "abit")

3

→ In a finite time one can only consider a finite number of models.

20: ~~10.40~~ But surely for the data well, we can't be sure that there is not another model that ~~is~~ will give better predictions. There is no way to avoid this. The incompleteness of the universal distribution is related to its evaluation ~~of all possible models~~ and taking ~~to~~ in finite amount of time to do this.

03

3.26 of talk: Discussion of Subjectivity

Not only a need Bayesian statistics to subjective, but it is desirable that any other method of inserting a priori info, be subjective.....

10

List a few Desires: ① Accuracy ② no need to divide data set into training set and test set ③ Never problem of overfitting or underfitting.

Apparent disadvantages: ④ Data need not be "stationary"

① Incompleteness, ② Subjective

Advantages: Includes pr models. [OOPS is only game I know that searches over p.v. Models. Some G.P. systems might.

For Random variation of L search, use (shot N.U.) to get fixed time spent on each trial! usually in computers, getting a fixed time interval, grows known cut off, is difficult.

20

for 4.31 of kolTalk. When people see the expression $T \cdot 2^L$ is the CJS of the problem, the time needed to solve it by search, they often become frightened.

People are often frightened by the expression $T \cdot 2^L$. It says that the time is exponential in the size of the solution. For many people this is equivalent to unsolvability. However, in A.I., the sizes of practical solutions are usually quite small, and the dichotomy between "polynomial" and "exponential" is not a useful ~~distinction~~.

30

If all the information one has is in the probability distribution, and the only way you can control the search is by probability and time for a trial, then the absolutely best way to order trials

13 $\frac{T_i^s}{P_i}$: smallest first. This will give the best performance. There is no better way to solve the problem.

Since the solution search is a good approximation to the optimum technique.

40

To repeat, $P_M(X)$ is the sum of the probabilities of all programs that could have produced X , ~~with machine M~~.

00:11:07

~~To use this distribution~~ It is easy to use this distribution for prediction: If X is a binary string, then the probability that 1 will be the next symbol ~~of X~~ of X is just.

$$\frac{P_M(X1)}{P_M(X0) + P_M(X1)}$$

$P_M(X1)$ is the probability assigned to $X1$ by the universal distribution, using machine M as reference.

2 1/2 min

The accuracy of the universal distribution associated with the accuracy of the universal prediction using the universal distribution with the accuracy of the universal distribution as a predictor is ~~highly~~ certainly not an important feature, there are other important features:

Some other important features:

- We can do prediction on non-stationary time series.
- ~~This technique is~~ we obtain an optimum solution to the overfitting (underfitting) problem because of "optimum fitting", there is no need.
- If ~~it is not necessary~~ to divide data into "training set" and "test set".

~~The error~~ we obtain very good error estimates for the future by analysis of the training set alone. Since we do not assume stationary data, we obtain a probability distribution for the future that is not necessarily correlated with the past.

~~Probability distribution of the program~~

- It is possible to use partial recursive functions to model our data.

While all of these are very nice, none seemed at first to be a serious problem: — that the universal probability was incomputable.

To my knowledge, no one has actually tried this. But the system I will describe later, will do it. What it gets better results from than using only recursive models, it turns to be from A.

The way I saw

Our "easy to describe" version of L search

possible candidates ~~in parallel~~ in parallel — by time sharing. If the universal distribution assigns P_i to ~~candidate~~ X_i then it will give a fraction P_i of the time share to that candidate.

While this version of L search is very fast, it uses lots of memory. Another ~~slightly~~ somewhat slower version of L search uses much less memory.

23: insert in 4.04

Why this is the ~~best~~ best solution to the problem! The kind of work that I do is to solve it.

- 1 6'
- 2 2 1/2
- 3 2:50
- 4 (5:10 / 1:40)
- 18' = 2
- 5 2'40
- 6 2'10"
- 2'45"
- 7 6'
- 3 1/2 min.

✓ I may not want to include this — I miss go on p 4 when I introduce PST's

10:

: ~~A PST~~ Given a problem, a PST is any algorithm that attempts to solve that problem.
 The domain of a PST, is the description of the problem.
 ~~The domain of the PST, is the description of the~~

problem. For ~~inversion~~ inversion problems as we discussed before, in which we are looking for a string x , such that the given function $f(\cdot)$

~~exists~~ ~~I know~~ Inversion problems are defined by a function, $F(\cdot)$ or number. For example, and a string s . The problem is to find a string x , such that $F(x) = s$.

If $F(x) = \sin x$ then to solve $\sin(x) = .5$, ~~what is x ?~~ While search itself is ~~a particular~~ one kind of PST that

0

can solve inversion problems. But there are many other PST's that can solve ~~inversion~~ ~~problems~~ ^{them} as well.

~~Another kind of~~ Optimization ^{defines} another broad class of problems.

Given a function $G(\cdot)$ that maps strings to reals; to find a string x such that $G(x)$ is as large as possible. ^{In an} important subclass of

optimization problems, ~~where~~ it is necessary to find the best x possible in a certain time limit, T . The problem is then defined by the function $G(\cdot)$, and the time limit, T .

10

There are a great variety of PST's that have been developed for optimization problems ... Artificial Neural nets, Simulated Annealing, and Genetic Algorithms are a few fairly general PST's of this sort. Zmin

36 [Insert at 8.32]

On OOPS: A big difference between OOPS and my "Phase 1", is that OOPS has been actually programmed, and an interesting result was obtained, showing the utility of forming sequences. In the particular set of problems ~~that~~ it solved, there was a factor of about 1000 speed up in a solution, — depending on order in which problems were presented to the system.

0

20. ① Read 2 part of abstract on How system works. P's 5, 6 & perhaps 7.
(Abstract taken from NIPS workshop Handout)

0 → 1) Define Inv, Oz probs, (from Sol 86?) | Eqs on slides

2 → 2) Descr. 1, 2 or 3 kinds of batch to solve ~~INP~~ INV, Oz probs.

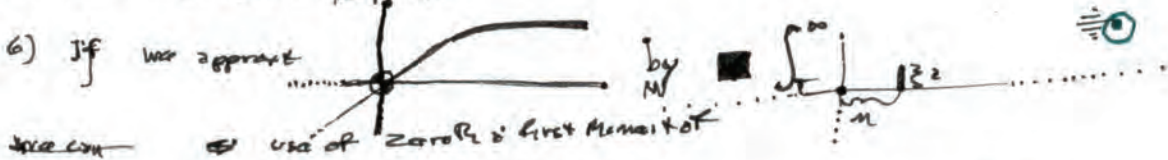
5 → 3) Discuss Phase 1 update: use of focus of tokens & of used eqns & used subfunctions, subeqns.
~~3) Problem is to find of solving; then "Updating" is support to do this.~~

c → 3.5 Descr. Soln of Operator induction ~~is not needed~~

Give Eqs: In Box forms $I = \sum_{i=1}^{n+1} \dots$
① $\sum_{i=1}^n \dots$

d → 4) Go into Phase 2 update: Show detail (derive both for h and h'
(i.e. negative as well as positive cases).

e → 5) Give "Integration Equations"



use of zeroth & first moment

df.

10 ... We can quickly obtain probab of each D.F. leading to ~~shortest soln~~



better approx by use of 2 point or 3 point D.F.

7) OOBS: Differences & Similarity to Phase 1 -
Optimum (Order problem solver, Derb. inst sets)
Shell/shell/func/func v.s. LISP

The 2 problems TOH is 1^2^n.
Third TOH for 2^n: no recursive soln.
Third 1^2^n get soln. soln after n=1,2,3,4. (recursive soln)
Third TOH = recursive soln for n=1,2,3,4 recursive for n=3.

How "Boost" works mechanics

"we got kind of "Mutation", but use learn rather than random change or replacement,

2 Boosts could give crossover - ~~but its rather wasteful because~~
many cases will end up with identical data.

Differences betw. Phase 1 in Oops in TSC construction. (perhaps)

Maybe not so wise to give details on many problems! Best to give details on a few problems.

0

α

β

γ

8

30

Methods for Kol talk.

P1 Convergence Theorem: If there were any regularities in a data set:

Univ. D.F. would eventually predict ~~using a relatively small amount of data~~ it would give about the best possible probability values for a given sample size. using a relatively small amount of data.

First convergence theorem was for self-similar data - like Time Series

Analysis for normalized models:
Lester Peter Gacs showed.
Marcus Hutter
- on order sets of strings
Question: Permuting System.

rich 130

Discuss Kol's
not having discovered
Algorithm
User name: Grey
Password: Piyoyek
Host name: Sante
URL: ~~http://www.sante~~
http://www.sante
It's server has
usable office
fact-names!
ray@sante

Discuss Kol's not having noticed the application of induction to ~~complexity~~ inductive inference.

Discuss Criterion for Goodness of induction:
rigorous?

1968: Willis paper.
How it was ~~very exact~~ very exact, recursive developments
of earlier paper.

win scp - double click
" of induction. ~~from discovery of proof of goodness~~

~~page Sante~~
~~password~~

Proof for time series induction
using Norzd & Univ. D.F.

Peter Gacs showed convergence theorem for unnormalized Univ. D.F.

How to have Univ. D.F. ~~convergence~~ interest
Mathematicians happy - ~~when accepted~~ accepted
reason - Levin

But in general, while there may be a "best" Normal D.F. - any ~~other~~
measure of D.F. is very much better than ~~any~~ unnormalized
measure it is derived from. The Normal measure
from a practical viewpoint, it's always better to use a normalized
measure for prediction.

~~Security~~ Security M. Hutter: A-by alphabet,
Various different loss funct. - over of them was of ~~much~~ much utility.

2.25.03

For P18
SOP 5+

17

Discuss. of Encew pitibility. — lower 25 13.

644

(18 missy?)

inserted.

Inc. in property of ~~forecast~~ ^{averagability} ~~en~~ ^{average} probability itself!

That it is not a small effect! That is water how many models you've tried

That ~~there~~ ~~is~~ ~~no~~ ~~one~~ ~~best~~ ~~model~~ ~~that~~ ~~is~~ ~~known~~ ~~to~~ ~~be~~ ~~the~~ ~~best~~ ~~one~~ ~~yet~~ ~~found~~ ~~and~~ ~~you~~ ~~would~~ ~~find~~ ~~it~~ ~~if~~ ~~you~~ ~~spent~~ ~~10~~ ~~minutes~~ ~~more~~ ~~looking!~~

Subjectivity — ~~Normal~~ ~~Statistics~~ ~~are~~ ~~Bayesian~~
7.10 - - 8.14 9.06 - - 10

Δ6 : **OOPS!** Computer ~~ratio~~ ratio of $P_C \uparrow$ for $1^2 \rightarrow 10^4$ v.s. 10^4 first
 as a function of no. of insts: say 15 insts actually used in trials: so $73-15=58$ not used.
 73 to "15"
 so let N_Q go from 10^4 to 15 & see how this ratio \downarrow ~~from~~ ~~from~~ 1000 down.

Normed Unvl. P.D. Go over that argt. Part ~~Normed Unvl. P.D.~~ Normed Constant for Unvl. P.D.
 Must $\rightarrow \infty$. This now seems impossible! i.e. If we have $\mu =$ normed measure $\approx \frac{1}{P_M}$

$P_M =$ unnormed semi-measure. Then $P_M < k \mu$ if μ has a finite den.

If the normed constant is bounded, then... Well, it can have different Bnd's for every
 different μ . say the p.c. of μ wrt M was P_0 , some normed constant could be

P_0^{-r} , r is a real > 0 . so "length of unit" of normed constant is $r \times \text{length of } (1/P_0)$

A large normed const. means that μ is "very often" as U ; but we know $\lim_{i \rightarrow \infty} P_C^i \rightarrow 0$.

However μ never has M^s . P_M be > 1 ? No!

The ~~point~~ weird thing is that if the normed const $\rightarrow \infty$ would it

I think it is true that the Normed constant can be very large (Solovay): in which
 case, the normed P_M can have much more rapid convergence to M than

to semi-measure.

But I really have to go over this!

It seems to say that $-\ln P_0$ (in the convergence) approaches ∞ !
 and clearly can't be > 0 !

Suppose we have a probabilistic algorithm that can be described in a certain finite number of bits, and this algorithm produces a long sequence of symbols according to its probabilistic rules. Then ~~if~~ we have a ^{general} induction system that gives probabilities for each symbol, in terms of the previous symbols. For a good general induction system, and a long enough sequence, the two probabilities given to the symbols by the 2 different methods, should be very close.

While the formulation ~~was~~ this criterion seemed reasonable, I was at first unable to prove it.

In 1968 I was asked to review a paper on Inductive Inference, by David Willis. It was ~~very~~ ^{difficult} though I was familiar with his ideas, I found the paper difficult ^{really} to read. It took me about 6 weeks to read the paper properly.

Willis had taken my unimproved system for induction and made it into ~~an~~ an exactly rephrased system. He had an ^{error} criterion it satisfied, which was a kind of error ~~criterion~~ but it was certainly not ~~strong~~ good enough to convince a person that the system was ^{always very good for prediction,} ~~any good for prediction.~~

He showed that the average ratio of the ~~correct to the~~ correct probability to the estimated probability ^{approached} ~~approached~~ zero as the length of data ^{sequence} increased. — (The individual probability ratios could, however, be quite large or quite small.)

I was then able to improve this result to show that the square of the differences in probabilities between the correct and the estimated values ^{the expected values of} approached zero as n increased from n .

I called this the convergence theorem.

Rob

really made it clear that the universal distribution gave very good probability estimates.

Perhaps show eqn.

When I sent in ~~my paper~~ ^{my} ~~paper~~ ^{recommendation} that Willis' paper be published with no revisions - but other 2 editors who ~~initially~~ were really quite ~~opinion~~ ^{opinion} ~~in this~~ had rejected ~~it~~ - Ray felt that it had little to add to my original paper!

As an aside: I later found out who the other 2 reviews were. - Ray was both ~~very~~ ^{very} expert in the field of the paper. I wrote Willis telling him what a great paper it was and suggesting that he sent it to another journal. He did that and it was published 2 years later.

Willis' paper may have been a bit ahead of its time:

About that time, the interests in the mathematical community were in what one might call "True Prob"

My original conference program was for a Normalized Measurement sequences of symbols.

Many mathematicians were unhappy with my ~~definition~~ ^{definition}

Normalized Distribution.

start 26 20

Introduction:

I'm going to start out with a description of the Universal Probability distribution - some and some of its properties. - Then a brief description of a system I've been working on. In Machine Learning - a program designed to take advantage of the described features of the Universal Distribution.

after incompleteness.

Insert on Subjectivity. 4.27

Another ^{apparent} difficulty with the Universal distribution ~~is~~ ^{is} its Subjectivity.

last
paper.

When ~~the~~ the Universal distribution is mentioned, there are 2 possible meanings of "Universal": First that the error will converge to zero rapidly if the algorithm generally the data has a ~~finite~~ small finite description. — This is true for all such ~~general~~ algorithms. This is what I mean by Universal Distribution.

Another interpretation of Universality is that we can use the same ~~distribution~~ Universal distribution for all problems. This is what is called a "halt Truth" — The ^{same} Universal distribution will work for all problems, but for most, ~~it~~ it will work poorly — the ^{errors} ~~will~~ will converge very slowly.

When I speak of a priori distribution or a priori information — what I mean is information available before the data in the problem is known. The universal distribution, and the information it contains, changes during the life of the statistician. → 7.14 - 29, 8.02 - 14

So? The subjectivity ... The fact that it is based on choice of which Universal machine to use — is characteristic of all prediction systems based on a priori probability distributions. The choice of universal machine and its instruction set is a necessary parameter in the system, that enables us to insert a priori information into it.

So ~~the~~ ~~subjectivity~~ The dependence of the universal distribution on choice of machine is not a Bug in the system — it, too is a Feature!

20

: Addition

For Seminar: Write a bit about OOPS. Just what was done, ect.
See p 15 for More on Seminar, including OOPS.

for End of kal talk:

Derb. Motivation for inv / 02 probs: Best Maxwell, Simon only discuss
class of chv. probs. — Best inv + 02. covered almost all problem types

Especially work on inv via probability!

(1985 All reason not physical relevant to A.S.!
How they move is rite. (But even that was wrong!)

How the "Updating" of which PST to use, also how to use
new PST's: U. A.V. D. can be used to get a good solution to
How best to update

How previous methods of updating just tried to 'do so what worked
before' — But we can actually get machine to search for
a better way to update.

All more Goodies are a replacement of ~~U~~ Univ. D. f.
(Probably we now run on A.S.)

For Seminar: More on OOPS (See p 15 on Rtg)

2.27.02	50	1.47	170 + 43	170 749
			3 to 58	3 to 55

Q: in
Phase 2 update:
any way to use
very approximate
2-param h functs
to use Negative
data ?

0

10

10

00 The system was designed to ... NIPS P1, P2, P3

NS-16 - ~~to end~~ to end.

Much of my work of ~~the~~ recent years has been in understanding developing and understanding the existing system that enables the General Comb. Prob. Distribution to learn from both ~~both~~ successful and unsuccessful ~~to~~ problem solving trials.

The last talk I gave at Royal Holloway was at a Symposium on the Importance of Being Learnable. It was a description of some ideas I had on ~~about~~ transfer learning - how learning in one domain ~~was~~ could utilize information from ^{other} apparently disparate Domains.

26

enables transfer learning both direct and transfer learning from both successful and unsuccessful

Subjectivity - That it is much dependent on ~~what~~ just what universal machine is ^{chosen} ~~used~~. While ~~that~~ it is known that differences in universal machines "only" introduce differences in probability that ~~may~~ vary by a "constant factor". This "constant factor" can be about 10¹⁰⁰. ~~The subject, while the constant factors of the order of one "it is not small."~~

30

Seminar: Quick fix!

Say I will describe a very general Prob. solving System that I've been working on for some time. P's 1, 2, 3 from Advanced | IRV, 02 probs. Then talk about proper Lsearch & its properties.

Another fool in our Problem solving Techniques is Learning Universal Search Algorithm which we will call LSearch.

00 : 10:47
 49.
 55.20 and 057
 11 00 Start of P₁
 11 12 Start of P₂
 11 16
 30 ~ 1/2 40.

How OOPS works! Stack lang - like forM. 2 or 3 stacks.
 Tokens are instructions but how tokens can be defined by program. ~~As a functionally~~ P₁ M = seq. of tokens.
 2/3 tokens to start. P₁ = 1/3 to start! P₁'s using
 - 1/3 Laplace rule new, down down.
 Each successful program is stored in Frozen Memory (ROM)
 can be referred to:
Boosting!

Next slide

c1 ... c5, 1/2, dec, boosting.

How OOPS works:

This Stack language similar to forM: 2 or 3 stacks.

Tokens are instructions! New tokens can be defined by programs.

Programs are sequences of tokens.

2/3 Tokens to start: probability of each to start = 1/3.

Laplace rule used to update probabilities:

Numerator = 1 to start down down down = 2/3 to start.

Increment numerator, down down down denominator after each use of token.

Probabilities of program = product of probabilities of its tokens.

Probability used to guide search.

Each successful program is stored in "Frozen Memory" (= ROM)

How Boosting instruction works!

Booster instruction looks on stack, gets address of program in ROM:

Boosts all tokens in that program by 2/3 (~~Down/Down~~ (Both numerator and denominator))

Boost is a "Mutation"

10

30

06

More on "In computability: Perhaps Give Example of MDC:

Pick a ~~set~~ ^{subset} of models: Pick best of those models (must be.)
Model error: set of initial models ~~is~~ ^{is} very probably did not include
model that is much better than any in the chosen set.

07

Use of Univl d.f. gives us a Pragmatic awareness of possibly very large error (2) Suggestions on how to reduce prediction error

08

Suppose we have a data set and do prediction. ~~As an example, we might choose a set of models~~ ^{on the} ~~from use of MDC on a dataset~~ ^{we pick the best} ~~or some other criterion~~ ^{model in the set using MDC. No matter how well the model seems to fit, it is ~~likely~~ ^{possible} that there is a better model outside the set of models considered, and ~~we cannot estimate the probability of this being true.~~ ^{there is no way to estimate the probability of} ~~it being true.~~ ^{this being true.} It gives an error of unknown size.}



10

SN

On Learning during Lurch, correlations betw. conds.

Say Cond_i has up to time T_{in} failed to solve probk.

If we can use this fact, to recompute the h curves of residual for Cond_i wrt probk. [If there are a bunch of conds that are "correlated" w Cond_i (wrt probk), Prms Prgr h curves will be the main ones updated — (in fact we may only look at the modulus of the h curves of conds that were "correlated" w Cond_i).

(Anyway: After Prms Prgr (partial) h updating, we will tend to find only conds not correlated w Cond_i. So Prms Prgr is ok.

desired "Say prms Prgr is ok" — long during Lurch.

12

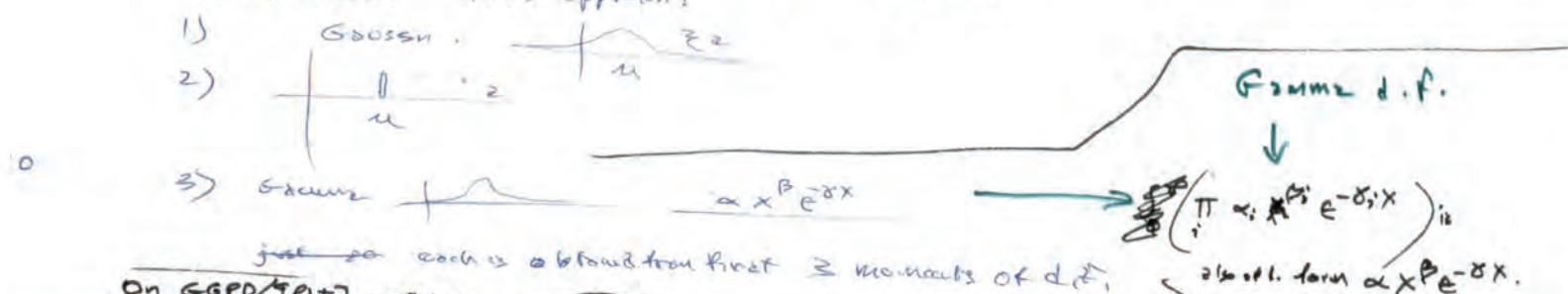
13

loop

Note Bare: The "correlations" of 12 are empirical results when h curves are recomputed. They are based only on the same old data as (PST, Probk, Tok) triplets of the past is option of JOY.1

20 : Make Expanded version of "T. report" with various improvements & discuss: Like for 2 param version of h, h' & how to use it for $\beta \in \mathbb{R}$ positive & negative cases.

04 Also 28, 20-40 ^{29.12-40; 30.30} looks very impl. & should be included. Part 4 "even now" expand it: Make sure I understand it.



On GGP/PST updates: 21.20

12 : 28.40 : More exactly: After unknown T_{ij} for failures of $PST_{i,k}$ (It has failed for $T < T_{ij}$), we have strong new data for

So we can get 149 rds "easily"

$$\sum_j \left(\prod_i \pi_{i0}^j O^j(T_{i,k} | PST_{i,k}, Prob_{i,k}) \right)$$

$$\sum_i \prod_i \pi_{i0}^j \left(1 - \int_0^{T_{i,k}} h(t) dt \right) O^j(T_{i,k} | PST_{i,k}, Prob_{i,k})$$

failure prob = $\sum_j \left(1 - \int_0^T h(t) dt \right)$

is prob of all known failures.

We multiply this by π . ~~original product~~ ^{seems better now failure prob} for success of cards (= PST's).

As we reoptimize ~~entire~~ O^j to get a new O^j on set of O^j 's:

This gives now h 's for many of the PST's. The PST's that happened to be "close to" PST_i (f. known by T tail we for problems) will usually get worse h curves (smaller $\frac{1}{\mu}$ say) than before, & so decay: its more likely that ~~the~~ PST's not "close to PST_i w/ prob k) will ~~be chosen~~ have the most likelihood of success for k next trial.

If we repeatedly chose most likely PST then w/o on it until (by simultaneous updates), the most likely card \hat{c} has been changed, & jump to new card. Then \hat{c} looks no more like Lsren, but more like WON! It certainly is it Lsren. \rightarrow 30.30

The Universal Distribution and Machine Learning.

insert title after title.

Start of ~~lecture~~ talk.

There will be two main topics of my lecture:

I'm going to talk about two main things: First, The universal

distribution and some of its properties: Its accuracy, its in-computability, ~~and~~ its subjectivity.

Last: I'm going to tell how to use this distribution to create very intelligent machines.

3.27 replaced by 30.20-22
3030 replaced by 30.10st
5.11 insert
28.08-16



It is often possible to obtain predictions using a truly a priori probability distribution (obtained before the data was known). Under these conditions there is no under-fitting and over-fitting - the data need not be divided into "training set" and "test set" - All data can be used for training and the expected error will be unbiased.

The data need not be stationary; ^{sub}sequences of data can be missing; The data can be multi-dimensional - extending finitely or indefinitely in all positive and/or negative directions.

3.3.03	00:01	300 150 30	530 150 60	3.4.03	320 200 32 57	190 720 48 47
--------	-------	------------------	------------------	--------	------------------	------------------

29.40 is its probably much better than Lach.

T. (apparently) BAD aspect! It seems too "elitist": Very little random, crazy ~~the~~ exploration.

Re: Implementation: In 29.eg, we will normally have several O^i 's for each "sub. str. ep.". ~~When~~ ^{When} unknowns are working on a particular PST, the relative wts of the ^{the} "Hypothesis" (most wt) O^i 's will change - This aspect of "updating" can be rapid/cheap/easy to implement. Actually finding new, better O^i 's is, of course, Much Harder.

20 : Re Grad students, ect! What are major problems in α ?

Easiest to work on is phase 1. To gain simple TSC's. Look at CJS's. Try to find Corollaries (at all levels, of all kinds) to α CJS. (Also develop a perhaps "R" system (recogn. functions) or/ better, using Corollaries.)

Perhaps actually gain "SaarbTM" for ANh. Argument to inst. sec: See if ~~we~~ we can use corollaries to get acceptable CJS's. $\{ \text{It may be that we simply need larger CJS's} \}$

Try to get at W. Publ's student's work. (Peter Bergmann). It's probably in German, hvr. Maybe one of MIT grad students can read German. To do this, first make up bunch of problems. Then to "Necessary" forms of α words will become clear. $\rightarrow 37.00$

\rightarrow Work on Form (S) for Sfuncs.

16 \rightarrow Some imp't, unsolved ~~and~~ incompletely solved problems.

1) Test L svcd ~~ones~~ α pc of α given cond. - but more can be many ~~equiv~~ equivalent Conds - would like to get a set of equiv cond's, one trial we vary large α (α no of equiv cond's). At present, nowhere to do this except: Use of Lops rule to assign pc's to Turing's α perhaps $\rightarrow 32.00$

19 SN: For my versioned report, write up just how SCPS/PST's complex, updated "as a whole": how trainer can help system by inserting factor set of v.g. PST's, ~~to be~~ extrapolated. NB: we are mainly interested in opten. problems α Rate over many denom.

28 For Col Lecture paper: PC error
Any finite derivable. PEM that gets an idea of α of later than $\frac{1}{2}$ must be for all finite derivable ~~data~~ data sources used
for incomputable. (If you ~~use~~ data sources buffer from random for all finite derivable data sources, must be incomputable.)
If α PEM gets k bits pc's v.g. for α data src.
 α is attracted to do well for δ . rest of α src. - for all data src's:
- Req α ϵ - Req α must be incomputable, i.e. given a flow w, finite denom of bits, and α data src's derivable by δ bits, we can construct α data src's that short out w. if known data src's α has error $\geq \frac{1}{2}$ for all subsequent bits reduced by δ bits.

30 ~~For~~ So! Any ADM CPM (computable proby Measure) bounded input - input PC must have PC error $\geq \frac{1}{2}$ for all bits. SPCC 35.00

Mitch
Glickstein
Glickstein
@ucl.ac.uk
University
College London.
020
7679
2000
Mick Glickstein
+44 work
020
7679
2888 (2888)
Council
M. Glickstein
@ucl.ac.uk
Internal phone
32888
~~0207~~
0207
679288
888
Mick
Glickstein

LOSS FUNCTIONS (20)

no: (Adapted) 3.1.19: Other ways to assign pc's to cards (other than "row coding")
~~can be good ways to~~ — Also t. method of 28.20 ff (long dunny Lesh)

may help, by not trying equivalent Cardy Out here fails

2) Trying to find "Conditionals" for t. pc's of tokens, to deal w. "Priz Scaling" problem. "Association", Gouzen of "Recognition Functions" in Phrasal.

3) How to get some effects of Boosty ~~is~~ instruction in OOPS, but more legitimately. When a set of tokens is "Boosted" There is an imp.

"Togetherness" (association) of t. set of token boosted. — Key "association" of (perhaps) ~~is~~ specific tokens, is not obtained in the simple AZ lang that I've desc'd.

3/5/87 **SM** On "Loss Functions": There is now Much extensive Resch on How to bet w. various "Loss Functions": This may be identical to my classic "AE" (?) Alg. Evaln. Problem! — A special case is SM strategies. A more General Case is an All-over Alg. for Running TM — to optimize some future loss function on it.

As for what Ollars are doing on 8/2: Do they assume "Stochastic TS"? What kinds of "Loss functs" do they consider? Look at that paper by

30 ~~3~~ Crasti-Bianci assoc w. t. NIPS/workshop I attended — is a review of what they were working on.

My initial/original impression is that of Kestly Hutter: That if one assumes % data will follow to Univ D.F. & bets accordingly, one will "do well".

(Hv, first calculus f. Un. D.F. (i.e. doing prediction) ~~then~~ optimizing t. loss funct., is an al. approach & would like to be able to optimize loss funct. Directly

— OR as in SM, How a way to evaluate (in some unbiased way) possl.

"Strategies" (i.e. "ways to bet"),

(SPEC 37.00)

\$3/02, 100 = \$120

- 28.20 ff plus 40.03 ^{way to solve} _{ways} _{to solve} _{ways} _{to solve}
- 1) Main problems: cross v.s. post: (special case of correlated cands)
- 2) Defining good contexts, to ↓ Scaling problems. This could be major step in Making SaaS TM (ANL) practical. → 39.00
- 3) Finding, expanding SaaS TM. — use OOPS for ideas — ^{sort} 26
- 4) Analysis of "Boost" in OOPS: can I find all essential features? — ^{sort} 26
- 5) Put them in optimized form? → 38, 30 — NONBAD, repetitions → 37.00-38.30 (looks v.g.)
- "Boosty" → Good way to implement S-FUNC
- SN In "Boosty", t. current = no. of eqvnt. repetitions of each token vs N_d .

A better way, t. no. of repetitions k (present denominator). Initially, $k=1$, but after TM has been running a while, we try to optimize it over past data.

(A very expensive option!) Try first w. $T=1, k=0, k=.5, \text{etc.}$

Examine other factors (S) of "Boosty" — In particular, t. cross-correlation between tokens & propositional. One way to do this last: If a pm has been "boosted" we immediately get its coverage, so that if any of t. boosted tokens is chosen, either in future (or past of t. present cand), other tokens w. boosted pm have ↑ in pc. We could have 2 adding effect, so that 2 off. boosted tokens are chosen, t. other boosted tokens have even more ↑ in pc.

T. idea of Boost (a kind of crossover/mutation) is that "the set of tokens tends to work well together!"

- SN T. troubles w. S-SARB TM: 1) Only final/successful pm used as deltas (no subdeltas) (Fuzzy because I used Skinnerian TSO ... No jumps of > 1 "layer of concs")
- 2) Scaling — pc's of solns began to ↓ rapidly, m. fl. var. of no. of concs. in Many — contents of various kinds could help Much.
- 3) Set of insts, was very small, not near Unit Versal.
- 4) Perhaps use ideas from OOPS

SN I may want to pm my latest version ("fix") of 2.141 — so what can put pieces into speculative English Text. Here, note that I really had no cut off — w. any corpus size, it would start to combine words that occurred often together. Corp. (size) must be large enuf to discover prefixes, so fixes and suffixes of which they are parts — but not so good that words are often combined to create ("sporadic") compound words (≠ German ©).

SN Main for talk/leader: "t. importance of being Unlearnable"

LOSS Funct (cont) .00
SM .00 ff

SM .00 ff

50: 32.90 : What about this: each strategy has $d_{i,t}$ length (units) of P_i
So w. β in Bank, we bet P_i on S_t . (This looks like \rightarrow general of Cover's
Universal Betting Scheme ... but uses a priori uniform d_i , i.e. $P_i = \text{constant}$.)

05 The trouble w. Cover's scheme (\hat{z} perhaps a fortiori, this is one) is that
it takes too long to get up to speed \rightarrow But, it may not if I assume $S_{i,t}$
is not stationary (\hat{z} look for reg. in its non-stationarity decr), AND,
do parallel betting w. other stocks \hat{z} allow cross-stock information.
10 [This is a v.g. Gouze of ~~multiple~~ multiple time series Analysis implied by
the Mxd Corps Thm]

The idea was that the stocks could have a certain amount of "shared" info,
plus a certain amt. of individualized info for each stock. A better way
(perhaps) for many stocks would be to make trades w. each ~~stock~~
stock having a "lot":



$\alpha \rightarrow \beta \text{ share } \delta; \alpha \rightarrow \gamma \text{ share } \beta$
 $\delta \rightarrow \epsilon \text{ share } \gamma; \delta \rightarrow \zeta \text{ share } \alpha$

Are there more General ways to dec. how info
is shared?

This looks like a v.g. way to deal w. the non-stationary mess \hat{z} t.
"Small SSZ" available in SM!

20: 05 That "upto speed" may not have been the problem; T. problem with how been
that yield/yr. was much less than my approx. scheme (strategy).
Perhaps at present time, such low yields would be acceptable, since
30 SM yields (as well as other security ^{best} yields) are very low now.

31: 01 Some "strats" could have large negative yield: It may be hard to fix strats so
they have bounded neg yield — ~~rather unusual~~ (minus original Bank)

So we are sure our yield will be large. — But it was bet in log domain \hat{z} after
"Gouze's Rule" we can't have yield $< -(\text{original Bank})$.

30: 30 Actually, the covers final yield was is better than some ~~strategies~~ (a stock yield).
I tried to converge water to that was slow — slowness of size of upprofit space,
which was \hat{z} dim object (M stocks). — The space of $m-1$ dimensional w. vectors.
It would converge faster if we bet on each stock (same amt.), or bet same amt. \rightarrow (Spec) 36.00

0 (3.1.40) There should be a nice compact way to say this:

- 1) for any CPM, there will be finite desirable sequences for which that CPM will have errors of $\geq \frac{1}{2}$ in prediction: i.e. ~~worse~~ equal or worse than Random choice meta
- 2) This is not true of in computer program measures: ex. Conv. Perm. errors

errors fact $\rightarrow \phi$:

How ALP differs from any particular CPM₀:

Say we have a particular CPM₀. It has been obtained w. data set D₀. Using ~~any~~ Typically, it will have been obtained by dividing D₀ into a training & test set, & will ~~then~~ give a certain value for expected future error, ϵ (assuming ~~that D₀ was stationary~~ that D₀ was stationary)

Using Same Data, ALP will obtain many CPM's - among them, perhaps, ~~some~~ CPM₀. Its predictions will be close to those of CPM₀, but it will have used all of D₀ for "modeling". ALP The expected error ϵ gives for the next prediction is unbiased. It makes no assumptions about ~~the stationarity of D₀~~ D₀ is stationary. Tho its error estimate will usually be smaller than that of CPM₀, and be "unbiased", it will be on average such that the true error due to model uncertainty may be much larger (?). ALP If more data is given, ALP will be able to make changes of wts of its models, in accord with this new data, usually

usually giving a better error than CPM₀, which remains invariant as new data comes obtained. **N.B.**

Go we could modify training set for CPM₀ & see better CPM₀ - no chance to generally experiment.

If time is available, ALP will be able to suggest new directions

to search for better models. \rightarrow (need not use "try, fast" routine)

Main Difference, ALP gets somewhat different models from CPM₀ - but smaller, unbiased error, is more robust about error claims, but that is CPM₀ matter. Also suggests new areas for exploration.

SUMMARY

HMM! in \mathbb{R}^N ALP, + search can be fully "biased" - while ~~any~~ ALP may

have unbiased error, bound $CC \ll \infty$ ALP (= RLP) can be very biased, I should think.

3.5.03
3.5.3
3.5.3

SM

10:34.40 exact m^2 pairs. This buys $\frac{1}{m^2}$ in each pair, then used Cover's (unchanged) system on each of the m^2 pairs. The idea here is that the densities of 4 species taken, is much higher than in Cover's original scheme.

(SN) In cover's original scheme: A possible way to compute it: for $m=10$, using 10 values of each vector component a_1, \dots, a_{10} ; So 10^{10} components to be updated each day. Certainly poss. is perhaps not to time consuming.

$m=8$ would certainly be feasible this way.

Using present day CPUs, it might be possible to update 8 pts. simultaneously!

An even faster convergence to buy $\frac{1}{m}$ on each of m stocks, "Buy & hold".

A classic strategy. — Final yield will not be as large, of course.

So, it may well be that by using sub-compl'w.

$m=1$ or 2 or 3 , (rather than the full no. of stocks being bet on) we will get faster convergence rates — not to mention a better

way than using the full set (is perhaps) much more reasonable computation time.

3-6-03	0.000	1.75	1.27	1.80	2.57	2.40	5.62	2.40	2.63
	$1\frac{2}{3}$	(1.05)	1.20	(1.35)	2.12	(1.15)		$3-\frac{1}{9}$	(1.03)

(SN) Res Advantages of ALLi Heuristics keeping several hyp PC codes —

Enables one to quickly pick another code when the ~~first~~ formerly

EMPT/ Best one "fails". (When this is done, one must recompute supply of all codes —

Look into just how precise (Physics.) — I did look into this some time ago, but

in view of its importance — look again! — It might be a fractional of long incursion

0:21.06! On Secrets!

02 One favored form; (2 input unc). T. original 3 input unc. had 4 inputs: (1) An input that modified the format of Umc, it continued to ~~repeat~~ continuously changing "A priori" in fact. (2) Q (3) random input.

03 Well! .02 - .04 is (1) my older 3 input unc (2) T. classic output D.f. of A "in view of Q": "Given Q, how many bits are needed to create A."

04 Did I ever realize that (1) & (2) of .05 were the same? — Pretty accurate exactly to start. The (1) input of .02 is not ~~really~~ ordinarily a univ. d.f., but $\approx O^j$ — [$O^j(A_i|Q_i)$] which need not be univ. (It may be well defined univ. — e.g. so every output has $pc > 0$.)

→ T. exact details of O^j will probably depend on how O^j was land! ←

Initially O^j will be for d functions — but assoc. w. each d - func, there may be > 1 way of implementing it, some ~~ways~~ ways to get known A_i 's, that extrapolate directly from the best pc O^j . — There will normally have much less pc . (Of course in "tree" s - prodn, there will have reasonably large pc 's.)

23 3.6.03 A poss. way to do this: Umc has 3 inputs: O^j, Q_i, R : outputs s - prod for A_i . we find O^j (2) & R to get $(O^j) + |R| = \text{min}$ for $Q_i \rightarrow A_i$. Then we try R w. same O^j , to try to get A_1 . Or we immediately switch to Q_2 w. same O^j , to try to get A_2 output.

29 What we want is $O^j \approx |O^j| + |R_1| + |R_2| = \text{min}$ ($R_{1,2}$: inputs for $A_{1,2}$ outputs) (I don't think there's any pt. in trying Q_1 & Q_2 in reverse order) Tom's goal is remember we remember codes $\approx (O^j) + |R_1|$.

30 We start by keeping track of O^j 's that give $Q_i \rightarrow A_i$: As a search for ~~the~~ R_i of larger & larger z , we had track to those earlier O^j 's. I think there is an efficient way to search the space for solns w. ~~small~~ \approx (larger & larger Q_i set. (8.25.03: Not obvious! 138.19 R 138.23.02 top at a soln.) So: we given $\approx [Q_i]$ set, we end up w. a set of O^j 's that do it w. "minimal" \approx best. T. pc of a A_i can be obtained by \approx using same O^j 's different R_i 's. Also different $O^j \approx$ various R_i 's. Research for the best codes for a given Q_i, A_i 38.00

v.g., imp vary idea!

optimize for small O^j to source $Q_i \rightarrow A_i$'s.

Nip

67 8 9 10
Th F Sat Sun Mon

Flight 653 PM
1859 PM

Cost 12

4403

935

File 239
Cost 12
16K <

00: 37.40! Can be done in a different manner, considering the 2^j 's of various O^j 's.
Essentially PC order for O^j , $[A_i]$ set.

03 So, this sounds like there is a rather simple search that is really Linear, that is about as good (flower cc) as can be done.

The mechanism of how I'm to implement this 3 input UMC is unclear. I did have another way to do it, but using a Lisp or ~~the~~ AZ-like machine, but I didn't like it. Would it be easier using Jvarkan's (OOPS) formalism?

02 **SN** Check up on OOPS' prefix mechanics: What's criterion for an output being a legit "trial"? Does it have to stop? - Does it have a stop state?
(His solves were always finite strings. Anyway, looking at OOPS system a search

I could use it for a 3 input UMC. \rightarrow Note \bullet 3TM (\equiv 2003TM) 318.32 \equiv stuff loading options suggests that OOPS system can be modified to 3TU

For a usual UIO UMC, how to search in $|O^j| + |A_1| + |A_2| \dots$ order is unclear.

all 3 inputs are prefix safe, \Rightarrow it's "usual" way to get this is to allow the system to ask for another token and libidum. \bullet If excluded, then backtracks to check branch unexplored. \bullet "usual way" is to allow the system to ask for another token and libidum. \bullet If excluded, then backtracks to check branch unexplored. \bullet "usual way" is to allow the system to ask for another token and libidum. \bullet If excluded, then backtracks to check branch unexplored.

A possibl. way, using a OOPS system: 1) The O^j pgm is a regular OOPS pgm.

2) Q_i is a read only input - it can be copied, but a copy can be modified.

3) R_i is perhaps an input like O^j where the machine can request new input by giving the address P_i^R , where P_i^R was the lowest address thus far for R_i .

How to implement the search routine of 37.29 - 38.03 w. such a UMC, is unclear. \bullet $pc(O^j) \cdot pc(A_1) \cdot pc(A_2) \dots$ \bullet $\sum 2^{-|R_i|}$ in $||$.

Actually, it may be rather easy, since its $|O^j| + |R_1| + |R_2| \dots$ has to be minimized, a prefix and all "prefix code" type inputs: Since we want Minz R₁, we don't have to be that order. (shortest, first)

It may be necessary to devise a special "Multiple tree search" w. 1 tree each for $O^j, R_1, R_2, R_3 \dots$. This routine feels like when to consider a "hyper" is a "next Q ". We can \uparrow $|O^j| + |R_1|$ either by $f(|O^j|)$ or, by considering short codes for R_2 as well. I'm not clear on the order in R_2 search.

BOOSTQ Just what features do we want/need that Boost suggests? Once a pgm has been successful, not only are tokens of this pgm more likely, but other "chunks" or inducers (w. perhaps gaps) \equiv other similarities, are more likely. These are "Mutations" or "Mutations" (\equiv Modifus) or t -successful PSMCS. If we use 2 pgms, we can cancel crossover (\equiv induction w. $SSZ=2$).

Re: Mutations: many mutations are induction using OSL (one case) - use subclones, etc. It's a type of Φ A induction! - So we want one or more O^j 's to do: O^j (mutation / mutated) spec

a P.P. on mutation, in view of mutated.

3TM 319.10
3T 323.21
for Bibl. Summary
Check more SERIOUS solutions to the problem

This is a good way of thinking about it.

41.00

0:39.02

Context of tokens: foundation w. scaling (of our kind).

Black/wh.

As I see it, "Context" is an enormously General Concept. A discrete Context (not Gray, stochastic) ~~is~~ Is (and hence, anyway) is a bunch of subsets of previous corpus. Each subset is a (discrete) context.

A circumtext is, for each token, a pd over ^{all} previous possl. corpi. — essentially it gives a pd. over all tokens (Epi:1) as a function of all previous corpi. This is sort of what Jo had in mind in OOPS — (but he really had no v.g. idea) → how to implement it,

Discuss this w. Hon

Other than "Boasty" is a few other probly Mod. insts.)

→ We could just do a QA induction to get a s-struct that gives the pc of any token, in view of the past corpus. — This is, of course, the "copy success soft post" update, is oblivious of the concept of "Optimization". Sorts a "Phase 1" (p.d.)

Anyway: I may want to try to do that induction "by hand" (i.e. M_E) — at least

in the early ~~TM~~ Top. of TM.

More Generally (if I think I may have had this in ~~my~~ mind some time ago — contributions

prob assoc. w. each past part-corpus is a pd. on possl. complete combinations of those

part-corpus. (I did write about this in recent post... where??) → Mainly early discussion of Difference bases

Essentially it is an ~~induction~~ pd on all poss complete corpi, it should be updated ^{by} pc's to every time f. corpus is augmented.
 My system is OOPS' device seems to assign tokens

→ Note that this is a "corpus of ^{successful} epus" — not [A]sets or [A]sets. narrowish

The large is a narrow with kind of "Context": A more General kind would ^{also} consider

to have a higher order Context of the problem itself (i.e. Q) is T. type of problem

to "indices" of f. Q (if any) — f. induced is indices, based on previous

Experiences.

30

So it looks like there are several places in TM ~~to~~ (or Phase 1) in which there would be a sub-branch on past corpi. We may want to do this "by hand" in ~~the~~ Young TM, but by TM in advanced TM (i.e. in Phase 2 we have an extra phase of "itrag").

If the IPC of the machine used is large enough, we can omit in view of the "By hand" induction is how to do it. So perhaps we should see how far we can go w. accessible IPC, it do hints a " ~~by~~ " Induction by hand" to bridge the "initial part"

→ (40.00) space

3.6.3 428 20:07 Back.

now from
from middle

40

1030 → 1
2220 → 24:30.

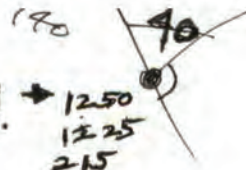
ca. MM map

12374 in Bus
at Priv.

825

11 P.

1250
1225
215



10:39:40: In 34.00 - 40 (Context of Tokens), I was getting a non-al. soln. of the Phase 1 problem - but what I was originally bringing in "Context" considerations was an gl. approach. - Pz may be ~~more~~ less "well-defined" than the non-al problem.

03 (33.03R): On the PC vs. kcost problem: In 28.20 ff (on corrected cards) it was felt Pz: Pz would not help because to write soln would get much less from assigned to it. Hvr, if all assigned cards have the same ratio of p to z^k , then it would be ok. i.e. a kind of "renormalization".

So 28.20ff may solve the z^k vs. pc problem. !

Worst shape
So, 1. problems in worst shape (ten 33.00 ff):

- ② Defining good contexts for scaling problem 39.00 is approach of probly linked
- ⑦ Analysis of Boost MOOPS 28.20 ff (first part 33.04 - 005) → 2 Mutzall
- ③ finding, expanding, fixing SAIRB TM: see 33.26

② Context: look at non-al soln, then look at various glens.

List many kinds of Contexts. Doing so should suggest many types

1) Most General non-al. "Context"! = PD on next token or entire completed corpus or part of unknown corpus - as a function of part of corpus, including Problem term (Q), Any indexing of problem &/o indexed indices.

2) Parts of D: PD' = an known corpus or Problem term or indices.

3) PD's on unknowns.

SN Methodology for "Context Discovery" (by Tromer/Mz). Look at situations in which I expect context to be important, & try to devise contexts that would help get hyperps for tokens (or larger chunks).

SN The Mechanics of L search for S functions of 37.01 ff, is a **GREAT JOY!**

1.87
1.75
1.53
12.469
cc
10.38
Post 2.02
Apr 10.5!

20 : 38.40 (on "Boost" as mutation).

Mutation is "size of" in GA, is somewhat different (But not much!):
 In GA, we might have as input: (cand, G) output D.f. on (cand, G). \leftarrow seems to be more useful.
 or we might just have cand as input, ~~and~~ d.f. on cand as output.
 In G.A. we can be (slowly) looking for a good (satisf.) O^i functions — P_{i+1} O^i will slowly change as t. population ^{it turns} changes — so our Mut. rates change in population ~~and~~ as it \uparrow in G .

Note that when I say "Boost" is like a "Mutation", I mean that to ~~mean~~ entire new cand pop, is a mutation of the "Frozen pop" being mutated.

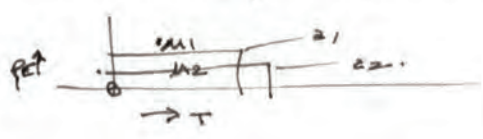
This Analysis of Mutation via the GA operator O^i , can probably be applied to the Problem of "Context" of 33.01 (for fighting the Scaling Catastrophe)

→ 42.00

SN on OOPS! T. only linkage betw. the 1^{st} problem & TOH was new OOPS!
 There was no ^{other} memory of frag. of use of tokens!
 Was there any (other) many of token fragus betw. 1st & 2nd trials?

SN In the ~~convergent~~ GPD, \rightarrow GPD₂ x $\frac{z}{M}$ (for mv. probs) is a prob. ordering.

What is rel. prob. of "best" betw. 2 cand: z_1, u_1, z_2, u_2 ?



if $u_1 < u_2$ then rel. prob. of cand₁ win is z_1
 prob of cand₂ " is $z_2(1-z_1)$
 Win prob. depends only on the orderings of the u_i & only values of z_i .

If the z_i 's are "small", the ordering is by z_i , indep of u 's: ~~eg~~ eg, z_1 vs. $z_2(1-z_1) = z_2 - z_1 z_2$ small.
 If z_i 's are large, ordering is by u_i and z indep of z 's. G.P. say all $z_i > p$.

I was using $\frac{z}{M}$ as rel order (largest first). This may not be so good!
 Perhaps use u_i order if z 's are large, use z_i order if z 's are small.

— But what if how to compare large z , ~~small~~ large u w. small z , small u ?
 So use z to compare, if z are small
 " u " " " " " " (over
 uses $\frac{z}{M}$ " " " " betw. mixed z 's.

What is threshold to change to $\frac{z}{M}$?

I don't feel that this analysis is so hot, but it does use a way of looking at the problem.

I. "all z_i are small" case is probably very common!

T. prob of cand being max = its z_i mult by $\prod (1-z_j)$ that depends on the ordering of the cand, z_i is a \downarrow funct of best ordering. The \prod factor is of most import for large z_j 's.



20: (4.12): T. main present problem seems to be "Scaling". General context is Boosty are 2 aspects of it: ways to deal w. it.

Can we regard the "boosted" set of tokens as a "context"? - Superficially, no -

these sets are common to all token preds i don't have t-specificity of "context".

"Boost" ~~subset~~ sets of tokens are common to all token predictions i are a property of t. "State" of t. token generating system - so in Real so use, they are a "Context".

07 But they don't seem un-logged for dealing w. "Scaling", since t. v. of Boosty. Boosty char. is t. corpus size, n. - i.e. T. prob hitting a particular set of tokens v. i boosty is $\propto \frac{1}{n}$ - will not be cost. - Randomly, as n \uparrow , t. no. of (sets of tokens) that would be helpful, also \uparrow . (as t. no. of related problems, solved in the past). Unfortly t. amt. of processes available is $\propto 2^{-k}$ not PC - which is sum of pc's of related (sets of tokens) $\sum_{i=1}^k \binom{n}{i} = 2^n - 1$ N.B.

12 Jürgen wants to use "context" (I guess) to deal w. scaling of boosty: as no. of perms in RAM (Frozen memory) \uparrow . $\sum_{i=1}^k \binom{n}{i} = 2^n - 1$

Looking at "Boosty" as a kind of Mutation is probably much better than regarding it as "Context"!

Anyway: Both Boost & Context can use QA induction: In Boost its $SSZ=1$ = Mutation. In Context its all previous (corpus \rightarrow token) \exists situations as $\{Q, A\}$ data set

So, boost actually does 2 things ① it \uparrow pc's of Tokens (fighting "Scaling") ② It \uparrow of a particular boost & w. corpus size, contributing to t. scaling costs to rise.

So, we would like to use (perhaps) use "Context" to narrow down which perms in Frozen memory to Boost. So, what we want is a "Context" that points at "related"

(Successful) perms of t. past.

Can I frame .22-.23 as a QA induction problem (as any kind of induction problem!)

I think I have to make t. Q's = problem defn, A's = perms that solved these problems. However, in A's I'll be only (for Jürgen's "Boost" detector) be interested in t. fragments of tokens which form.

So I want want a problem to maximize a density defn on Tokens

29: (12) Actually .07-12 may not have a scaling problem. As n \uparrow t. no. of successful perms w. token frags is useful for t. present problem, also \uparrow . We may have "Scaling problem", but of much smaller magnitude than t. normal scaling problem for pc's of Tokens. There is, however,

32 problem that t. no. of perms. (t. cardinality of t. perms. set) \uparrow w. n and is $\propto n^3$ so Boosty cost is $\propto n^3$. - So there is a serious scaling problem, but due to t. non-additivity of expanded cc.

A poss. way to deal w. problem of .32-.34! If a fraction c ($c < 0.001$) of the perms are repetitive: i.e. say each perm occurs $\frac{1}{c}$ times, then it would be well to eliminate perms that had the same set of tokens. We might have a bunch

0: 42.40 : of N clusters, 2 cluster centers.

On "cluster" centers: Each pt. is regarded as a Bag of tokens.

One way we can have a "cluster center", then each cluster differs from its "center".

I could represent each PDM by 2 pt. in R^p -dimensional space: r, s to card. of tokens. Each pt. has integral coords! Usually zero & occasionally 1; sometimes more. See NIP/ID 4.40/D.16.02 for station Clustering - Page 44.00

9N on Scandp PDM ID: ^{from} 836 backwards.

825: On time-varying O2 probs: I was unable to remember soln. at 1st pt. ^{See} ETP: 766.10

816 on Scaling.

811 Goal Priority of prob solving.

811.30 on TSO construction ← Goal.

807.20: Part of proof of corr. PDM for QATM

806: N Dm. ALP

802.02: Human backtracking/Praty revision → 800.03 → 796.20

.21: "Phases I, II, III" of ~~TM~~ ← this probly ≠ praty idea of phases 1, 2, 3.

788.12: SUMACS - 768.20: Sequential coding for QA in Timeseries.

777.00-18 No. param > no. data pts. 778.00 "Curve fitting"

771.07-08 Bigunsplit problem at Mat time, NIPS 12.13 may be relevant - ± Praty & solved it.

765: HMC how many Corts?

764: ANN & RANN - How to get ALP values

762.00 Multimed AAE (Active Algor. Evoln), etc. : Praty very because of some es problem of Difrent "Loss functions" ← superficially, NO! It ALP comes to write P.D. it should be used to optimally control "Loss function".

761.17: Developmental G-PD - PST Grammar - from early beginning. Review

753.11 SM Strategy as Goals 01/02

See of NIPS 660 backwards ^{2/18/03} Next

(A1) Things to put in Revision of Report: (also see 3.01)

40: on Noraz Univ D.F. : Just what did Levin show?

39: MEM : error in my ideas

30: ~~to~~ Q. mostly list.

18.25 GA.

10.25 SOY, STEIN ~~to~~ AAE, etc

4 Clustering: Purports to be genl soln. for continuous N dim vector data. → 44.00 ^{with}

3.01 see (A1) (on next version of report).

CLUSTERING

ID
NIPS

This was in response to a proof (maybe at NIPS 2002) that clustering was impossible to do in an exactly rigorous way.

General clustering: Given a data set, a set of points, find several parameters. \vec{x}_i is m dimensional

A "cluster" is usually a way to ~~pick~~ a) pick a point in the space, b) find a ~~center~~

Scalar function $F(\vec{r}_j, \vec{r}_i)$ (\vec{r}_j is a data pt. \vec{r}_i is "center of j th cluster")

$F(\vec{r}_j, \vec{r}_i) = \max$ overall second angle values. Perhaps $\sum_i F(\vec{r}_j, \vec{r}_i) = 1$ & $F(\cdot) \geq 0$
So $F(\vec{r}_j, \vec{r}_i)$ is a ^{kind} probability that \vec{r}_j is in cluster j .

Anyway: The idea is that this is a way to decribe a set of points.

How to describe

~~How to describe~~: If every data pt is assoc. w. any j & all data

points ~~each pt.~~ If there are j clusters & n pts., each pt. can be

described j ways. We want F 's such that total info in the desc of F

functions F plus F j \vec{r}_i . ~~plus~~ plus the info in the desc of F .

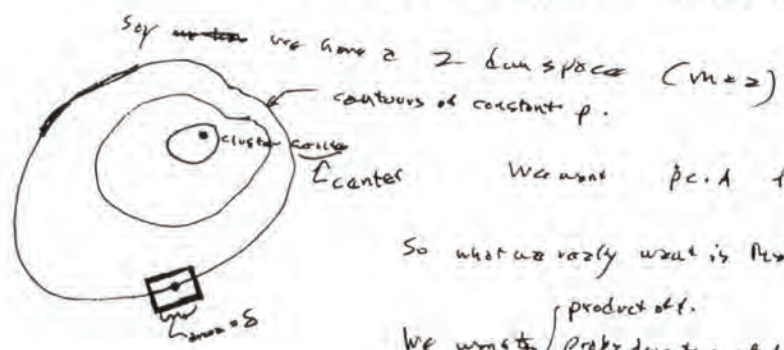
The pc of each data pt is some of its pt's w/rt. each of the j clusters:

We want product of ~~pc~~ pc of each pt times (pc of j centers) \times pc of $F(\cdot, \cdot)$ to be max.

So can compare it w any other clustering of these pts.

The entropy is a way to describe set. ? — Yes! .22 - .40 shows how.

.21 is uniform! In a cluster, $\sum p_i = 1$, but how does this enable us to describe each of the pts. w/rt. their center of the cluster?



What we want is the prob of describing a data pt. So that it is within an area, A ,

We want $pc \cdot A$ for small A . $[(pc \cdot A) \text{ will be } \propto A \text{ for small } A.]$
$$\int_A F(\vec{r}_j, \vec{r}_i) dV = \int$$

 V is a volume of element in (\vec{r}_i) space

So what we really want is that product of prob densities of the data pts to be max, when multiplied by pc of the $F(\vec{r}_j, \cdot)$ descs.

The actual ^{pc of A} desc of the data set to be by accuracy can be obtained by multiplying above pc by δ^n (for n data pts & we want to know each pt. to within area of δ).
but the δ^n factor is indep of the clustering method! So ~~we~~ choice of clustering method does not depend on δ .

KTK
K.10.10.10

≡ 143 1/2

Clustering

so: (NIP/ED 4.90
D.16.02) → An alternative view: Not as a k-means partition:

We have as before, a function $F(\vec{r}, \vec{r}) \geq 0$ ~~→~~ $F(\vec{r}_i, \vec{r}_i) = 0$,
← center of cluster

$$\int \dots \int F(\vec{r}_i, \vec{r}_i) d(\vec{r} \text{ volume}) = 1$$

find.

We have a data set. We hypothesize that these data were generated by k (cluster) centers.
In accord w. ~~to~~ $F(\cdot)$ d.f. Each center C_i has pc of p_i of being chosen

SN Min "Clustering" problem: It's ~~undetermined~~ BAG induction: "Two kinds of... induction" paper

T. General idea: We have a BAG (Continuous or Discrete). We express it.

BAG by a set of k "centers" & a "distance function" from k centers. (This distance function integrates to 1 for each center: is, normal)
The dem. of BAG: each element has k ("codes") pc's add. each code gives its distance pc.
from each center. To bag a center & its words.

We want a set of k centers & a "distance" function, so that k pc of dem is k Max.

Each center requires a pc - total of them pc's is 1: we assign these pc's to Max. pc of BAG.

(Continuous or discrete distance functions can be regarded as "Mutations" of each center

This suggests many poss. forms of $F(\cdot)$ function. ~~→~~ COPS' best "mutation" -

Uses pc's of ~~of~~ focus of center plus (low) fixed pc's for other tokens] Other mutation by pos

used in GA & GP] probable ("") S-grammar.

is assigned a pc: $\sum pc = 1$ ~~with~~ - ideally, each center grammar should be able to express any string, but most w. very low pc]

SN Also Note 36.20, on an Advantage of ALP over "those best GAs".

10:35:40 : On Computability v.s. {Universality as a predictor}

01 \square any CPM (Computationally PC Measure): One can design a definite seq. in finite time,

02 \square PC \leq CPM gives PC of .5 or less to each bit - some error is always $\geq .5$.

How probable are these "irreducible seqs"? : There is only 1 seq. in finite CPM desc, that will be as bad as $10(-.02)$. However, if we allow ϵ CPM to have \rightarrow PC errors $\rightarrow \epsilon$ ($\epsilon < .5$), for each ϵ , there is a larger set of seqs. i. no. of seqs $\rightarrow \infty$.

It is probly poss. to measure in some sense, the no. of seqs for error w . PC error $\geq \epsilon$.

0 If we have a CPM w. observed \square ms error of ϵ for a given sequence, can we say anything about the expected future error, in terms of PC of desc of seqs that have future (ms) errors of ϵ .

47.00

19 on a simple proof of the unsolvability of the Halting problem:

A version that I (Sincerely so) understand:

Say \exists PGM that can determine halting of a finite string.

- 22 1) examine all PGMs < 1000 bytes long.
- 23 2) run all of them that halt.

~~22~~ 22 - 23 is < 1000 bytes.

If it halts, running it \rightarrow infinite loop of calling itself.
 " " doesn't halt, PGM running to rest of PGMs must halt.

This is related to: "find halting PGM of length $N < 1000$ by w. longest run time" (any PGM runs to ∞ which is ∞ loop, not with or that does longest number)

So the no. of PGM seems to be:

- 29 1) If this (29-30) PGM halts, runs, then points to ∞ loop.
- 30 ~~29~~ " " (#) doesn't halt, then stop.

Which is the PGM I couldn't understand! i.e. it was unclear as to what input to the \exists PGM that determined if a string \rightarrow stop or not?

33 So PGM: $\alpha(x) = 0$ means $x \rightarrow$ stop $\alpha(x) = 1$ means ∞ loop.
 34 1) If $\alpha(.33) = 0$; then $\rightarrow \infty$ loop
 2) " " = 1 : then stop

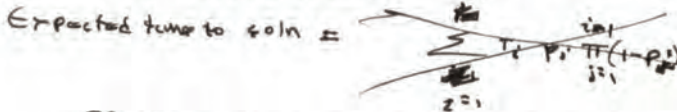
Now: What did I go wrong? It wasn't clear to me as to what input to α was to be.

I had not conceived of the recursive PGM, .33 - 34!

So, essentially 3 ways: (.23 - .23)R, (.25R), (.33 - .34) all about the same.

$\beta = 15$
 $2:30 = 14:30$

T. first G-H turn:



Expected time to solve \approx
If we change ordering of z adjacent z 's in $\Sigma \rightarrow$

$$\sum_{i=1}^N \left(\sum_{k=1}^i p_k \right) p_i = \sum_{j=1}^{N-1} \left(\sum_{i=j+1}^N p_i \right) p_j = \text{expected time}$$

We want to compare expressions: Only $z \in \Sigma$ terms of r result in expression change.

It would be good to include this in a report as an appendix, - R is it is an essentially trivial (but messy) to prove. Much bookkeeping needed!

A point about the updating phase z update, involving h is h . $h'(t)$ is the prob. density of finding z at time t . This automatically says that no solns were found at $t < t_0$

Prob. of failure to time t_0 is? $\prod_{i=1}^n (1 - \frac{1}{n} p(\frac{z_i}{n} t_0)) \approx \prod_{i=1}^n (1 - p(z_i t_0)) = \prod_{i=1}^n e^{-p(z_i t_0)} = e^{-\sum p(z_i t_0)}$

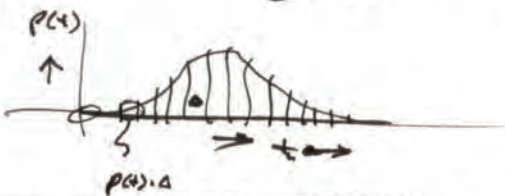
This seems wrong!

T. Prob. of failure at $t=0$ should be $1 - \int_0^\infty p(t) dt$ $\rightarrow e^{-\int_0^\infty p(t) dt}$

These expressions are same if $\int_0^\infty p(t) dt < 1$.

In Question is my idea of what $p(t)$ means.

- related to z Q's
- ① what's prob. of soln at time t_0 ?
- ② " " " no " up to time t_0 ?



pos. conclusion: It is poss. to use $h'(t) \approx h(t) (\text{+ const})$; $h'(t) \approx h(t) - h(0)$

or $-r(t) e^{-\int_0^t r(t) dt}$; $e^{-\int_0^t r(t) dt}$; In either case, the first r found is

to divide of t second - I may need "constants of integration" in both, here.

I think the integral starts at $t=0$ & goes to $t=\infty$,

so $h'(t)$; $\int_0^\infty h'(t) = h(\infty) - h(0)$.

$\int_0^\infty r(t) e^{-\int_0^t r(t) dt} = \left[-e^{-\int_0^t r(t) dt} \right]_0^\infty = 1 - e^{-\int_0^\infty r(t) dt}$

backwards! \rightarrow no!

correct! $e^{-\int_0^\infty r(t) dt} - 1$

derivative = $(\int_0^\infty r)$

note $e^{-\int_0^\infty r}$

sign $-$

$h(t) \approx h(\infty) - h(0)$ would be correct.

$R'(t) \approx -R(t)$ $1 - e^{-R(t)}$

35

$R(0) = 0$

$r = R'$

related \rightarrow

SM

0:00 Variation on Cover's method: Have a grid on strategies (linear, non-linear...)
 → wt. on each or a grid. T. a grid is PC that each strat is Best?
 Update by each strat having 2 mean → var to yield each year. Pick update
 by prob that each is Best (≡ max yield). Maybe not want to update —
 Monte Carlo is 1 way. Perhaps have 2 or 3 outputs for each strat, & do an
 exhaustive search. For 3 outputs & n strats, this is 3^n cases. — so
 $3^{20} \approx 10^9$ so $n=9$ for 18 ~~strats~~ strats gives 10^9 trials (a lot)

(This is a QATM Phase 2 update problem!)

→ A really impl. How fast is to do a phase 2 update: A more specific idea is what
 we really want updated! — what is Real GRC is.

→ So: IS is OK to make the PC wts. correspond to the PC that a given strat. is the "Best" in the
 very considered?

→ The Phy Ret behaves much more like an updatable pd. is the curve giving PC of each
 poss. yield. — the pd. of yield. The "Phase 2 update" is a formal mathematical return
 of the yield space.

A By Problem in Cover's system with Q of just what + initial wts represented is
 how one could use previous & perhaps || info, to sharpen them. up.

Es he presents to system, to initz (wts are uniform). They might represent expected (by world "per yr." say?)

or exponential $\exp(\text{expected yield/yr})$ & 1 initially since exponent = 0.

We write initials to wts. at $t / \left(\sum \exp \right)$ (sum of ln year yields of individual stocks).

45.19
0: ~~10~~ : Re: Computability, accuracy etc: for revision of kolTalk-comp Journal.

APL uses approx to get ~~max~~ modes: Runs is no "training set". All data is used for "test set" — I think the resultant error is unbiased; but, it's an awfully complex distribution function. (?)

In APL, the error is to an extent not "known", but with usual sources. We obtain error dist. from entire data set. This error should be minor for \downarrow as $CB \uparrow$ is \downarrow PC of data p. But we don't really know how big it is! — in what sense is the error unbiased?

In the kolTalk, I said that we never have a usable upper bound on error. Well, the mean error in the data set, for APL gives a kind of upper bound on error. The expected error should \downarrow as \downarrow p.

On the other hand, for any CPM, (and any APL, w. finite CB is a CPM), we know there are sequences for which that CPM will do catastrophically poorly!

- 1) So, there is the approach of .00-.04: related to Conv. Num.
- 2) Runs is ~~idea~~ of .13-.14, on uncertainty of APL.

on this subject.
131.28... to
135.00... to
147.00...

A reply of .00-.03 to .13-.14: The untractable seqs are unlikely. The most extreme intractable seq. Runs is only one. Seqs that are untractable to extent of error of ϵ in PC's become more numerous as $\epsilon \downarrow$, from \uparrow number $\epsilon \uparrow$ further!

One idea of error in approximating univ. dist. That as $CB \uparrow$, error \downarrow — but at no point do you know how large it is. I think as $CB \uparrow$ the PC's become closer to 0 & 1.

Also APL can often tell when one approx. is better than another (less expected abs error), one can't get size of error, but one can know how to get error & know which ways are better than others.

Apr. .00-.04 does give a real kind of error measurement that seems (?) to be unbiased. (Meaning of unbiased is unclear here: perhaps wrt. true generating probab measure?)

Also, the Conv. Num applies to $CB \downarrow$; What if $CB < 1$? Then even if \downarrow is a CPM, we don't know that $\frac{P_n}{n} > P_n(n)$ (i.e. $P_n > P_n(n)M$).

The main idea I guess is that shorter codes mean our must be getting closer to the univ. dist. & the more likely that Conv. Num. is relevant.

There is a long seq. in Prose notes about just why shorter codes are better. — it may be that .34 was the "minimal seq".

Consider the Corollary of Conv. Num.

Spec:

147.40: Well, ~~at first glance~~ at first glance, t -correlary is not relevant. By finding shorter codes for t -corpus, implies ~~we have an approx to~~ we have an approx to t Unl. d.f. But ~~might~~ ^{Worse} have a pc for all seqs that is a "constant" times worse than M but t -constant may be $< P_M(\mu)$.

Note In trying to approximate t Unl. D.f., we know that as $C \rightarrow \infty$ we will $\rightarrow t$. d.f., But this is still not t -same as knowing M . With P_M we get t -"best best" but we still don't know

how good it is! — Yes no do! We get, in all cases a Unbiased estimate, since our model is a priori \geq entire corpus \geq test set.

Also as $C \rightarrow \infty$, our apparent error \downarrow so our real expected error will \downarrow !

\rightarrow look/adjust to justify of ALP.

A t -Cons. P. suggests that usually it will be a sig. (perhaps t base poss., fall into it in approx).

Opposing ideas of knowledge of error in t -error: There is idea of how near knowing how large we err. from Unl. D.f.

In measuring error, perhaps best use KL distance, since it is concerned w. ~~error~~

rates — so when we get small error, it's still relevant, \geq can show that we still can have no distance to go!

General conclusion: That usually, incomputability of U.P.D is not relevant.

By using a priori predn. one cannot get good a priori estimate!

One " " " " " " far larger than this way, also.

ALP, hvr. (1) Suggests other models: P.r. models.

(2) Can do OSL (which MDL \geq MML is most often used & not deal w.)

(3) Tells us that we can usually do better. (But is never sure of this).

"Better" can be deterministic?

This is new from ALP fact Prob. Randomness is never confirm.

(4) If we have 2 models \geq one has shorter dev. time t & P !

Why "is" shorter one better? ALP says use with-mem. w/ly as length of code.

Empirically, we will find usually find that short code methods gives smaller error.

(The not always: — it may have a very short coded t can a' more "complex" error)

Well b. shorter dev. is closer to ALP (approx as monotonic \uparrow).

\geq ALP is v.g. — perhaps is good enough poss! !

IS \geq \geq \geq as I can say about it? It's not bad — but...

What to do if we have 2 difent a priori (2 difent real machines) \geq wait!

Well, this is equiv. to difent "pre corpora" \geq t for difent time series, one has difent reqs: The disagreement can be resolved.

Unless, in some way add the 2 "pre corpora" together! A Resolution!

So I'm still unclear on how to solve this one is better than other.

10

13

10:149.40 : Be sure to mention in (list of desirable): Optimum tradeoff betw. Complexity of model & apparent fit of model to data.

Under these conditions, all of the data ^{can be used} for testing and ^{Error Estimator} ~~estimation~~ will be unbiased. There is an optimum tradeoff betw. complexity of model and ~~apparent fit of model to data~~ ^{apparent fit of model to data} ~~no underfitting or overfitting, but~~

perhaps no mention of ~~SSZ~~ ^{unc. v.s. model uncertainty?} ~~uncertainty~~ ^{usually}

In practical applications of the UFD, incompatibility is not relevant. we use approximation to the UFD that are computable and we obtain unbiased error estimates ^{usually} ~~usually~~ ^{we use approximation for them.} These estimated errors are usually smaller than those obtained using ~~other~~ ^{statistical methods} since we are able to use all of the data for testing.

soften this
a Bit - it's not sure it always works this way

The intractability of UFD implies that there are usually ~~many~~ ^{many} models of the data that are better than those we have used in our approximation to the UFD, ^(but) ~~and~~ that we can never know how much better these models are. For example, if we were given a pseudo random sequence of bits, and we didn't know it was pseudo random, we would probably conclude that the probability of 0's and 1's was about .5. If we were able to discover the pseudorandom rule, this more accurate model would give us probabilities very close to 1 for each prediction.

In ~~all~~ ^{empirical} prediction (not only approximations to UFD), we are usually ~~in~~ ⁱⁿ ~~uncertain~~ ^{uncertain} that ~~the~~ ^{the} ~~best~~ ^{best} model we found thus far, is particularly good. It is ^{almost} ~~always~~ possible that if we searched for 10 minutes longer we would find a much better model ^(though) ~~than~~ ^{than} the pseudo random case ^{above} just mentioned is an extreme case.

The uncertainty of ~~the~~ ^{the} ~~best~~ ^{best} model ~~is~~ ^{is} ~~characteristic~~ ^{characteristic} of ~~all~~ ^{all} prediction ~~problems~~ ^{problems} for empirical data. It is a characteristic of probability itself. ^{holds true for} ~~all~~ ^{all} prediction ~~problems~~ ^{problems} ~~(and the uncertainty of how much better models exist.)~~

SN Side Q: In attempting to calculate UFD: (if one is interested in a sub cover of E, can one only consider a finite no. of codes? Say we have found a code for a data that is of length L; Or: we find a model of length L that gives pc of 2^{-L} to E versus 1 to total code length L + r. Do we only have to consider codes of length $L \leq L + r$?

Many very promising Models will take very long times to evaluate. Others will be ~~partial~~ ^{partial} recursive, so that in a finite time, one can't be sure one has covered them ~~completely~~ ^{completely}.

In the Rolf paper, "Error" is used in 3 ways: (1) SSZ error (2) expected error for a particular CPM, corpus pair (3) Deviation of CPM used from UFD ("best" model).

Models for which we have only spent enough time to evaluate partially,

... but for Model uncertainty... give poor predictions. At any point in our approximation of UPD, we will be using a weighted sum of the predictions of the best models we've found thus far. We ~~can~~ ^{2nd} know the expected error in P_{est} .

will have a good, unbiased estimate of error in prediction, for that approximation. At ~~any~~ ^{first} point, however, P_{est} will be models we have not ^{yet} considered and models ~~that we have not spent enough time to~~ ^{for which we have} ~~partially~~ ^{fully} evaluated.

These unsearched-for models give rise to a difference between P_{est} approximation and the true value of UPD. The incompleteness of UPD means that while this difference must approach zero, as we spend more and more time

evaluating models, we ~~cannot~~ ^{know} know at any time how far our approximation deviates from the ideal UPD. An example: Suppose we were given a pseudo random seq. This difference can sometimes be very large - as in the following example:

Since we always have a good estimate of prediction error for each approximation, this incompleteness does not in any way inhibit our use of the best approximation for prediction. It does mean however that no matter how long we've spent evaluating models, there is always the very real possibility that if we spent 10 more minutes in search, then we would find a much better model. This is true of all methods of estimating probability - not just the UPD.

(SN) It's clear that there are 2 meanings of "error" here. One is error in prediction. One is deviation from true probability.

Perhaps draw 3 models of data ① M best generated data ② P_M to UPD. ③ an approx to UPD ④ Any CPM trying to approximate M .

at 4.21 of Kollet's types There are several kinds of errors involved in statistical inquiry. (Some not out of time was spent for following)

To see how they are related, let us consider a stochastic algorithm M , that is able to generate stochastic sequences. M assigns a probability to every possible sequence x - $P_M(x)$ is the probability that M will generate x .

Let $P_M(x)$ represent the probability assigned to string x by the UPD. Let $P_M^+(x)$ be an approximation to $P_M(x)$ obtained after search time t .

~~Let~~ $P_M(x) \approx P_M^+(x)$, the approximation obtained after an infinitely long search.

The convergence theorem tells us that as x becomes very long, the difference between $P_M(x)$ and $P_M^+(x)$ becomes, on the average, very small. Unless we know $M(x)$ (and we rarely do) we cannot know how large the deviation is.

0: $P_M^T(x)$ approaches $P_M(x)$ as T increases. We know that $P_M^T(x)$ approaches $P_M(x)$ arbitrarily closely as $T \rightarrow \infty$, we can never know how close $P_M^T(x)$ is to $P_M(x)$. It is in this sense that $P_M(x)$ is "incomputable".

While the exact value of π , as a decimal fraction, is not computable, it is regarded as a "computable number". This is because as we make approximations to π of greater and greater accuracy, we know how much our approximation deviates from the true value of π . It is in this sense that π 's "computability" differs from P_M 's "incomputability".

ABCDEFG

All methods of computing probability from empirical data are similar to approximations to P_M , in that we can never know how far they are from the best possible approximation. The "incomputability" of P_M is a reflection of a property of probability itself — our inability to know its value to any useful degree of precision.

The way to do this is just what kind of errors was normally measured in a CPM.

For digital sigs, however, when a digit occurs, how close is its pc assignment to the one? For prediction of continuous quantities, what are the probability densities assigned to a element of the data sequence.

One thing UPD gives is mean (in pc density of elements (continuous or discrete) that actually occurred. This is like coding length of a data seq. Is it an "unbiased estimator" of code length (or mean code length)? The log of each element of a corpus... Presumably some distribution, probably Gaussian, so we get an unbiased estimate of mean pc/element (or pc/symbol).

So (mean) pc/element is usually unbiased estimator — we can, here, know its var. (or var of $\ln P$).
 [unclear] $\frac{1}{N-1}$ or $\frac{N}{N-1}$ — probab $\frac{N+1}{N-1}$ (0.35)

3.18.03

Factor points: for continuous predn! say we want to code $a \pm \epsilon$ to T . Will be like $\alpha + \frac{\epsilon}{\Delta}$. α is "fixed code" pc approx $\frac{\epsilon}{\Delta}$? α is α .
 But suppose we are comparing 2 alternative codes for $a \pm \epsilon$. For discrete coding, normally we have a "trade-off" between cost of model & cost of corpus w.r.t. model. Thus we seem to have a "variable" cost corpus in terms of model.

Maybe "no problem": If we have 2 competing models (mod_{1,2}):
 The cost of code of $a \pm \epsilon$ will be $\alpha_{1,2} + \frac{\beta_{1,2}}{\Delta}$? or $\alpha_{1,2} + \beta_{1,2} \ln \frac{1}{\Delta}$
 or $\alpha_{1,2} + \gamma \ln \frac{1}{\Delta}$ (I think the part dealing w. Δ getting very small should be common to 2 models)
 So maybe $\alpha_{1,2} = \log_2 A$. (= $\frac{\log_2 A}{\text{Model}}$ $\frac{\beta_{1,2}}{\Delta}$) Model.

So perhaps "no problem"!

35: (23) So if bits/symbol is a measure of precision of a prediction: This is expected to decrease as ϵ approaches a limit (for stationary data). Here, I think what UPD gives is perhaps an unbiased estimator for $\ln p_i$ of next symbol. ($\ln 1/2$)

At least of my understanding of this would be to see how it fits linear regression.

UCB

nice?

0:152.40: Actually pc of next symbol: or $E(\ln P_{n+1})$ would be better - More General.

Consider linear regression: After we get a set of models, we get a PD on poss. values of t next element.

Actually, the pc of t next element is always about the same. T. interesting thing is Var!

Of poss interest: Varc using "history" v.s. varc. using UPD or U PD approx.

Say we have a bunch of x_i w. a Gauss. D.R. about 0, & $\text{Var} = \sigma^2$.

D.R. = $\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{x_i^2}{2\sigma^2}}$ to desc. n of t pts: $pc = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{\sum x_i^2}{2\sigma^2}} = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{n}{2}}$

= $\left(\frac{1}{\sigma\sqrt{2\pi e}}\right)^n$. If our width is Δ , we mult this by Δ to get $\left(\frac{\Delta}{\sigma\sqrt{2\pi e}}\right)^n$.

So we are interested in t . expected value of $\ln P_{n+1}$. (P_{n+1} = pc of next symbol).

$\ln\left(\frac{\Delta}{\sigma\sqrt{2\pi e}}\right)$: $-1/n$ is the main interest: Not exactly! \uparrow is for next symbol.

\uparrow is for mean previous elements. 54 mbols.

of μ , t . creator of data

A simple case: we have n , x_i values that are normally distributed about zero, but we don't know t . mean is zero. We study n values, decide to Varc $\frac{\sigma^2 n}{n-1}$: T. expected error for x_{n+1} is

$\sigma^2 \cdot \frac{n}{n-1} + \text{noise}$ $\frac{\sigma^2}{n}$, because we think the mean is not zero but $\approx \frac{\sigma}{\sqrt{n}}$. \rightarrow Unusually, the variance due to error is linear $\approx \frac{1}{n-1}$ instead of $\frac{1}{n}$.

for n Gauss var mean 0, varc σ^2 , t . expected value of $\sum x_i^2$ is $n\sigma^2$.

$\frac{1}{1 - \frac{1}{n}} = \frac{n}{n-1}$
 $\frac{1}{1 - \frac{1}{n}} = 1.06$
 $\frac{1}{1 - \frac{1}{n}} = 1.517$
 $\frac{1}{1 - \frac{1}{n}} = 3.846$
2.921 mm
115 389
down 95"
1.43 mm

156.33

So: Discuss 2 kinds of communication error types! Mean $\ln P_{n+1}$; σ_{n+1}^2 : variance of loss.

FLP theory says that whatever approach method you use, you could probably do better.

T. Incompleteness of UPD tells you that there is no way to know how much better.

The σ_{n+1}^2 error: T. prediction method gives a d.f. for x_{n+1} , if \hat{x}_{n+1} is the mean of this d.f. Best $\frac{\sigma^2}{n-1}$ variance value of $(x_{n+1} - \hat{x}_{n+1})^2$ is a common measure of error in prediction.

For a discrete alphabet, one error criterion is the mean value of $-\log P_{n+1}$ assigned to x_{n+1} by the prediction method.

For both of these error criteria, if we use a purely a priori UPD, we can obtain an unbiased error using all of the data in our sample. In this respect, UPD gives results comparable with other a priori prediction methods - except

that they tend to be better than UPD because the entire sample is used for testing the models, whereas a priori -

No data is wasted for training. In most statistical analysis

only a few recursive functions are used for models. UPD uses the full range of recursive and/or partial recursive functions.

By using a weighted sum of all of the prediction models tested, we obtain a somewhat better error prediction error than that obtained with any single model. using the best single model found thus far.

13

19

20

30

NIPS

20 : \boxed{SN} On unc output: If its a U.I.O. wire, and "stop" state is used to indicate end of an output, then the system has essentially \exists symbols, 0, 1, stop.

So why is more complex "prefix code" formalism?

03 \boxed{SN} For any CPM, we can easily construct a ^{deterministic} seq. for which each bit CPM always assigns p_i E.S. to each of its bits (an "untractable seq"). Hvr. This is not to kind of error that Conv. Perm. limits. It is, t. kind of error ~~we can get an~~ "untract error for our estimator" ~~this is~~ con. d. m. p. - To map \exists to a potential "reader". To get an untractable seq. that becomes untractable after n bits -

How long is its down? Say we have CPM approx to U.P.D. Its down length is $\approx \log_2 T$ bits. T being CB . So, t. untractable seq. binary has n bits + $\frac{1}{2} \log_2 T$ bits. $\approx \log_2 T$ where n is t. entropy/bit of U , T source. T then $\log_2 T$ ($\approx \log_2 T$ is bits to desc T)

is needed to tell when untractability begins. Then t. seq. can be untractable w.r. n , $bc(n) + bc(T) + o(n)$. So total down = $n \cdot h + \frac{1}{2} \log_2 T + o(n)$.

4 If we just use the CPM to produce first n bits, then use t. "opposite of CPM" to predict (after n), then bc of seq. is $bc(T) + bc(n) + o(n)$.

$bc(T) \rightarrow \infty$ as $T \rightarrow \infty$, but is a "small" ∞ . (argbly "small")

16) Hvr! $\cdot 14$ bitcost is $\approx n$, so for long sequences, it can ~~be~~ have $\ll bc$ from $h(n)$, t. bc of t. down of t. "most likely" corpus.

18 -16-17 are not exactly relevant to anything!

I probably should figure out just what "untractability" is "untractability" mean in view of .03-18!

20 Anyway, UPD does give untract estimate w. no "waste of training set" data. Their cost in computing can be larger, hvr., since we do have to examine a lot of candidate cases. approx. We can, hvr, limit this \approx ~~for~~ for h pps as Reasonable does, by a priori choosing a (large) class of CPM's to consider.

26 I guess .21-.23 is t. main "weakness" of using UPD, rather than stochastic prob.

27 \rightarrow (or perhaps ~~that~~ stoc "augmented by OS Seq.")

The Listing of Pams in pc order may not be so hard: Much of my Prob has been on continuous prob, hvr. Any method that does prediction then has \geq MS error (or some other error distribn) can be used. Trouble is, we end up doing Perm in "not pc order", i.e. we do a set of $\#$ probn methods \exists we get ms error for each. We can get t. smallest error size of t. smallest error m.t. set, but this looks like we are not doing trials in "pc order" at all!

34 If a "PAM" for real induction involves getting t. peak pc, we can get Perm in pc order - t. cc of t. cond. involves time needed to find peak.

On 5 Nov, I did find a paper of a guy who did pc Perm. I think it took a long time to get. He used various post opten techniques. - But even so, we really don't know t. pc of cooper w.r. pam until its done!

T. Idea of trying a Perm on t. corpus (Perm's in pc order) Perm obtaining pc of corpus w.r. t. pam doesn't give pc order for cooper. (38.00)

3.21.03
NIPS

$2 \frac{1}{2}$ cups H₂O to 1 cup brown Rice.
Boil then → Micro at "2" for 45 min.

156
155 1/2

Things to Do:

- 1) Write Ivan C reply
- 2) Finish review of Kol Lecture: I don't need this for publication, but I do need it for my own remembrance & understanding. (Also for FSAIR talk)
- 3) Write Jurg. about differences betw. OOPS & α :
 - a) OOPS has less "memory" betw. ^{types} nodes: Only @ Boost.
 - b) How α has more ~~memory~~, memory - i.e. direct memory of previously solved probs. in format of re-iterations.
 - c) How Boost is a kind of "Mutation": Distinction - not - a problem type.
- 4) Make more "final" corrections to IDSA "Report".

Re: Wallace & MML: T. correct way of thinking is that χ^2 data induces a P.D. on \mathcal{H} models. By picking the "Best" (by χ^2) model, we "discard into" \mathcal{H} a set of χ^2 discarded datapoints on how sharp the peak of the P.D. is.

I think this "picking the best" is pretending it is "true" is a common error in scientific Reasoning.....

- 2) Main uses of \mathcal{H} discarded into:
- 1) Better proxy values
 - 2) Proxy ^{revision} Reasoning.

NIPS ● ●

1991 3.22.03
F

0:154.40: So: Problems w. S-functions & Levenberg Thru:

I certainly know a lot of reasonable PEM's & for roots - But they do seem to suffer from

diffy of 155.27-.40: It appears that diffy involves ideas that knowledge of the PC of it

code is "very sudden" - it's not a matter of of decreasing by a factor of 2 after we code a bit

Well, Suppose we just try the ~~code~~ ^{code} in param order. We can regard the problem as a variety of INV problem, in which a final output is either better or worse than previous peak.

What would a search for a v.g. PEM look like?

T. Search is in $\approx \frac{PC(\text{param})}{CC(\text{of param temporary on corpus})}$ order.

It certainly looks like a search

In doing S-functions on strings rather than v.e.s, I used that 34 bit v.m.c.: 137.00-138.29

In this case, the trials were in really $\frac{PC(\text{computational of code})}{CC(\text{total of code})}$. I could do it such a way, or at all, which is perhaps what I had in mind before I discovered (137.00-138.29)

Actually, .04-.09 doesn't directly apply to QA's on keys. - tho I see how it could!

Say $Q \rightarrow A$ are ~~vectors~~ vectors. We try various (parameterized) ~~S~~ S directs from Q to A . We tend to param in which the "MSE" error is min (would be 1 way)

Examples of QA's:

Q	A
dig. Name of person	photo of person analog
photo	name of person digital
real vector	real vector

Industrial process!

dig. } Who is President Bush? dig. }
 IS 124 integer? Yes. }

(SN) How to get arrays of costs & all other things! see (33) for v.g. way to do linear costs.

E.g. for (linear prodn. costs) Assume a parameterized form of spread function. Then, over a large set of problems, optimized to params. This would beaten? (like a vicious \odot), since we then need to know arrays of "parameters"; However, we can use O.f. (old faithful) to get v. values; They need not be very exact. Also, ~~more~~ more params can be used to give better param values (recursion - successive approx).

33:153.19 Actually, situation is far better! The analysis of var of 153.13-.19 words in

n dimensional linear regression. We can get unbiased var estimates caused by 2

effects ① var of parameter $\hat{\beta}_m^2 \approx \sigma^2_{\text{observed}} \times \frac{n}{n-m}$;

② Since costs have error, var of number is $\approx \frac{\sigma^2_{\text{m}} n}{n-m} \approx \text{costs} \times \frac{1}{n-m}$.

Therefore by picking most params \rightarrow this is minimal we can know no. of params w.o. knowing arrays of params!


NIPS

Bette James & Associates


Business Support Services - Since 1977

00: 156.40

This generalizes to some extent to Non-linear, if it's "locally linear":
Diff cases of "Not locally linear"

- 05 1) We are at a local (anti-global) peak - Not a peak corresponding to true M.
- 06 2) We are at a point of discontinuity or  where no second derivative exists.

(1) & (2) are common approximations of N.L. funct. In Agura/Stats we can use non

nice $\frac{1}{1+e^x}$ type functions 

Linear Curve fitting: $\text{Func} = \sum_{i=1}^n F_i(x)$
weighted costs, etc.

I don't know if 15313-19 applies to linear "curve fitting" (when not .05-.06 do not interact), non linear Cur fit, tho it would seem to, from cursory analysis.

Can analysis of 15313-19 apply to other curve fitting problems: Clustering?

In which we have to deriv (parameter) to "distance" function

NB 15313-19 gets var for each competing Model, so it can self regulate var: $\sigma \propto \frac{1}{\sigma^2}$

Can we translate this into into equivalent approx of models?

~~$\frac{1}{\sigma^2}$~~ $\approx \left(1 - \frac{2m}{k}\right)^n \approx e^{-2m}$ is approx pc of model. - for z in cont linear model.

So e^{-2} is pc of each cost! (could this be used as a "ruff'n'divvy" measure of PC of continuous params?)

~~Not so simple!~~ $n\sigma^2 \rightarrow n\sigma^2 \cdot e^{2m}$

$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}$ from 15308 PC of copies w. var σ^2 is $\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n$ (n distinct copies)

$\sigma^{-n} \rightarrow \left(\sigma^2 \left(1 + \frac{2m}{n}\right)\right)^{\frac{n}{2}} \approx \left(\sigma^2 \left(1 + \frac{m}{n}\right)\right)^n \rightarrow \left(\frac{1}{\sigma}\right)^n \rightarrow \frac{1}{\sigma^n} \cdot e^{-m}$

So each param has $pc = e^{-1}$

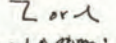
tho precision of measurement of data ($\approx \Delta$) this is $\left(\frac{\Delta}{\sigma}\right)^n \cdot \frac{1}{\sigma^n}$

I certainly don't feel confident about "Significance" of this e^{-m} result!

It's a statement, not of single about a single node, but about to set of nodes w. m costs.

A different way to use this linear regression result: It was obtained (I think in 1962!)

When I integrated all poss/ linear models. Also integrated σ^2 's I think. I may have used uniform p.d. for costs (even tho it didn't converge). - T. idea being that I could choose

2 "flat" cost d.f. so that it would give convergent as "uniform (divergent) d.f." 

Anyway, I could do to some thing w. non-linear problems (having d.f. of .06, but not .05) or to Clustering analysis (using d.f. "distance" function).

→ $\left(\frac{Spec}{156.00}\right)$

NIPS

$$n(\sigma^2 + (r_1)^2)$$



0:157.00 : Ent (linear regression, (un-til) case! : where 2 points in m-space; Probability in space of cells, r = radius.

The n space is a rectangular space of outside points at a distance d from it.

(d = 0): The point in (n+1) space is a point in m space, so P is distance d from origin, It is at distance $x = (d^2 + |r|^2)^{1/2}$ from any pt. r , in m-space.

If 2 d.p.'s have zero first moment. Sum of 2 d.p.'s has sum of second moments.

$$\frac{1}{x^2} \frac{x^2}{x^2} = 1 \quad (x = \text{distance from } P) \quad \text{so} \quad \int_{-\infty}^{\infty} x^2 \cdot \frac{1}{x^2} dx = \infty \quad (\text{b7b!})$$

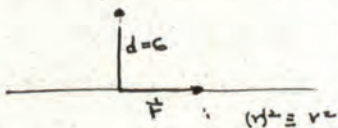
(No!) wt. = $\frac{1}{\sqrt{2\pi} d} e^{-\frac{x^2}{2d^2}}$ (so mult) by x^2 and $\int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2d^2}} dx$?

We have 2 bunches of pairs of form

$$\frac{1}{\sqrt{2\pi} \sqrt{\sigma^2 + r^2}} e^{-\frac{(x-r)^2}{2(\sigma^2 + r^2)}} \cdot e^{-\frac{(x+r)^2}{2(\sigma^2 + r^2)}}$$



Second moments of 2 ($\sigma^2 + r^2$). wt. = $\frac{1}{\sqrt{2\pi} \sqrt{\sigma^2 + r^2}}$



at pt r , have moment product $d^2 + (r_1)^2 = \sigma^2 + r^2$.

$$wt = \frac{1}{(\sqrt{\sigma^2 + r^2})^n}$$

$$\int_{-\infty}^{\infty} dx \frac{2(\sigma^2 + r^2)}{(\sqrt{\sigma^2 + r^2})^n} \cdot \left(\frac{1}{\sqrt{2\pi} \sqrt{\sigma^2 + r^2}}\right)^n \quad \therefore \int_{-\infty}^{\infty} \frac{dx}{(\sqrt{\sigma^2 + x^2})^{n+4}}$$

$$\int_0^{\infty} \frac{dx}{(2^2 + x^2)^{\frac{n}{2} + 2}} \quad \int_0^{\infty} \frac{dx}{(2^2 + x^2)^2} = \frac{\pi}{2 \cdot 2^2} \int_0^{\infty} \frac{dy}{(1 + \frac{y^2}{2^2})^2} = \frac{\pi}{2 \cdot 2^2} \int_0^{\infty} \frac{dy}{(1+y^2)^2} = f(2)$$

$$2^{-2} = 2^{-2} f(\frac{n}{2} - 2) = 2^{-n-4} f(\frac{n}{2} - 2)$$

$$\int_0^{\infty} \frac{dy}{(1+y^2)^2} = f(2)$$

P 294 Gauss: $\int_0^{\infty} \frac{1}{(x^2 + a^2)^n} = \frac{(2n-3)!!}{2(2n-2)!!} \frac{\pi}{2^{2n-1}}$

$$2n!! = 2 \cdot 4 \cdot 6 \dots 2n = n! \cdot 2^n$$

$$(2n+1)!! = 1 \cdot 3 \cdot 5 \dots 2n+1$$

$$\int_0^{\infty} \frac{dx}{(x^2 + a^2)^n} = \frac{\pi}{2} \frac{(2n-3)!!}{(2n-2)!! \cdot 2^{2n-2}}$$

$$1 \cdot \frac{3}{2} \cdot \frac{5}{4} \cdot \frac{7}{6} \dots \frac{2n+1}{2n}$$

$$= \frac{(2n+1)!}{2^n n!}$$

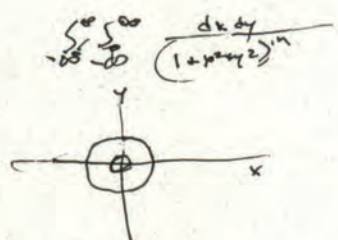
$$\frac{n}{2} \frac{1}{2} \approx \ln n + \gamma$$

$$\prod_{i=1}^n (1 + \frac{1}{2i}) \approx \exp \sum_{i=1}^n \frac{1}{2i} \approx \exp \frac{1}{2} (\ln n + \gamma) \approx e^{\frac{1}{2} (\ln n + \gamma)} = e^{\frac{1}{2} \ln n} e^{\frac{1}{2} \gamma} = \sqrt{n} e^{\frac{1}{2} \gamma}$$

N.B. G & R p 295 bottom of (23!! / 25-1!!)

$$\int_0^{\infty} \frac{2\pi z^2 dz}{(1+z^2)^n}$$

$$\int_0^{\infty} \frac{z^2 dz}{(1+z^2)^n} \quad \text{See G & R p 295 (4.5.)}$$



Note Clustering is a BFG induction problem! What corresponds to Var in time series analysis?

One way: f. f. functions are narrower in 4.3 (143 to 20), the total second moments of the f functions are smaller. "Hr. What is criterion if one introduces more (or less) cluster centers?" No! f. second moment about the mean (= first moment).

Also Note! Index $G^2 \rightarrow G^2 \frac{n+1}{n-1}$ discussion of f. n, G^2 model - was $n \geq 1$ or $n \geq 2$ param model?

Nips.

0: 158.40: If it was a 2 param model, its extrapolation to infinity would be $\approx \frac{n+1}{n-1} \approx 2$ — quite decent from t. Old Maxim results of ~1962, 1972. It certainly seems like a 2 param model.

* one. Flooding of notation for numbers: (Growth) : A 32 bit string gives a set of nos. of = pc: It is a ruffly "log" distrib, but cut off on the low ends. It is symmetric over $\pm x$. It is a reasonable d.f. in many Physics problems. [In many cases, we know, a priori, whether its + or -.

For other cases, the d.f. below. - (int) wouldn't be so dense — More dense than say 100 to 101, but nothing (like fractals)!

I think the moral is that for each situation, one must learn what kinds of d.f.

to use. This "learning" can involve logical reasoning.

[5N]

T. "incompleteness" of UPD is a sad having, it has nothing to do w. applicability of UPD. The really serious problem is that to UPD has to be "learned", that coding considering constraint ~~its~~ its form, but ~~the~~ general ~~is~~ support aspects of its form, but essentially, much of its ~~form~~ form has to be "learned".

T. Math community is enchanted w. the idea that it is known "within a constant factor". This fact is usually of little value by itself. $\rightarrow 160.00 \dots 10$

20

For d.f. an / integers ^{positive} $\log^+ n$ is not bad. For distribution on roots from 1 to ∞ , perhaps $\log_2^+ x$ is tolerable (no using log base 2 is somewhat arbitrary)

Linn said it doesn't converge for (n^x) . Not quite reasonable!

$\frac{1}{n}$ div, $\frac{1}{n^a}$ conv. $\frac{1}{n \ln n}$ div, $\frac{1}{n (\ln n)^2}$ conv? $\frac{1}{n (\log_2 n)^2}$ converges, $\therefore \frac{1}{n (\ln n)^2}$ converges.

$\frac{1}{n} \frac{1}{\log_2 \log_2 n}$ diverges $\frac{1}{n} \frac{1}{\log_2 n} \cdot \frac{1}{(\log_2 \log_2 n)}$ conv.

$\log_2 \log_2 n$ v.s. $\ln n \ln n$ constant $\log_2 n = \ln n \cdot \log_2 e$
 $\ln \log_2 n = \ln \ln n + \ln \log_2 e$

$\ln(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 \dots$

So (?) $\frac{1}{n} \frac{1}{\ln n} \cdot \frac{1}{(\ln n)^2}$ div, but $\frac{1}{n} \frac{1}{\ln n} \cdot \frac{1}{(\ln n \ln n \ln n \dots)^2}$ conv?

If taking these ~~the~~ multiple logs, we consider only large values of n. — small values give log of imaginary values.

Ratio of terms: $\frac{n+1}{n} \cdot \frac{\ln(n+1)}{\ln n} \cdot \frac{\ln \ln(n+1)}{\ln \ln n}$
 $\frac{\ln(n+1)}{\ln n} = \frac{\ln n + \ln(1+1/n)}{\ln n} = 1 + \frac{\ln(1+1/n)}{\ln n} \approx 1 + \frac{1/n - 1/2n^2 \dots}{\ln n}$

" \bar{x}^α : $\frac{(x+1)^\alpha}{x^\alpha} = \frac{x^\alpha + \alpha x^{\alpha-1}}{x^\alpha} = 1 + \frac{\alpha}{x}$ converges if $\alpha > 1$.

0: 159.40 : 159.11-15 (on necessity of freq approx) seems reasonable! So how can we have
 t. $\frac{m+m}{n-m} \approx 2$ ~~is~~ ~~is~~ criterion for sequential prod of reals?
 Note Koz 26 on "Akaike"

03 T. ~~the~~ Prust of 159.11-15 \bar{x} that + density function params must be "rind" - i.e. optimized
 over a large set of sub.cases. ~~There will~~ There will be a hyper order of def. that must be
 chosen a priori, but probably this will have a "reasonable" ERM - that seems

"obvious":
 In a practical case, one could choose a apriori on density and metz on distrib. on density
 of params during decisions. It might be first necessary to decide on what
 to range was ($-\infty \dots +\infty$; $0 \dots 1$; $0 \dots +\infty$; $1 \dots +\infty$; Integers or reals, etc.)
 It may be that there aren't a big (large no. of cases,

Consider the way one normally does statistics! We have no data, a bunch of measurements
 of the length of a certain object. From a priori physics & statistics, we expect a normal
 dist. about some mean. Using no more a priori than that, we take the data fit a
 Gaussian to it, i.e. fit a μ, σ^2 for it.

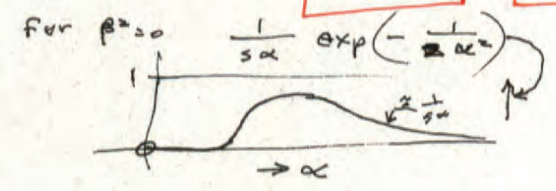
Say we choose a uniform a priori for $\mu (-\infty, +\infty)$ & $\sigma^2 (0, \infty)$. (or $\sigma^2 (0, \infty)$)
 The probability of the model μ, σ^2 will produce n data $\{x_i\}_i$ is $\bar{x}, (\bar{x} - \bar{x}^2) = s^2$

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\bar{x} - \mu)^2 + s^2}{2\sigma^2}} \right)^n$$

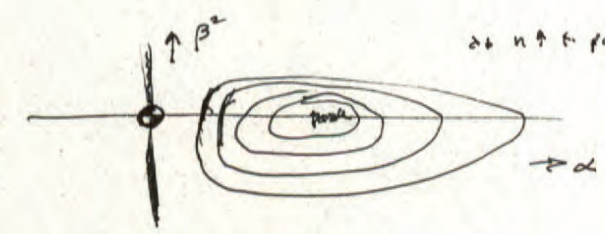
This is max when $s^2 = \sigma^2$ & $\bar{x} = \mu$.

Just look at $\frac{1}{\sigma} \exp\left(-\frac{(\bar{x} - \mu)^2 + s^2}{2\sigma^2}\right)$ is a function of μ & σ^2 . $= \frac{1}{\sigma} \exp\left(-\frac{\frac{(\bar{x} - \mu)^2}{s^2} + 1}{2\frac{\sigma^2}{s^2}}\right)$

Use new variables $\frac{\sigma^2}{s^2} = \alpha$ & $\beta = \frac{(\bar{x} - \mu)^2}{s^2}$ $\frac{1}{\sigma} = \frac{1}{s\alpha}$ $\left(= \frac{1}{s\alpha} \exp\left(-\frac{\beta + 1}{2\alpha}\right) \right)$



for constant α ; we have $\left(\exp\left(-\frac{\beta}{2\alpha}\right) \right) \cdot \frac{\exp\left(-\frac{1}{2\alpha}\right)}{s\alpha}$



as $n \uparrow$ the peak gets sharper & sharper. Its spread of α is β^2 is essentially, $\bar{x} - \mu$ is σ^2 . importantly unlikely for small α .

If the peak is very sharp, the exact
 a priori for μ & σ^2 is not important -
 It need only be "smooth"

We also have "Marginal" d.f. of α & β^2 .
 Marginal d.f. for α , integrates out β^2

NIPS

SM

00:160:40: → Note that a prior is important mainly if $\Sigma \geq 1$ is small!

Also: If a predn. method gives a very narrow output (small Σ or p_i in p_i) This means it's very confident but nearly right! Hrr, ~~the~~ / it's. width "width" is unbiased, This may be what we want!

$\Sigma \approx \exp \sum p_i \ln p_i$ a probly wt. = $\prod p_i^{p_i}$.

So this is the Q: Is the unbiased width of D.F. related to the wt. of the param? ← (i.e. its pc)

To do induction, getting pc of params using ideas of 160.03ff, is possi, reasonable, but often diff. Getting $\sum p_i \ln p_i$ may be easier: It seems easier in linear regression/curfit. L: "width" of dif.

$$\int_{-\infty}^{\infty} \frac{1}{\sigma} e^{-\frac{x^2}{2\sigma^2}} dx$$

$$\int_0^{\infty} z e^{-\frac{z^2}{2}} dz = 1$$

$$\ln \frac{1}{\sigma} + \sigma \cdot \frac{1}{\sigma} \rightarrow \frac{1}{\sigma} e^{\sigma x}$$

$$\int_0^{\infty} x^n e^{-\frac{x^2}{2}} dx = 2^{-\frac{n+1}{2}} \left(\frac{\sqrt{\pi}}{2}\right)!$$

Get p 337 to p 41 note.

10

SN SM v.g.!

To compute info in a strat: Usually a strat will contain several continuous

params & its deriv will involve a few discrete params. The continuous params

can be evaluated by running the data repeatedly on "near by" p's, & getting

empirical "Hessian". The discrete params may be mainly a priori & ~~empirical~~

cost no probly. Some discrete params will be ~~empirical~~ have been empirically chosen

& have associated bit costs.

Example: We inform it prior over n bits of info in the strategy (obtained this way)

& its empirical yield is y , then its expected yield is $y \cdot 2^{-n}$.

"Obtaining info in Hessian" is not clear — perhaps guess at a priori width of

of distribution — of ranges of continuous params.

20

In my latest strat (that seemed to work) there were (I think) 2 continuous params:

1) Buy Threshold 2) Smoothing const (n yrs); 3) was rather non-critical.

1) was a bit critical w.r.t. large losses. Yield was rather dissonant in that param! & gains were would vary inversely w. inv of 1) (There were other params (incl. no. of stocks in the "set being considered")

Re: 2) : when doing lower thresholds & higher thresholds, try to find aux. properties of both that & reliability of acceptance/rejection. Since this is an AND operation, if \uparrow noise (badly) & its a posteriori ($\rightarrow \uparrow$ cost) also new data needed to test.

Anyway, it could \uparrow yield a lot

30

Some "other aux properties": Correlations (+ or -) w. other stocks/indicators.

& Self correlation (obtainable if I really bet "on-line" & as a "broker")

t- "Akaike factor" ~~is~~ .26
Looks like Good explain, understanding.

00:60 of all PST's very change. As one works, $T_0 \rightarrow T_r$ (recovery time).
We always record the PST's with current T_r . When a new PST become "best best",
we switch to it. We can switch back to a previously discarded PST —
returning to work (Time spent on it).

Another way to think about it: for each G value, there is a $\frac{PC}{CC}$ ordering among the PST's,
which is mostly "optimum". Now, I don't see how this would help solve the problem: perhaps:
perhaps: for each G value, there is a Σp_i : total prob that one could solve a
problem with best G value.

ANY way 162.27 ff: (Just how Phase 2 does OZ problems) is at present, very unclear

10 Conceivably, a better understanding of it would enable me to ~~write~~ outline a
Soln. to the "Time varying OZ" problem.

It is a problem that Must be Solved. $\rightarrow 165.00$

"Things to do" 155.1.00

153.20 - 154.18: Derbs errors in PEMS: whether to values of General values is in brief for UPD is unclear.

151.26 - 152.12: Error is incomputability $P_M^T(x); P_M(x), M(x)$. } That "incomputability of UPD"
This is perhaps a directly utilizable unloop. } is unimportant, is clear

20. Re: "Unbiased error" for $P_M^T(x)$: What this means is unclear: Say we have to spend all budget
to use one doing digital prodn: Unbiased 2 pc for each symbol. Say on some occasions
it is "unbiased": Then, if we assign real use to each symbol we can do this in any
ways: It would be a seem unlikely that all assignments would end up being
unbiased! (?)

162.26 Looks Good!
[508] perhaps easy way to analyze/explain/understand σ^2 of regression $\rightarrow \sigma^2_{obs} \times \frac{N+m}{M+m}$
Say σ^2 is true σ^2 of θ . $\sigma^2_{u} \cdot \frac{N-m}{N} = \sigma^2_{obs} = \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}$ is expected observed variance part: $< \sigma^2_{u}$
Since costs are wrong, σ^2 for θ is $\sigma^2_{u} \cdot \frac{N}{N-m}$: T. product is $\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N-m} \cdot \frac{N}{N-m} \cdot \frac{N}{N-m}$
It may be that this just gives "Akaike's effect" — if does not give cost of telling law
many costs are not out / form of t. var. regress. eq. $\rightarrow 164.30$

30 Looks Good!
Her. 136 may be on a set that is always usable for concurrent parts of a Plan dem.
It assumes uniform a prod for costs, probly uniform a prod for cost (?)

to go from 1 parameter to a linear cost: t. singl param can be looked at as
confid to constant, 1, so it genes. easily! $\rightarrow 164.01$

24 The exact mechanics of how the uniform applied out costs gets to Akaike's result is unclear.
It may be simply integrating like in Maxim 1962! (see 153.00 for a try at it!): I think there
may be a very simple way to see it by relating Matrices in the space — which is what I was
thinking about in 1962 (or 1972). Perhaps the idea is that all span matrices can be

00: 163.40 rotated so they are diagonal.

01: 63.34:

$$\sigma_u^2: \text{observed } \overset{\text{mean}}{\text{var}} = \frac{(n+m)\sigma_u^2}{n}$$

$$\text{predicted var for } X_{n+m} \Rightarrow \sigma_u^2 \frac{n+m}{n}$$

$$\therefore \text{predicted var} = \frac{\text{observed var} \cdot \frac{n}{n-m}}{\sigma_u^2} \cdot \frac{n+m}{n} = \sigma_{\text{obs}}^2 \cdot \frac{n+m}{n-m}$$

[still have to work out details! Consider vectors in M space; Data z vector in n space. Model on m dimensional subspace. We try to fit model vector as close

as possible to data vector. This is "projection" of data vector in to M space.

T. ^{True} model, \vec{u} ^{data} is a vector in n space, plus a random ^{subspace} Gaussian n vector:

So the projection of \vec{u} vector, is its correct m vector plus an m dim. isotropic vector.

$$\vec{u} = \vec{u}_m + \vec{R}_m$$

\uparrow m dim \uparrow m dim \uparrow n dim

Why the orthogonalization!

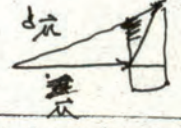
$$\vec{u} + \vec{R} = \vec{d}$$

\uparrow m dim \uparrow n dim \uparrow n dim

\vec{u} = correct model in m space.

$\vec{u} + \vec{R} = \vec{d}$ = data vector.

$\vec{u} + \vec{R} = \vec{d}$ = projection of \vec{u} into n space.

sq. error betw. \vec{u} and \vec{d} is m (\vec{R} has $\sigma^2=1$ & independent) 

$$|\vec{d} - \vec{u}|^2 = n$$

$$|\vec{d} - \vec{u}|^2 = n - m$$

\leftarrow projection \leftarrow imp. result.

\vec{d} deviates from \vec{u} in only $n - m$ dimensions.

smaller observed from \vec{d} error. $(n-m)$ (n)

How do we take \vec{u} and \vec{P}_m into n space (new data pt)? They random + same. (?) \vec{u} , however has an ~~extra~~ extra dimension.

\vec{d}_m is \vec{P}_m random + same? They still have only m dim. The new data pt + n vector is obtained from old data n vector, by adding ^{sub} a small vector to it, i

01: 164.34

Re: Akaike! For long data sequences, the contribution of the precision of the coeffs, to the overall p.c. of the code, Becomes dominant over other coding costs.

i.e. \rightarrow look at expressn. for best in Sol 86.

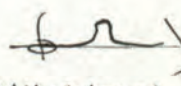
00 (Spec) Req: Phase 1 OZ problem "Soln"; T: Lsck soln is really not so great! We have time limit $T \in (CB)$.
 01 If $P(\cdot)$ is t. PC of P if $P(PST_i)$ is t. PC assigned to PST_i ; Then we spend $T \cdot P(PST_i)$ on PST_i .
 Note, how fast $P(PST_i)$ is to probly (perhaps) that PST_i will be best in time T , (not in time $T \cdot P(PST_i)$). T "meaning" of t. soln. is that no spent at least $P(PST_i) \cdot T$ on t. "Best poss. PST_i ".

06 Say we had t. $P_i(T, G)$ curves for all PST_i 's wrt t. particular problem of interest
 How much time to spend on each PST_i \rightarrow ~~max~~ t. max G obtained from all trials was Max?
 09 It may be (not unhelpful) that t. soln. is to put all time into one PST_i ! - But they would
 0 take greater expected G than distributing cc among t. PST_i 's "more evenly".
 \rightarrow One (perhaps serious) Advantage of distributing cc among t. PST_i 's is that we get more
 into about them, (to be used in future problems).

One "way out" of .09, is that for many PST_i 's doing OZ problems, they make many trials. In view of Phase trials \rightarrow ~~lots of long~~ lots of long contake place a $P(G, T)$ curves of
 how PST_i 's can change so we can reably jump from a PST_i to another
 "washing" w.o. losing much cc.

18 So: The Moral is: If I would be satisfied w. t. Phase 1 soln of OZ of
 19 00-.01, then 2 sample "Phase 2" version get $P(G, T_0) \equiv P_{T_0}(G)$ curves
 20 \downarrow T. destdnt. OZ problem deen)
 21 t. cost is very similar to soln. for INV probs as desc'd in ξ_2 of
 22 "Preliminary Preliminary Report".

On t. other hand, t. .00-.01 soln for OZ Phase 1 is not so hot a
 idea as 006-.18 suggest ways ~~to~~ to improve both Phase 1 & Phase 2.
 Improving Lsck over OZ problems is a harder problem than WON for INV probs.

30 (SN) On t. other hand, if we use 2 or 3 param $P(T)$ curves. ()
 as we probably will, t. optimum ~~solns~~ ^{solns} may be poss. Br INV at least.
 (Also note, WON solns so far do not allow long during trials or
 32 even between trials!) \rightarrow (for WON \rightarrow 166.00)

WON

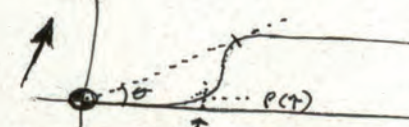
20:165.32: A useful, won soln, for Gaussian P(T) curves; or gamma d.f.

Gauss: $P(T) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(T-T_0)^2}{2\sigma^2}\right)$: Params σ, T_0, σ^2 .
 ↑ T_0 mean

Gauss $P(T) = \frac{a^k T^{k-1} e^{-aT}}{\Gamma(k)}$: Params a, k, b .
 ↑ Normal const: involves $\frac{1}{n!}$ or $\frac{1}{\Gamma(k)}$.

doesn't have proper origin at $T=0$.
 does have nice origin at $T=0$.
 finding first & second moments is diff, hrr.

2 tasks: which to work on to finish ~~the~~ ^{first} (OR next):

10  This is the ~~the~~ soln. See 169, 35-120, 07 for details, proof

Work on curve w max Θ
 ↳ It is a very "ruff" soln. - Needs much work. But I may be more comfortable positive to do it since I have recognizable shapes for P(T) curves. (3 params)

after working on P(T) curve for time \uparrow
 We have new curve $= P(T-T)$
 $T = T - \tau$
 so Θ increases at first.

eventually, Θ hits a peak and then begins to \downarrow . The peak Θ is when slope of $(.10 \curvearrowright)$ curve is max - after Θ is max, we jump to another PST's $P(T)$ curve

or look at it this way:



as we measure curve, we always have a Θ that gives slope of steepest line thru present point, that intersects $P(T)$ curve.

Note for Gauss d.f.

$f(T) = \frac{1}{\Gamma(n)} T^{n-1} e^{-T}$

$f(n) = T^n e^{-T}$
 $f'(n) = n T^{n-1} e^{-T} - T^n e^{-T} = n f(n-1) - f(n)$
 $f(n) = n (f(n-1) - f(n))$

$\int_0^{\infty} T^n e^{-bT} dT$ normalize: $z = bT$
 $\frac{1}{b^{n+1}} \int_0^{\infty} (bT)^n e^{-bT} dbT = \frac{\Gamma(n)}{b^{n+1}}$

$\frac{M_0}{2} = 1$ 2 normalized functions $\frac{b^{n+1}}{\Gamma(n)} T^n e^{-bT}$
 $\frac{M_1}{2} =$ zeroth moment \rightarrow \int
 $\frac{M_2}{2} =$ first " " \int

Zeroth moment of $z T^n e^{-bT}$ is $\frac{z f(n)}{b^{n+1}} = M_0$

first moment is $\frac{z f(n+1)}{b^{n+2}} = M_1$

second moment $\frac{z f(n+2)}{b^{n+3}} = M_2$

$z = M_0$
 $\mu = M_1 / M_0$

$\sigma^2 = \frac{M_2}{M_0} - \mu^2 = \frac{M_2}{M_0} - \frac{M_1^2}{M_0^2} = \frac{M_2 M_0 - M_1^2}{M_0^2}$

30

WON

20:166.40:

Consider 2 param PCT's:
Ordering choices by $\frac{a}{u}$ is not bad.



Trouble is, if works working in many PST's in (1), we don't get any results until time μ .

Say I have a bunch of PST's like Part 2. What order should I try them? I guess Part there is no gain in switching between PST's before μ is spent.

My guess is $\frac{a}{u}$ ordering. Since $a \in \mathbb{R}$, no customer has a very large a w.c. large $\frac{a}{u}$.

09

e.g. 2 expected time to solve for $\frac{a_1}{u_1}$ followed by μ_2 - μ_2

0 -

$$\mu_1 \cdot z_1 + (\mu_1 + \mu_2) z_2 = \mu_1 \cdot (z_1 + z_2) + \mu_2 z_2 = \mu_1 z_1 + \mu_1 z_2 + \mu_2 z_2$$

11

$$\mu_2 z_2 + (\mu_1 + \mu_2) z_1 = \mu_2 (z_2 + z_1) + \mu_1 z_1 = \mu_2 z_2 + \mu_1 z_1 + \mu_2 z_1$$

So $\mu_1 z_2$ v.s. $\mu_2 z_1$ or $\frac{z_1}{u_1}$ v.s. $\frac{z_2}{u_2}$ is expected;

T. first Gamb H.
Mem. may
Prove Part
more exactly!

Then I'm not quite sure Part 10 v.s. 11 comparison is legit!

15

For INV, problems, this suggests that ~~they~~ trials in $\frac{a_i}{u_i}$ order would be best

Also this or z (i. Param) is not zero, we write want to use $\frac{z_1}{u_1 + 26}$ or $\frac{z_2}{u_2 + 36}$

17

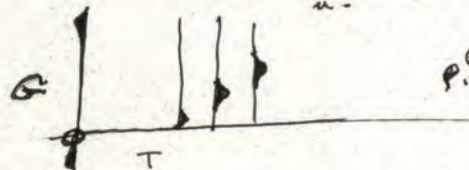
for ordering is take $\mu_1 + 36$ for testing.

So \rightarrow 169.00 for which better soln!

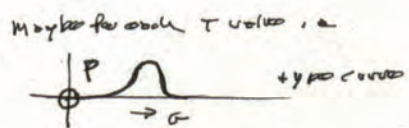
T. For z , seems to "solve" the INV problem!

20

what about z probs? I would guess Part 1. $P(T, G)$ curves would be much different from $\frac{a}{u}$.



$P(z)$ is in vert. direction



Maybe for each T value, a type curve.

"stochastic complexity"
Is good at .09 for min expected time of soln? \rightarrow 169.19-20

One strategy: for each PST, do a T_2 such that we expect some G .

We may go a bit past that T .

30

Perhaps better way: Pick PST w. max expected G in $z = T_0$ ($= CB$) (\equiv PST₀)

As we work on z , T constraint time is $T_0 - T$. We simultaneously ~~search~~ ~~for~~ ~~best~~ ~~expected~~ ~~G~~ ~~in~~ ~~receiving~~ ~~"all"~~ ~~PST's~~ to see which one will give best G in receiving $T_0 - T$ time.

We jump to PST that is best.

As we work on PST₀, $P(T, G)$ curves of other PST's changes since $P(z)$ (optimized) function that maps 2 PSTs in to $\{P(T, G)$ curves $\}$ changes as we work on PST₀.

For time being, use 165.00-01 for Phase 1, z is use 165.19-.22 for Phase 2 z

The discussion of z on 165.00-.40

For INV probs; z is 165.19-.17 (cash v.s. : its not quite such, but is $\frac{4.6}{SPAC}$)

Can't use something like for Phase 1 induction? \rightarrow 168.25

Context.oo

9th.

20:0 : on "Context": in Phase 1. I had been thinking of; for each token, it would have different pc's in each "Context". So if contexts were each a component of the dimension, w. k contexts, each token has an assoc. 4 dim vector of pc's.

It looks perhaps better to do it in reverse: for each Context, have a vector giving pc's of all tokens. #

At any rate, we have a k x L matrix: k contexts, L different tokens. I think Jaeger in OOPS that of vectors of pc's of tokens, being stored for trial access: This matrix of context vectors were tokens.

More Generally U. had OOPS "computing" the pc of each token, each time a new token was chosen. Sounds very general. Hrr, he only had maybe 5 units that did token pc matrix; & he only actually used 1 (boost) in successful pems

The "context" matrix also include SSZ info. (Number/dimension)

If >1 context applies to a situation, its have to ~~add them~~ add up the vectors because the SSZ's can overlap - this gives "double counting".

The OOPS "boost" gives a weight of a context, the system (as is) does not associate any features of the problem w. an idea as to which "q" to "boost" (a "q" is a frozen pem). When I mentioned to the low pc ~~at~~ of ^{frozen pem} particular frozen pem of it.

Machine is "Mature". Jay suggested that the system would have to develop ways to "narrow down" to ^{frozen pem to be "boosted"} "boosted".

I should go back to read my stuff on "Gated context" - just how it was defined is just what M(early) discussy of "oops" compared it "computation of pc's of tokens" w. my "context". At first his method looked better: but then my Gated idea of "context" ^{begs into look u.g.: try to find thresholds.}

its Gorc is - "what is it trying to do"

25: 167.40 : on ^{172.09} ~~170.09-07~~ ^{170.09-17} for Phase 1 induction: For phase 1 induction one simply lists ⁵ a founds in pc order. This can be somewhat modified by Lep's rule is the discovery of sub-structures.

Sub-structures can become ~~more~~ more useful if we have several domains of induction, w.

"Indexing" on the Q's. This means that the token d.f.s for each index can be different yet they can have common sub-structures (a common regular to focus as well. ?)

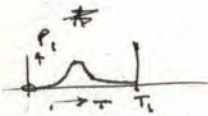
Note in trying to find $O^j \Rightarrow \sum_{j=1}^k O^j(A_i | Q_j)$ is max,

A token is allowed to forget all of $O(A_i | Q_j)$ pairs: Hrr, $\sum_{j=1}^k$ the pc of O^j will be indep of this data.

NIPS

WON

SO: 167.17: Generalize to:



Which order to my team?

$$\int_0^{T_1} p_1(t) dt \text{ : assoc w. } p_1 \text{ we have } \int_0^{T_1} p_1(t) dt \text{ expected soln. time. } \int_0^{T_1} p_1(t) dt \text{ probly of } p_1 \text{ soln}$$

Anyway: After solving ~~that out~~, say we have 2 tasks α_1, α_2 w. finite time cutoffs, ~~can~~ T_1, T_2 , resp. We can break α_1 & α_2 to

$T_1 = T_2$ resp, so we could do:

α_1 $t=0$ to T_1 then α_2 $t=0$ to T_2

α_1 $t=0$ to T_1 then α_2 $t=0$ to T_2 then α_1 $t=T_1$ to T_1

α_1	1	2	3	4
α_2	3	4	1	2
	3	4	1	2
	3	1	2	4
	3	1	2	4

SW Some time ago, I was considering

$$f_1(t) = \int_0^t p'(t) e^{-\int_0^t p(s) ds} dt \text{ v.s. } e^{-\int_0^t p(s) ds}$$

$$f_1 = -\frac{d}{dt} f_2$$

$f_2 = e^{-x}$

Anyway, it was not.

Alternatively $f_1 = p'$; $f_2 = p$ so $f_1 = \frac{d}{dt} f_2$ ~~there is a minus sign~~
 $\int_0^T p'(x) dx$
 interplay
 $f_2 = x$, $f_1 = \frac{dx}{dt}$, but ~~there is a minus sign~~

Maybe $E = \frac{\int_0^T p(t) dt}{T}$

6 ways: 3 1 2 4
 which is best? 3 1 4 2

So corresponding not exist.
 int. f_1, f_2 formulation was con
 defining $f_1 \equiv f_2$ is ~~more~~
 The party of success at T and failure from 0 to T
 looked like 2 Laplace & EM

$$f_1 = x' e^{-x}, f_2 = e^{-x} \text{ : I was thinking } x' e^{-x} \text{ but it isn't!}$$

19 **N.A** in ~~coeff~~, I have to consider "Expected time to soln", but ~~it~~ should not
 20 include cases where time to soln ~~is~~ is.

22 $f(t) \text{ (of (1)) is correct for getting expected time. } \int_0^{T_1} T p'(t) e^{-\int_0^t p(s) ds} dt.$
 $= \int_0^{T_1} T p e^{-p} dt.$

$\Rightarrow \int T p' dt = \int T' f = -TF$ $-\int T f' dt = TF + \int f$ $| f_2 = e^{-p}$
 say $-f' \equiv h$
 so we want $\int h(t) dt = \min.$

If I do cond 1, I will win a fraction of time $\equiv P_1$
 $= \int_0^{T_1} h_1(t) dt$; It will take out average time $\int_0^{T_1} T h_1(t) dt = CC_1$
 To do cond. 2, then cond 2 $P_{1,2} = (1-P_1) \cdot P_2$
 $CC_{1,2} = T_1 \cdot (1-P_1) + CC_2$
 $CC_{2,1} = T_2 \cdot (1-P_2) + CC_1$
 $CC_{1,2} = \frac{T_1(1-P_1) + CC_1}{(1-P_1)P_2}$ $CC_{2,1} = \frac{T_2(1-P_2) + CC_2}{(1-P_2)P_1}$

Looks reasonable
 Simple is not always better!
 he writes!

$$\int_0^{T_1+T_2} (f_1' + f_2'(T-T_1)) dt$$

$$\int_0^{T_1} T f_1' dt + \int_0^{T_2} (T+T_1) f_2' dt = \int_0^{T_1} T f_1' + \int_0^{T_2} T f_2' + T_1 f_2 \Big|_0^{T_2}$$

so $T_1 f_2 \Big|_0^{T_2}$ v.s. $T_2 f_1 \Big|_0^{T_1}$ or $\frac{T_1}{P_1} \text{ v.s. } \frac{T_2}{P_2}$ or $\frac{f_1 \Big|_0^{T_1}}{T_1} \text{ max}$

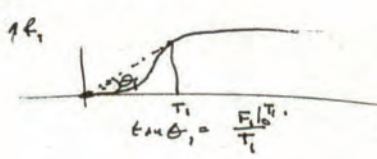
As
 See 174.00 ff
 for worry
 about our
 "proof"

NIPS

WON: $B_{\text{max}} R_{\text{ro}}; .09 + (.03 - .07)$

6:169

R_{ro} $\frac{f_2 | T_2}{T_2}$ looks fine!



03
Very simple jump to 172.00

In this case of $t_1 \neq t_2$:

T. best way, do f_1 until past T_1 , $\frac{df_1}{dt} = \frac{F_1(T_2)}{T_2}$, then do f_2 out to $T=T_2$

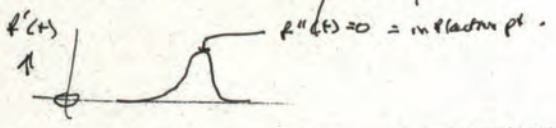
Then do f_1 & f_2 similarly as to keep f_1' & f_2' the same.

This technique will work in any no. of curves. — it's particularly easy if they all ~~are~~ have $f'(0) = 0$. Better, since $f'(0) = 0$.

Note now $f'(t) = p(t) e^{-\int_0^t p'(t) dt} = p'(t) e^{-(P(t) - P(0))}$

So this looks like an ideal soln. for INU problems!

Soln is particularly easy if f' only has one reflection pt. i.e. $f'' = 0 \Rightarrow \frac{d}{dt} f' = 0$

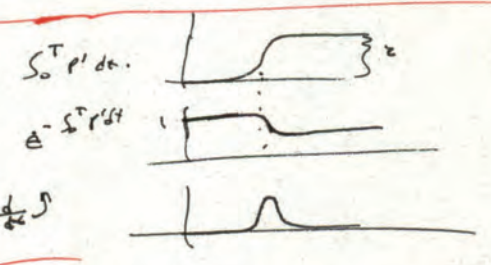


→ 172.00

Q: I have recently assumed that $f'(t)$ had R_{ro} simple shape,

but look at $-p' e^{-(P(t) - P(0))}$ it \neq simple.

look at $e^{-\int_0^t p' dt}$. Then take its derivative, but first look at $\int_0^t p' dt$.



So f' does look ok, but it do have a sign change!

We will get TM to find optimum f' , f directly, rather than by finding p' , p first (169.20-22)

T-params up as two most important. $\frac{\partial^2}{\partial t^2}$ determines (∂^2) order in which PSTs will

be ~~involved~~ recruited. Sky $\frac{\partial^2}{\partial t^2}$ is in correct order so $\frac{\partial^2}{\partial t^2} \geq \frac{\partial^2}{\partial t^2}$.

Then, when PSTs has started, we will be timely working on all PSTs $\Rightarrow \int \partial^2$.

They will all have same slope.

Actually, Parallel L such may be regarded as a special case of $(.03 \rightarrow .07)$. we work on all cards so that their $\frac{p}{cc}$ R_{ro} far away about the same, which \Rightarrow

Then, f is obtained by code length directly. In $.03-.07$, we could vary the code length to be $(\int_0^t f'(t) - f(0))$, in which case we get identical same as L such.

O.h. now for 02 problems \rightarrow 171.00

Spec 171.00

NIP

WON approach to OZ probs.

0: 17040 : OZ problems:

Say we have \geq PST surfaces $P_1(T, G)$, $P_2(T, G)$.

We have $KB \geq T_0$. Chances $T_1, T_2, \geq T_1 + T_2 = T_0 \geq$ then 2 resultant $P(G)$

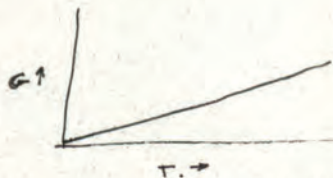


We can do a ruff & dirty so on by using 2 params (μ, σ) approx.

A ruff ordering is by μ_1, μ_2 .

so maybe μ_1, σ_1 is not such a good index:

Use just μ_1 to start. $??$ No!



Given 2 curves at $(0, L)$, what is expected value of max?

T. probability/density of 2 peak at G_M is $P_1(G) \cdot P_2(G)$

We want max expected value:

$$P_1(G) \cdot \int_0^G P_2(G) dG + P_2(G) \cdot \int_0^G P_1(G) dG$$

$$\int_0^{\infty} G F_1'(G) \cdot F_2(G) dG + \int_0^{\infty} G F_2'(G) \cdot F_1(G) dG \leftarrow \text{maybe should be mult by } (-1)$$

$$d(xyz) = \left(\frac{dx}{x} + \frac{dy}{y} + \frac{dz}{z} \right) \cdot xyz$$

$$d(x \cdot yz) = dx \cdot yz + x \cdot d(yz) = dx \cdot yz + x \cdot dy \cdot z + x \cdot dz \cdot y$$

$$d(xyz) = dx \cdot yz + dy \cdot xz + dz \cdot xy$$

$$xyz = \int x'yz + \int y'xz + \int z'xy$$

$$G F_1 F_2 \Big|_0^{\infty} = \int_0^{\infty} G F_1' F_2 + \int_0^{\infty} F_1 F_2' + \int_0^{\infty} G F_2' F_1$$

$$\text{So } \int_0^{\infty} G (F_1' F_2 + F_2' F_1) = \underbrace{(G F_1 F_2) \Big|_0^{\infty}}_{\neq} - \int_0^{\infty} F_1 F_2 \quad \text{unlucky!}$$

I guess I want $\int_0^{\infty} G F_1'(G) \cdot F_2(G) (1 - F_2(G)) dG$

+ $\int_0^{\infty} G F_2'(G) \cdot F_1(G) (1 - F_1(G)) dG$

$$= \int_0^{\infty} G (F_1' F_2 + G F_1 F_2') - \int_0^{\infty} G F_1' F_1 F_2 - \int_0^{\infty} G F_2' F_1 F_2 - \int_0^{\infty} F_1 F_2 \quad \text{(2.5)}$$

My suspicion is that its most likely that usually we will want to put all of time into our "Apparently best" PST. We may do this until (using II updating) it is clear that that is no longer a "Best bet".

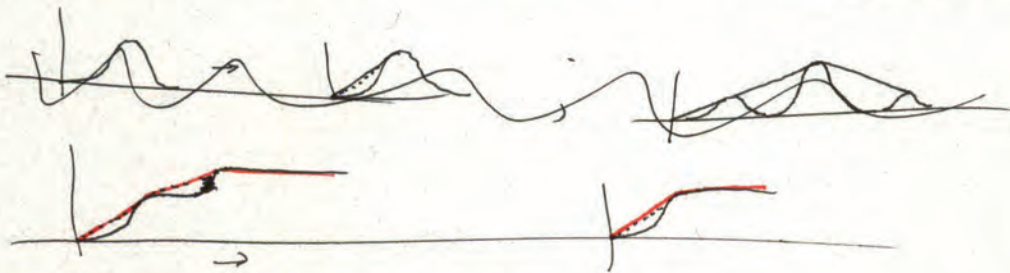
If we don't have much data or strong ideas on the $P_i(G, T)$ curves, then we may want to spend a small amt. of time on many PSTs (with present problem) to get some evidence of which will be best.

(SAC)
→ (72.20)

4.4.03
NIPS

WOP (ENV problems)

170.15
0:170.07: ENV probs: T search process is rather simple! first convexify all $P(T)$ curves:



categorize slopes into discrete bins.

09
Work on Bins in slope order. At least slopes first.
That's it! : A min amount (Δ) is spent on each curve; $\Delta \gg$ transition cost to a new curve. Thus Δ may cause a jump over > 1 bin. i.e. after working on Δ curve, it is put for time Δ , it is put into a new (later) bin.

14
15
Convexity means that slope is monotone (not \uparrow) in T.

16
17
while we are working on a dotted st. line in 0.5, we are really "cheating", but on 1.2 average, we will be doing the right thing. The time between Bumps will be small times relative to Δ .

total time spent on solving problems. So R is "bin size" will not be small. \rightarrow (174.00)

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

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

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

171.35
172.20-40
172.40

Re: "Trials": In general, most OZ ~~methods~~ using trials are iteration methods - (SS ≥ 2) or Ray ≥ 2 Past trial to get new trial (Gen. Alg., Sim Anneal, Hill climb ~~using~~ using extent w/out previous sets): Non-linear Past ~~trials~~,)

Different optz methods use different "approx" to do ϵ extrapolation. Most all methods use fairly "Greedy" (little if any "look ahead").

For ~~the~~ PSTs of this sort (very ~~trials~~ on corpus), we can use $PST(\epsilon)$ Past ~~trials~~ ^{what appears to be "best" trial} ~~trials~~ ^(best trial)

~~Previous~~ "most promising" future trial. - which is quite Greedy!

We can modify procedure to do L.A. ~~at~~ ~~steps~~ 1, 2, 3 ... jumps into future.

Non-trial "optz methods": {linear regression using solution of history}, some ANN methods, but mainly

logical analysis. I haven't got much updated about this. I don't expect my ~~to~~ ~~be~~ ~~able~~ to do it until it has tried to work logical problems as part of its "Math" training.

So: Pruby f. soln for INV problems of 170.00-40; 172.00-15 is ok.

The fragmentation of OZ probs viz $\in P_i(S, T)$ curves of 171.32-40, 172.20-40, 173.00-12 is a general approach that's pruby correct "insight" - that Updating is an imp. part of it. Entire Strategy, "in spirit"

Since $\in P_{i+1}(QA)$ is domain extract very from normal OZ problems, ~~and~~ ~~is~~ ~~closer~~ ~~to~~ ~~INV~~ ~~problem~~ ~~soln~~) I ~~think~~ ~~I~~ ~~can~~ ~~use~~ ~~t.~~ ~~recent~~ ~~WON~~ ~~approach~~ ~~of~~ ~~for~~ ~~Phase~~ ~~1~~ ~~QA~~. The idea is ~~to~~ ~~try~~ ~~to~~ ~~find~~ ~~a~~ ~~code~~ ~~part~~ ~~to~~ ~~corpus~~ - which looks like INV problem. So try to find out how ~~to~~ ~~QA~~ ~~problem~~ \Rightarrow INV can be mapped into INV WON type soln.

$F(x) = True$; find x : $H_2(F) \Rightarrow x$

$2^j \cdot O^j(\text{corpus}) = H_2$ $H_2(\text{corpus}) \Rightarrow O^j(\text{corpus})$ $2^j \cdot O^j(\text{corpus})$

Alternatively find $p \in M(p) = \text{corpus}$: try p 's in pc order (Revis \Rightarrow 3 input vnc approach to induction)

.27 is it really bad at all! It could be f. best way: That all induction methods involve

finding trials in t . corpus, that have to be expressible as "short codes".

SN A poss. trick: Some short codes take long time to verify: A ppm test could look at \geq short codes for t . corpus & ~~give~~ give good reason to believe it would code t . corpus, but take a long time... Something one can look at a ppm's ~~src~~ ~~test~~ if must converge, or one can find indications while it's running on how long it will take to converge.

WON INU probs

(172.17)
169.35
169.40

Worry about \Rightarrow WON proof for Inv of 169.35

Print its layout so that P_1' first then P_2' first vis. P_2' first then P_1'

02 This gives + functions $\int_0^{T_1} P_1' e^{-s_0 t} dt + \int_0^{T_2} (T+T_1) e^{-s_1 t} P_2' dt + \int_0^{T_2} P_2' e^{-s_0 t} dt$
constant factor of $S_1 =$ probab factor of S_1 at time T_1

$$P_1' e^{-s_1 t} = -\frac{d}{dt} e^{-s_1 t}$$
$$P_1' = -\frac{d}{dt} F_1 \quad ; \quad F_1 = e^{-s_0 T} P_1'$$

$$\int_0^{T_1} F_1' dt + \alpha_1 \int_0^{T_2} (T+T_1) F_2' dt = \int_0^{T_1} T F_1' dt + \alpha_1 T_1 \int_0^{T_2} F_2' dt + \alpha_1 \int_0^{T_2} T F_2' dt$$
$$F_2 \text{ then } F_2 = \int_0^{T_2} T F_2' dt + \alpha_2 T_2 (F_2(T_2) - F_2(0)) + \alpha_2 \int_0^{T_2} T F_2' dt$$

$$\int_0^{T_1} T P_1' e^{-P_1} dt + \int_0^{T_2} T P_2' e^{-P_2} dt$$
$$P_1 \equiv \int_0^{T_1} P_1'$$
$$S_0^{T_1} T F_1' + \alpha_1 \int_0^{T_2} T P_2' + \alpha_1 T_1 (F_2(T_2) - F_2(0))$$
$$S_0^{T_2} T F_2' + \alpha_2 \int_0^{T_2} T P_2' + \alpha_2 T_2 (F_2(T_2) - F_2(0))$$
$$\alpha_1 = e^{-P_1(T_1) - P_1(0)} \quad \alpha_2 = \exp(-P_2(T_2) - P_2(0))$$

N.B. $\int_0^T P'(t) dt$ can be > 0 but can be > 1 .

$$\int_0^{T_1} T P_1' e^{-s_0 t} dt + e^{-s_0 T_1} \left[\int_0^{T_2} (T+T_1) P_2' e^{-s_0 t} dt \right]$$

Should $\alpha_1 = 1 - e^{-P_1(T_1)}$

$$\int_0^{T_1} T P_1' e^{-P_1} dt + e^{-P_1(T_1)} \left[\int_0^{T_2} T P_2' e^{-P_2} dt + T_1 \int_0^{T_2} P_2' e^{-P_2} dt \right]$$

$$\int_0^{T_1} T P_1' e^{-P_1} dt + \alpha_1 \int_0^{T_2} T P_2' e^{-P_2} dt + \alpha_1 T_1 \int_0^{T_2} P_2' e^{-P_2} dt$$
$$= - \frac{e^{-P_2}}{s_0} \Big|_0^{T_2} = -\frac{1}{s_0} e^{-P_2}$$
$$= 1 - e^{-P_2(T_2)}$$

$$\int_0^{T_1} T P_1' e^{-P_1} dt + \alpha_1 \int_0^{T_2} T P_2' e^{-P_2} dt + \alpha_1 T_1 \int_0^{T_2} P_2' e^{-P_2} dt$$
$$\beta_1 + \alpha_1 \beta_2 + \alpha_1 T_1 - \alpha_1 \alpha_2 T_1$$
$$\beta_2 + \alpha_2 \beta_1 + \alpha_2 T_2 - \alpha_1 \alpha_2 T_2$$
$$\beta_1 + \alpha_1 \beta_2 + \alpha_1 T_1 (1 - \alpha_2)$$
$$\beta_2 + \alpha_2 \beta_1 + \alpha_2 T_2 (1 - \alpha_1)$$

$$\beta_1 - \alpha_2 \beta_1 + \alpha_1 T_1 (1 - \alpha_2) = (1 - \alpha_2) (\beta_1 + \alpha_1 T_1)$$
$$\beta_2 - \alpha_1 \beta_2 + \alpha_2 T_2 (1 - \alpha_1) = (1 - \alpha_1) (\beta_2 + \alpha_2 T_2)$$

$$\frac{\beta_1 + \alpha_1 T_1}{1 - \alpha_1} \text{ v.s. } \frac{\beta_2 + \alpha_2 T_2}{1 - \alpha_2}$$

$$\frac{\beta_1}{\alpha_1} + T_1 = \frac{1}{\alpha_1} - 1$$
$$\alpha_i = e^{-\int_0^{T_i} P_i'(t) dt} = e^{-P_i(T_i)}$$
$$\beta_i = \int_0^{T_i} T P_i' e^{-P_i} dt$$

169.37-40 $F_2(0) \equiv \alpha_1$ or $1 - \alpha_2$

This looks like a complicated WON solve I got in X 1990!
 $1 - \alpha_1$ is the probab of success at T_1 .
 β_i is expected sol time at T_i .

On the other hand, T_1 max $\frac{f}{s}$ slope result at 169.37-40 seems quite reasonable. See if I can show it is correct!

WON: [INV: .00 ; OZ: .19] Summary

00:17440: Va Long 169.35-.40;

duration not F_2 : From 168.10-11 some recursion!

$F_2 = \int_0^T p_i(t) e^{-\int_0^t p_i(s) ds} dt$; $F_1 = - \int_0^T p_i(t) e^{-\int_0^t p_i(s) ds} dt = 1 - e^{-\int_0^T p_i(s) ds} = 1 - d_T$

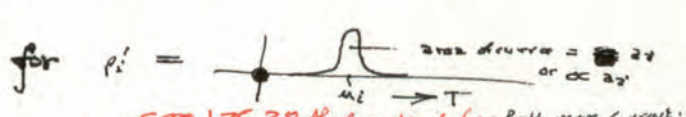
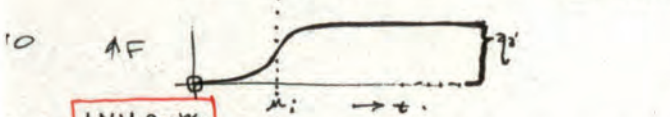
04 Look at (174.36R): $1 - \alpha_i = \delta_i$; $\alpha_i = 1 - \delta_i$

$\frac{\beta_i + (1 - \delta_i)T}{\delta_i} = \frac{\beta_i}{\delta_i} + \frac{1}{\delta_i} + T_i$

I do see a total $174.30R$, but I suspect the reasoning is somehow wrong - Real α_i is more likely to be riter

08 HVR: The Aht. of 167.10-17 gives good reason for $\frac{d_i}{u_i}$ as a v.b. approx for $\frac{d_i}{u_i}$ ordering.

20 $z_i = \int_0^{\infty} p_i e^{-\int_0^t p_i(s) ds} dt$. I think it's a prob of success at time ∞ .



11 for INV probs T. method of 170.03-.07 is provly correct.

See 176.29 for Mod. fn: Probly more correct. It assumes the success of PST's are independent

2 that curves are "f" curves $F(t) = p(t) \exp(-\int_0^t p(s) ds)$; $F = -\exp(-\int_0^t p(s) ds)$

When doing Updating during Srch, After each update, can vary the curves (172.00) Then necessary from to Bin in 172.09 instead of past work on each. - Rem can have work on best bin that has content of 1 or more PST's.

19 For OZ probs: If we do Update during Srch; (173.32-.40, 172.00-.40, 173.00)

20 Choose single PST w. best expected yield for time variable; work on it until update time. This need not be best decision criterion (see 176.12-.20)

21 That Do update then loop to .20 The update should be able to discover "correls" (i.e. when certain PST's do badly, certain others will put worse curves affrd to Rem.)

22 Routine .19-.21 is \approx for INV also, except that faster updates, as we will often jump from one PST to another as we move along the F curve & we jump into a different Bin (172.09).

27 So .08-.27 may be a good start for phase 2 INV & OZ. Could it also be final Phase 2? Soln?

0 For SFunc search in Phase 1 Q.A. 137.00-138.29 (Simpoture. - using Jung's formalism)

Also of much import: "Context" 168.00-.20 is a good recent discuss. Try to find other, older discuss of context, the mechanics of how they work; their inductive meaning (Just how they modify a prod, is just when a prod they modify)

- Another bit Q: How does Q ATM Inv Bernoulli seq. type replys!?

ONE BIG Trouble w. Lsrch is that it doesn't consider give cc credit for eq. v.b. ll codes - which is one way to realize Bern sep. (Lsr's rule) by ALP.

See how to proof that $z^k \geq ALP$ with "constant" works - try to fix Lsrch so it can take advantage of ll codes.

SEE 176.29
-.40
This closes up e.d.f.ty
T. Curve of .04 is fine
is not much distant from the $\frac{d_i}{u_i}$ soln.

Minimum Problem

→ WON

NIPS

Time Varying σ_z :

Refer to previous mention of this problem on ID 825...

425

92.40.55.100
59.7 plastic tray
51925
42.50

WON: Perhaps go back to work on AND nets.

I had a proof/criterion for AND/OR nets. — to show it was a meaningful problem.

02: Time Varying σ_z . Say our Game is $G(x), K(t)$: F is a known function of T .

To solve P_{i2} : we have to $P_i(G, T)$ curves for each P_i for state G , $G(x)$.

Was modifying Russ curves $\hat{P}_i(G, T) = P_i(G \cdot F(T), T)$.

The shapes of the P_i curves will be quite different from the usual $P_i(G, T)$ curves because of G necessarily & eventually w. by T .

For each P_i at each T , we will have different $P(G)$ curves. We want to

select the "Best" curve.

Using convexity, highest expected value of G is not ~~right~~ unless G has been "Linearized". If the D.P. for $P(G)$ is "narrow", expected value is D.H. —

Otherwise, its ~~not clear~~.

Actually, the "linearization of G " is part of the problem domain. One has to be able to tell admittedly how much better one G_1 is than another G_2 .

Otherwise, its impossible to choose between $P(G)$ curves.

So, anyway, one simply picks just P_{i1} and T that have the "Best" $P_i(G, T)$ curve

$= P_i(G \cdot F(T), T)$ curve.

meanwhile, one is updating, & if it is clear that a new i, T should be

chosen, one jumps to the new i, T .

Essentially, the "Linearization" soln. of 15 means that I must be able to decide between "Batteries" $pc = .3, G = 10$ v.s. $pc = .6, G = 5$
or $.6, G = 6$
or $.6, G = 4$

4.7.03: $\int_0^T T G e^{-G} dt = \int_0^T e^{-G} dt - T e^{-G}$

So 174.300

$\frac{\int_0^T e^{-P_i} dt - T e^{-P_i}}{1 - e^{-P_i}}$

$\frac{\int_0^T e^{-P_i} dt}{1 - e^{-P_i}}$

we want Min (this is probably a very old & convex!) & probly of success at time T.

$\int_0^T e^{-P_i} dt = J(T)$
So: $\frac{J(T)}{1 - J'(T)}$

$S(1 - e^{-P}) = T = S \int_0^T e^{-P}$

We want Max

$\frac{1 - e^{-P}}{\int_0^T e^{-P_i} dt} \approx T$ for small T.

So $\approx \frac{1 - e^{-P}}{T}$ for small T.

.36R which is $G \cdot T |$ soln.

T. corresponding to "Convexity w. st. lines" is to find PFA. Smallest T at which $\frac{1 - e^{-P}}{\int_0^T e^{-P_i} dt}$ is max.

Do this for all cards. Pick card for which P_{i2} is max, do it first — Maximize until $G_{i2} = G_{i1}$ at next best card. Then do next best card. Then move along both, so as to keep $G_{i2} = G_{i1}$. When $G_{i2} = G_{i1}$ at 3rd best card, do 3rd best until we reach that max G_{i2} . Then work on all 3 cards on 1. cofc.

$1 - e^{-P} \approx P(T)$ if $P(G)$ is small.

NIPS

01, 8000 bits
8000 bits of symbols.

Exp:

DO Precision of TM in 3 parts ① "overview univ. of its properties"

② Soln of QA problem (Phase 1 soln) 2.5 L search? Do better QA

③ Won solns for ENV, 02, time varying 02

SN T. soln of time varying 02 is also soln for un-time varying 02 \rightarrow $\frac{2}{3}$ \rightarrow (see 175.19 ff)

SN Re: "Incomputable ALP" ~~the~~ Lecture: One can often tell if one Model or Model set is closer to P_m than another. As well as tell how "good" each Model set is for prodn (of next symbol). So perhaps Heuristics.

N.B. that 179.30 result \rightarrow 176.3 result \rightarrow "very old won result" suggests that Algebra is ok. - T. first 2 were obtained distrib. way.

STATE of TM:

1) General mode of Phase 1; Phase 2 operation seems like Good idea such (including updates during)

2) 3 input one for QA problems seems good to 2 state or.
2) Getting it to work. Beon exp. type problems - remains to be done

• 3) How to introduce Context: for L search in QA problem. See 168.00-40 \rightarrow This is Good.

There is much outline work on "Context": I should see to some of it. it does have imp. ideas

• 3) TSQ needs to be written.

\rightarrow 183.00

~~I don't know what this~~ 2 parts needing more work:
Context 2b; TSQ 3: Context may need more work: I'm not sure of state of problem.

TSQ may be ok. I could just use fairly large machine & try to do Algebra text.
For Advanced work, try that Book that RAMANUJIN learned that from.

More immediately! 1) Finish off 10) Lecture

2) Put out next edition of "Report": Add to my revisions: Use Soln., also various ~~the~~ indices to my notes ~~to~~ to expand various parts of "Report".

Re: Solns to won for AND, OR nets:

1) For OR nets! Re: "soln" is used for ENV & 02Z probs.

T. solns used assumes ~~the~~ P's are "INDIP" (certainly not true)

2) can we define a set of tasks that are indep, that somehow are present in dependent tasks?

In general, probably No, because to space of code equivalent problems is a very high dimensional space, & the set of indep tasks seems (as vectors in Hilbert space) seems (like much less nb)

Anyway: some ideas I had on how to estimate representational (179.00)

1) List of solved, unsolved problems + refs to work on those probs
2) find of the list. try to update it.

001

002

10

16

20

30

31

- 0.177
- 1) Clustering Clustering Cands into sets & using "representative Cands" for each set as "index" Cands, w. highly correlated other members of each cluster
 - 2) Cross corr matrix betw Cands.
 - 3) Updating during Simul, If t. Current Cands begins to look bad, update them with new ones. Cands should also be updated. Any way, the idea of using best trial for awhile then doing updating, seemed a use. way to deal w. t. problem, but I'm not certain about it. "best" or even it is correct. Updating again but do not correlate. very well, so it sounds like a very useful idea

0

SN To what extent can I use my current "soln" of t. OR problem to solve "Cure Cancer"? As I have no way to present "soln" to OR postulates index Cands, which act as a reference starting & return not for "Cure Cancer"; The update betw trials (Cand updates) idea helps w. OR; would it help w. "Cure Cancer"? Updating during trials is perhaps more relevant.

20

22

SN For WON problem: for OR nets; Criterion of soln is:
 a) T. technique will solve ~~problem~~ net fraction of time, F . $F \rightarrow$ maximized
 b) Opt. time betw soln. is obtained, $\frac{E}{\text{Expected}}$ time to soln. is Min.
 My soln. satis find key name.
 For Won singular, less 2 criteria can be used for soln., is soln.
 That soln. always exists is by using two to two OR Cands, is a theorem.
 Thus proby of soln's value is easy to compute: forward tasks, have will be a pct of completion at $t=0$. From these values, it is easy to compute pct of soln at $T=20$ for an hour. \rightarrow 180.22

30

Re: OR problems (2 tasks, say). If $P_i(t)$ is density pct of soln at time t (w.o. condition that P_{12} is smallest T for soln.) then if we do 2 tasks in order, is new P function will be (if task 1 is first) $P_1(t) + P_2(t-T_1)$ — T_1 being end time for task 1. However, I have no picture of how $\geq P_i(t)$ could occur if i.e. how can a poss. soln. occur at > 1 point? — It could for separate Cands, (tasks), but not for a single task! — Unless a single task is, say, composed of independent (or even dependent) sub tasks. (no, independent can't depend on success or failure of previous task)

35

So t. proby funct for a single task will always be one for first time soln. so $\int_0^\infty P_i(t) dt \leq 1$.
 For 2 tasks, we have normal $\sum P_i$ (first time pct) for tasks followed by $(1 - \int_0^\infty P_1)$
 $P_1(t) \cdot P_2(t - T_1)$. doing tasks in order, gives $\int T P(t) + \int (T + T_1) P_2(t) \rightarrow \begin{cases} 179.00 \\ 178.11 \end{cases}$
 $\equiv \int T P_1(t) + P_2(t) \cdot T_1 \cdot P_2(t) + P_1(T_1) \cdot \int T P_2(t)$
 $\equiv \int T P_1 + \int T P_2 + T_1 P_2 - \int P_1 \cdot \int T P_2 + T_1 P_2$

20:178.40 :
$$S^T P_1 + (1 - S^T P_1) S(T + T_1) P_2 = S^T P_1 + T_1 S P_2 + S T P_2 - S P_1 \cdot S T P_2 - T_1 S P_1 S P_2$$

$$S^T P_1 + S T P_2 \quad | \quad + T_1 S P_2 - T_1 S P_1 S P_2 + S T P_1 - S P_1 \cdot S T P_2$$

$$T_1 S P_2 = S P_1 \cdot S T P_2 - T_1 S P_1 S P_2$$

$$T_1 (S P_2 - S P_1 S P_2) = S P_1 \cdot S T P_2$$

$$= T_1 S P_2 (1 - S P_1) = S P_1 \cdot S T P_2$$

$$= S P_2 [T_1 (1 - S P_1) - S P_1]$$

$$S P_1 = A_1 \quad | \quad S P_2 = A_2 \quad S T P_1 = B_1 \quad S T P_2 = B_2 \quad | \quad \int_0^T T P_1 dt + \int_0^T S P_1 = T P_1 |_0^T$$

$$B_1 + S A_1 = T_1 P C T_1$$

11:179.35
$$S P_1 + (1 - A_1) (B_2 + T_1 A_2) = B_1 + (1 - A_1) (B_2 + T_1 A_2)$$

$$B_1 + B_2 + T_1 A_2 - A_1 B_2 - T_1 A_1 A_2 = T_1 A_2 (1 - A_1) - A_1 B_2$$

v.s.
$$B_2 + B_1 + T_2 A_1 - A_2 B_1 - T_2 A_1 A_2$$

$$T_2 A_1 (1 - A_2) - A_2 B_1$$

$$A_2 (T_1 (1 - A_1) + B_1)$$

$$T_1 A_2 (1 - A_1) + A_2 B_1$$

$$T_2 A_1 (1 - A_2) + A_1 B_2$$

$$A_1 (T_2 (1 - A_2) + B_2)$$

$$= \frac{T_1 (1 + P_1)}{P_1}$$

$$+ \frac{A_1}{P_1}$$

$$P_1^3 = \frac{2 \pi}{\text{month}}$$

$$\frac{T_1 (1 - A_1) + B_1}{A_1} \stackrel{\text{Gorc.}}{=} \left[\frac{T_1 (1 - \int_0^T P_1(x) dx) + \int_0^T T P_1(x) dx}{\int_0^T P(x) dx} \right]$$

$$= \frac{T_1}{A_1} - T_1 + \frac{B_1}{A_1} = \frac{T_1 + B_1}{A_1} - T_1$$

$$S P_1 = T P = S P + T S T P$$

$$T P = S P + T P$$

$$\frac{T_1 + B_1}{A_1} - T_1 = \frac{T_1 + \int_0^T (T + T_1) P_2}{S P} \stackrel{\text{Gorc.}}{=} \frac{T_1 + \int_0^T P_1}{P_1} = \frac{T_1 (1 + P_1)}{P_1} + \frac{A_1}{P_1}$$

20:
$$S P_1 = -e^{-p} \Big|_0^T = 1 - e^{-p(T)}$$

$$S T (-e^{-p})' + \int (-e^{-p}) = -T_1 e^{-p}$$

$$+ \int T P' e^{-p} = S e^{-p} = \frac{1}{2} T_1 e^{-p}$$

$$S T P' e^{-p} + T e^{-p} = S e^{-p}$$

$$\frac{T_1 (1 - S P e^{-p}) + S T P e^{-p}}{S P e^{-p}} = \frac{T_1 e^{-p} + S T P e^{-p}}{1 - e^{-p(T)}}$$

$$= \int T_2 (e^{-p})'$$

$$\frac{S e^{-p}}{1 - e^{-p}} = \frac{1}{S e^{-p} - 1}$$

$$\frac{1}{S e^{-p} - 1} = -1$$

$$\frac{1}{S e^{-p} - 1} = -1$$

$$\text{Which is } \approx 176.30 \text{ R}$$

$$\frac{T_1 + S T P + T S P}{S P} \quad | \quad S T P = S P + S T P'$$



00: Res: No use of WON to replace LSrch in parts of T.M. : This seems like a very implausible decision. While WON is in a sense "optimum", it is only as good as: $P_i(T)$ curves.

Similarly LSrch is only as good as the optimal. But guides it.

LSrch seems to have the advantage in that we can visualize an upper bound (at least) on best of discovering a "Best" soln (\equiv down cost of that soln)

For WON: we will try PST's in some reasonable order (which is what LSrch does)

Actually: from WON 23.00 (4.11.03) $F_n = P_n \left(\frac{I}{P_n} \right)$ ($P_n \equiv$ pc of PST_n) $G(x) = 1 - e^{-x}$

If we LSrch to be optimum: we can easily get F_n if we assume $\alpha = \alpha P_n$ or $\alpha P_n = \alpha P_n$ and the pc of F_n curves are $\alpha = P_n$; $\alpha = P_n$ or αP_n

0: Turn out that we don't have to choose between WON & LSrch!

From 179.26:30: Define $F(t) = e^{-P(t)}$

Then, "Updating" how to determine exact purchase for $F_n(t)$ (for particular problem/PST_n)

It can use some standard form for F_n is determined subject "increments" or α, μ, β, z .

for the desired F_n . If we used an "expressible" form of $G(\alpha, \mu, \beta, z, T)$,

Then, if the update uses $\alpha = P_n$; $\mu = P_n$ (G is distorted in both G & T directions, expand)

Then (α, z is not relevant) LSrch is optimum; Time spent on PST_n is $\propto P_n$.

So Essentially, "Updating" can determine if it wants to use LSrch or Not - LSrch is a

- Special case of WON (OR) so WON will automatically use LSrch if
 - Update is successful
 - LSrch is Best
 - We tell update to P_n so $\alpha = P_n$
- We DO have to tell update pc of PST_n is, hvr. $\rightarrow 188.17$

22:18 SN Rigorous analysis of "OR" Nets: AND/OR nets.

In the case of OR tasks problems: say we have 2 tasks, w. assoc F_i curves.

If $F_1(t)$ is prob of soln by $t=0$, then $(1 - F_1(t))(1 - F_2(t))$ is prob of any soln. by $t=0$

This is indep of order in which we do trials.

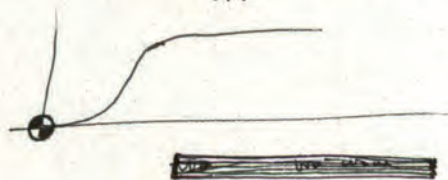
If we do tasks first) then tasks containing prob, then of t a fraction of successful trials,

we will have an expected soln. time. We want to ordering of tasks that minimizes P_n is expected time.

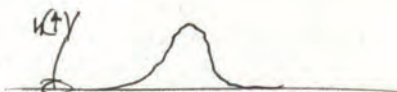
This is what was (3) done in 169.35 - Plan more accurate by 174.30; Plan simplified in 179.30

31 For any AND/OR net of finite size, having arbitrarily complex serial, || (sub)tasks; we can easily compute the prob of soln of any and set, to any or set, a computer search of combinations to get prob of soln. of entire net at $t=0$. One definition of pc of soln: one spends 1 msec on each task, until task is solved. A certain fraction of maximum trial soln will solve it. This α is indep of order of trials. For each algorithm that orders trials on tasks & parts of tasks, t, α will be the same but on the successful run, we will have a unique value of expected completion time. We want to ordering \rightarrow that \rightarrow time is min. \rightarrow Maybe NO! $\rightarrow 181.00$

00 (180.90): A possibl. counter Arg. to 180.31-40: That very long solns can occur & make "Expected completion time" (ISPs is possibl?).



T. dens. by d solns is



we want min $\frac{\int_0^{\infty} T h(t) dt}{\int_0^{\infty} h(t) dt} = \text{[shaded box]}$

Well, we can ~~prob~~ have distributions $h(t) \rightarrow$

08 \rightarrow zeroth moment exists, but not 1st moment (i.e. first moment = ∞) e.g. $\int \frac{dx}{(1+x^2)^2} > \infty$ $\int \frac{dx}{(1+x^2)} = \infty$

(Also zeroth & first may exist, but not 2nd moment) \leftarrow R. P. is not a real house.

For t . Gauss. & d.f. of all moments exist \rightarrow But note 181.13-14

3 If we ask T on to optimize 1st first 3 moments for our ~~same~~ date set, it will very probably be able to find good finite values.

13 If $h(t)$ is of form $P' e^{-P't}$ can ~~it~~ (0.08) be true?

5 say P' is bounded, and ≥ 0 . This \uparrow function would seem to have all moments!
 i.e. $\int_0^{\infty} t^n P'(t) e^{-P't} dt = \int_0^{\infty} t^n P'(t) dt$ is bounded for all $P'(t)$. | $\int_0^{\infty} P'(t) dt$ would not converge!

17:180.21: On Lsrch v.s. WON for ~~INV~~ problems! Note also that P_1 soln. satisfies GHT1 for small T . (176.36R) i.e. it's max $\frac{P_1}{T}$ ordering for small T 's.


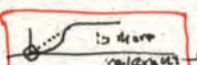
10 On Lsrch v.s. WON for OZ probs! The won soln. is 175.19-23. Also Note 171.32-40 172.00-40 173.00

Will this ever be \approx or \equiv Lsrch soln? — which spends time & space on each OZ PST_i or PC_i .
 Won picks "Best" PST_i & works on it until "update time" we may switch PST_i 's after update time: (The "T remaining" will change after each update).

Actually, for both INV & OZ probs, I should be comparing Lsrch & WON after update has occurred, when I have a $P_i(t)$ or $P_i(t, \epsilon)$ for each PST_i .

In Pz case, Lsrch gives $\alpha = \int_0^{\infty} P'$ ordering = (mean = zero moment) or $\frac{1}{2} = \left(\frac{\int_0^{\infty} T P'}{\int_0^{\infty} P'} \right)$ ordering = first moment / zero moment.

While an ordering for Lsrch is about rate (I used it for pre colling of "Good" PST_i)
 $\frac{1}{2}$ simply orders PST_i 's — it doesn't offset by PC_i 's to Perm — as are needed in Lsrch.

For ~~the~~ F' of form , $\frac{1}{2}$ is initial slope for Won ordering.
 Maybe $F = S P' =$  is done independently — But it amounts to same thing.
 There is a way to order probabilistically by using only $[x_i, y_i]$ info. (Detailed analysis within ~~the~~ last few pages)

Of course $\frac{1}{2}$ is optimum order for GHT1. It is exact if $\alpha^2 \geq 0$.

I think this is order used by won if T is small (176.36R)
 So for INV probs this should be comparable or better than Lsrch.
 My impression is that it was known = 2 param f or f' dist. giving $\alpha \geq 0$, this makes ideal GHT2 soln possible. — a Lsrch is ideal GHT2 soln also — but $\frac{1}{2}$ order may be better.

NIPS

concurrent simultaneous

10:181.40: So WCN such ~~with~~ with // Update looks fairly good & is likely to be better than Lsck - It is certainly close to it.

For OZ probs, WCN smk seems v.g. - probably ~~usually~~ better than Lsck, because WCN tends to spend more time on very promising PST's - unlike Lsck distributes CB over many PST's - even re ~~over~~ only over re Part to be v.g. [see 175.19-21 for version of OZ]

Other remarks on WCN for OZ: 175.22-23, 176.12-17 (on/information of G)

The discussion of "Corrains" looks good for INV, is relevant to OZ as well (177.31-178.03) One pos advantage of Lsck: that by trying many different candidates, Update has more data to work on

For my version of "T. report" I want to write a bibliography of WCN for PR, AND & OZ.

Also, ~~the~~ actual expository writing out: main ideas & arguments as to why it ought to work - why it may not work & ideas on "how to fix it".

SN δ functions: 3 realization methods: 1) 3 input vms 2) Lays rule, like AZH1 3) The "Quote" expressions in AZH1 that express δ funcs can be described by a Lispish (functional) Lang. [This (st) is not quite clear! - It may well be that Lisp automatically does this: It could be a string (= pfm) that executes that string (appm).

SN In OZ trials I never ^{almost} have "failures": A PST will usually get some G: 1'

Occasionally, however, there is not enough time for success trial: $G = -\infty$

Thus $-\infty$ can be one of the G values: Update can use this into for long.

I'm not yet sure how the probabilized should be mixed in densities of G


i.e. pc of $G = -\infty$ v.s. pc of $G = g \pm d$. Actually, we have the same problem in

INV problems! (?) No: At $cc = T$ the PST either succeeds, or fails to succeed.

Its δ G-D function $O^j(\cdot)$. Actually, the O^j function could look at

The situation in OZ: INV ~~can~~ do correspond. - In INV, the $O^j(T)$ curve

gives densities for success at T, & interpreted densely ~~prob~~ values for the success

at times $> T$, & cut off. 

probability of ~~no~~ no continuation up to T_0 is $1 - \int_0^{T_0} h(t) dt$: This assumes $h(t)$ implies no soln before t.

Similarly $O^j(T, G)$ give \rightarrow prob densities for $G \pm d$, and a prob that $G = \infty$ for T.

(We multiply proby densities by actual probabilities. Do we want to accord

same (logarithmic) wt. to both of them? \rightarrow (35)

NIPS 18.33-36 Critiques & updates optimization of O^j on basis of SOP: here, see ibid [NIPS 18.36 + 19.06]

35 \rightarrow Its 0.9! say we had n_d proby densities & n_p probys for our corpus. for δ density with Δ , the actual proby of model is $\left(\prod_{i=1}^{n_d} P_d^i \right) \cdot \Delta^{n_d} \cdot \prod_{i=1}^{n_p} P_p^i \equiv \alpha$:

α is a regular probability. What we want to Maximize is $\frac{\alpha}{\Delta^{n_d}}$. The $\frac{1}{\Delta^{n_d}}$ factor is indep of what models we try, so its not relevant to the maximization process. So its o.k. - we are not "adding apples to oranges".

Nips

177.16: "Context": In induction problems, we always have ^{have} to same a ~~script~~ ^{script}, no matter how many examples we have worked in the past. If we had an of time, Lisch would be the best way to go, ~~to code for the entire corpus.~~ ^{to code for the entire corpus.} Here, we are always short of time.
 The best way to code a corpus in very short time is to identify codes. (It's any large code, here!)
 Another trick is to use codes (with recursive functions) ~~to~~ w. to trap, or wish they occurred in the ~~past~~ ^{past} coded corpus — which is the usual of AZ141 way.
 There are other tricks, but they all depend on ideas about saving time & helping to a too small CB.

How is "Context" relevant to the above? Well, in various past "contexts" certain coding tricks have been "useful."
 How does this idea accord w. J's ~~use of~~ ^{use of} "successful" (?) use of "Context" via the "Boosting" instruction?

LN LZ compression uses OSL! (One Shot Loop)

Manly, AZ141 uses ideas that define that have been used in the past, are likely to be used in the future, as well. Are all "regys" expressible as "definitions"?

It is my impressn. that all types of regys that one observes/ invents, are obtained initially by a Lisch (or a process w. a same efficiency as Lisch f. w. Pen a factor of 2x depth of concept). So a rather simple scheme for inductive coding ought to work optimally! i.e. for each new problem, search a fairly large amt. ("oversearch" is meaningless at this pt., since the CB is part of the problem & con. ... T. problem of what CB to use for induction — a updating, is unclear: The 50% soln. is making full use of updating ~~inv. probs~~ Inv. probs, but it's not clear what CB to give for 02 problems originally — a for (non-update) 02 problems — in which case it's 50% soln is meaningless (or useless).)

I suspect that we will need "hyper level frng" to decide how much to spend on updating 02 problems (in which the CB is not given). Chess, w. 40 moves in this, total, is a mild kind of problem of this sort: CB for each search has to be estimated by TM so as to optimize some overall score.

In .13-.21, it would seem that "Context" is irrelevant. — yet w/o ^(context) context considerations,

we seem to have a serious scaling problem — that p.c.'s of concs ↓ as corpus grows. It may be that we intuitively in context do not to overcome the scaling problem.

One ditto in my thinking: I think of a "new problem" as involving nothing more than finding some new concs, that satisfy a certain constraint (possibly a "gray" search constraint) — an "02" problem, (like).

Mips

BACKTRACKING: How to do it! .25-.30

When to do it, ~~in a particular case~~ - .31 ff
How much time to spend on Backtracking: 185.15-.25

00: 183.40: Evr. QA problems is f. Update problem don't seem to be like that - in fact we are looking for a global soln. to "all problems together".

One nice look at it in Reg way: one has a sequence of problems: P_1, P_2, \dots problems to find a good model ($\approx O^d$) for $[Q_1, A_2]_{i=1}^k$. TM uses its ~~past~~ experience on the first k problems to solve $k+1$. - This is \approx what "OOPS" does.

Its also what I wrote about in trying to find a good function that goes from input of all past k problems & their solns, to t a proby distribution on the solution of t problem (Q_{t+1}) and its soln.

10

In General on "Context": When one discovers/defines a new concept, it is very to know Under what circumstances this conc is usable. This is "context" & it is as important as knowing/inventing/discovering the conc. itself. When a conc. is first defined/discovered, its SSZ is small & so is SSZ for its context is also small - perhaps \rightarrow see (183.14: LZ coding!)

OSL will have to be used, initially, for emulating/exercising/trying the contexts that occur w. the newly defined conc. Also, we want negative data for each conc., when it "fit" but was inappropriate: Fortunately, my present methods can work ~~with~~ when only positive data is available

There are a standard set of context types that will be initially explored!
Note 185.26!

.10 ft may be a difrent approach/data of "context" than my previous approach.

It seems to have a difrent theoretical background. ~~At least~~ We only try to code the corpus, using concs that have been useful in the past, on this course. It is a kind of W.d. Backtracking, one can get badly stuck!

20

"Sequential coding!" (little if any "BackTrack"). \leftarrow This is a serious criticism!

A poss. way to do "Backtracking": Say we are doing OSL & we look for things \leftarrow I don't see how OSL could be Modified to help

24

In past text would give probn of / recast past...

25

I think I know how to do Backtracking: One normally has ~~several~~ several (alternatives ("branches")) codes for the corpus - using difrent data. To continue the corpus, one will choose one ~~the~~ branch & try to continue, carry its data & pc's of concs. If things don't work out well in continuing on that branch, one does "back out fork" & chooses ^{initially} a lower pc branch to continue.

30

.31

While .25-.30 does tell how to do backtracking, it doesn't tell under what circumstances to do it! In the Sci Community ~~Backtracking~~ ("heavy revision") is done infrequently - Do it is occasionally necessary - It is usually a "big deal" & causes much fighting tearing of hair, ~~and~~ & human attacks etc. The idea is, that the predictions don't seem to be as accurate as we expect but they can be. What is our Criterion for this? In physics, we expect probn to be "exact" within precision of instruments, & we continue to \uparrow accuracy of instruments. In Chemistry, we

often can't control things exactly, in which cases we settle for less precision
In biochem, & biology, we know that we don't know even all of the conditions that affected
an experiment. & we expect even less precision

In Biological experiments, we have relatively little control in the experimental results
and not of much accuracy & often are not reproducible (except at low precision)

As one ~~under~~ begins to understand the "Domain" one begins to know what
the uncontrollable/unknown variables are - & estimate how important they are.
So as to get better idea as to how accurate one's predns should be!

e.g. say ~~Temp~~ ^{temp} is a variable that we couldn't formerly control or measure

So it caused much variance in chem/phys results. Then we are able, not
to measure it, but to be able to estimate when 2 temps were about the same.
So we now expect 2 experiments w/ same Temp to give same results -
or to both be in accord w/ theory!

**In General, TM will learn when to expect by accuracy &
when to expect low accuracy in predn. This will be by looking at past
& seeing under what cond, by accuracy was obtained. This ~~prediction of~~
Expectation of my accuracy will be wrong many times, but it will
be enuf to tell us how much time to spend on Backtracking. → Note 189.00**

Is .15-.20 an adequate soln. to the problem?

T. backtrack problem is: Should I spend time T looking for extensions
of old code or Backtrack? i.e. - which is more efficient use of
my time, T?

(Concept) well, if I forget, I'm kind of full of what I used to do whenever I discover
novelty I would spend a lot of time trying to see just what situations in the future,
that conc would be relevant. (like 184.10 → 2.20) → Also Note 186.12-.15!

Also I tried to see if past conc. could simplify (∴ shorter coding) problems
of the past - problems that I had not used as part of the snippets that gave conc.
(was "solved by")

This suggests that normally, I do not use my entire corpus (certainly true!
Particularly note This is the "Encyclopedia Problem"!) So perhaps, in the Encyc.
I begin to look for more cases in the past. Hvr, my not using the complete corpus
for every predn. could be because it was too time consuming. I suspect this
for any predn. I usually select out what I think is the relevant data of the past,
& make predns. on past basis. After I find a good conc. for the latest
predn., I think its structure may suggest other areas in the corpus to look

4. 18.03
Nips

116-35

How ↓ of ~~var~~ corresponds directly to ↑ in mean density of process
is ↑ of pc of symbols that occur

This is discussed in revisions to ~~the~~ talk

DO: 185.90 for referent data that conco might help code.

{ 183.15 to 186.00 } → Does this give fairly complete approach to TM - or at least to "Phase 1"?

T. idea of "Sequential coding w/ back track" is attractive.

→ Any way, try to write this up as exactly as possibl. : I don't want to forget it.
Also try writing up a naive TM Phase 1 in as much detail as possibl. so I can discover main project bottlenecks

FN One way to deal w/ pc v.s. 2ⁿ disparity in Leach! If stochastic systems are properly descrbd, it may well be that they are about the same. Using tokens of various pcs as in AZ141, may help... but can it solve the problem (completely)?

12: 185.25 Re: 185.26 ff: essentially what I did was create a set of obs for "contexts" in which conco would be likely to be usable. These I'd have one or more abstracted conc I discovered (or for all tokens)

th. system would then be a bit like $S_{i+1} = S_i + \epsilon_i$ "path dependent"

4.19.03 On ↓ of variance v.s. ↑ mean boost of next ~~token~~

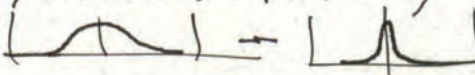
If we assume (for varc problem) a uniform spread of t . Variable being predicted, \leftarrow No! Expt. value of pc of next symbol \uparrow by factor of k .
Then ↓ of σ by a factor of k (i.e. $\sigma \rightarrow \frac{\sigma}{k}$)
Two things of interest: ① The "uniform spread" doesn't mean that its uniform from $-\infty$ to $+\infty$. It does mean that th. spread ~~is~~ doesn't vary much in the "region of interest" → 187.20 for exact treatment

② That we use Gaussian coding for $[X_i]_{i=1}^n$ $X_i = \bar{X} + \Delta_i$ ($\Delta_i \in [k_i - \bar{X}]$)

so pc of sequence $\prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp(-\frac{\Delta_i^2}{2\sigma^2}) = (\frac{1}{\sigma \sqrt{2\pi}})^n \left(\exp(-\frac{\sum \Delta_i^2}{2\sigma^2}) \right)^n$
 $= (\frac{1}{\sigma \sqrt{2\pi}})^n \left(\exp(-\frac{n\sigma^2}{2\sigma^2}) \right) = \left(\frac{1}{\sigma \sqrt{2\pi e}} \right)^n$ since $\sum \Delta_i^2 = n\sigma^2$.

so expected value w/ ~~prob~~ of pc of next symbol is $\propto \frac{1}{\sigma}$

NB The demo of 24-27 ~~does~~ not depend on X_i having a Gaussian Df. i.e. it ~~will~~ probably work if σ is the varc of any d.f. that has a varc.

First assume 21-23 then, if we narrow t. d.f. by a factor of k , the probab density at the "usual" data pt. will \uparrow by a factor of k .
But it follows that $\int_{-\infty}^{\infty} f(x) dx = 1$ & density must \uparrow by $\frac{1}{k}$.  → 188.20

Re: "Context" : I think an important part of my thinking about TM involved my expectations of what future problems would be. It was impt to have idea as to what I was aiming at, when I acquired my early education. This was done by "Looking" at the unsolved problems of Science from a "popular" pt. of view. → 187.005a

so: (Spec.) 186.40: 186.36 - to was, I think, ~~over~~ imp. in RTM, since I was mainly interested in "future world". (Horizon Logic). Since I'm not working on RTM just now, this may be irrelevant to the present problem, But it's an imp. part of the way a Scientist normally thinks. More a part of early education, than later... later. Scientist really knows many of the main problems in his field of expertise.

When a Scientist proposes a radically new theory: To a group of fellow concs us do, he shows how these concs can be used to better code data of the past.... This may or may not involve backtracking. For other scientists, it probably involves backtracking. For this scientist may have obtained the "New View" w.o. backtracking - i.e. he may never have really accepted the current view that he wants to (demolish / discard).

12
13
1. main things that QATM will Not do: (??)

- 1) Solve "Cure Cancer": Tho I think it could be taught to solve probs like that by touching first h₁, then h₂, etc. (Cure Cancer is a RTM-type problem)
- 2) RTM (see 13) \Rightarrow (13) \rightarrow [Any Large ^{prob} problem can get TM to act like RTM]
- 3) Propose & do "Experiments" (An "experiment" is cc used to gain info - not directly related to problem soln.) \rightarrow (t "info. is to be used "later")

Propose that we design a more accurate instruction
Suggest a new region of universe to look at.

14
20
Q: What about EncycTM! This is an imp. kind of TM problem. T. "Encyc" can be RW, in which case, TM would do "experiments". Even w/o. Encyc being RW. If Encyc is Text, TM's indexing of t. Encyc would be used to ~~some~~ ^{precision} \uparrow ~~value~~ of future predns. So this would be an "Experiment".

On the other hand, if one has a large OZ problem, part of it a v.g. PST Cooks very well be "purely m-p gathering" - only indirectly related to t. OZ.

So I'm not at all sure about 12 ff. E.g. to "Cure Cancer", if TM really understands all the problem was, it may very well find a PST that would look like the "Cure Cancer" routine that I've outlined. "Cure Cancer" really is a time limited OZ problem - but TM must have access to RW. T. criterion of "Cure" is that when TM accesses RW news or Google search; It finds that cancer has been cured - that ways have been found to get almost all cancer patients to get well.

QATM, INV, OZ, Timasovs predn, Bay induction are all ~~special~~ special modes of behavior of TM: using a common GPD. Is ~~Encyc~~ Encyc TM an other or these "special modes"? - if so, just how does it work?
I did write a fair bit on Encyc.TM. : well that is this large data base (Encyc); that TM has access to. ~~There is a way~~ Say we want to do prediction, or seq. of α (is that in Encyc. busy)
One way is to find short codes for α , using access to Encyc.

10:00: One way may have a seq. of QAs is a final Q to answer.
 TM has to learn how to use the Encyc. — so that TM doesn't have to read the whole Encyc to find useful info. TM could answer Q's by giving addresses in Encyc that are relevant & then doing some operations on these references.
 We'd probably have to use some kind of TSO to teach TM how to use Encyc.

0.05 → Find previous work on this Encyc idea. ... I don't want to repeat old work.
 So: I want to find (at least) 2 things
 1) Discussions of Encyc problems - (187.19 - 189.06) — ID 585.21, 36 (V.G.)!
 2) Discussions of context. (183.15 - 184.09) → 186.12-15
 (184.10 - 186.00) seems like a very good discussion of context.
 The earlier papers may be also useful.

0.11 → (Encyc) cont.: In one version of the problem: the "A" is always a direct quote from the Encyc text.
 Or, it may quote a part of Encyc.: At a higher level, TM will use the Encyc along w. other info — perhaps statistical info from Encyc & other sources. At the most general level, it will be solving OZ, Inv, QA ... problems: use roots to Encyc as part of the code for a soln.

0.2 → For finding older work on Encyc: (≡ Enc): Alternatively! Deduce the problem carefully; this will suggest soln. (Note that "indexing" of text can be part of "updating time".)
 Originally the Enc problem arose when there was large corpus & a problem that was reduced to a small part of it. In theory, it could code the whole corpus to solve the problem, but this would be wasteful of cc. (It would be OK, if $CB = \infty$, hwr @)

20: 186.35: Re 186.21-23 We can assume uniform spread from $-R$ to $+R$ in this case.
 For $R \gg G$ (update's) the results will not depend much on R , and results are indep of R vs $R \rightarrow \infty$. This argument about $R \rightarrow \infty$ works O.K. for this particular problem — but not merely for many other problems!

0.27 → [SN] In using WOI w. ll updating for both INV & OZ problems, T. System doesn't try ^{much} variety of PST's: T. result is that update to basic data for update is poor & updating itself is poor.
 To remedy this: Spend $\frac{1}{2}$ of time ^{on WOI} $\frac{1}{2}$ on updates, $\frac{1}{2}$ on regular Lush trials

0. → [SN] [NB] Re context codes are part of code of corpus. They modify pc's of other codes. They are what J. talks about when he says the system compares the pc's of the next token.
 If we allow these context codes to have adequate arguments (like "type of problem", & even "frequency" ...) then they will be general context.
 → Is there any way that context correlates to cc? — essentially a non-"pc of corpus" (type of consideration?)

10: 193.40 = A poss. formulizn! We have a function from \vec{X} to $(0,1)$ interval.
 01 $F(\vec{r}, \vec{X}_i)$: We want $\vec{r} \rightarrow \sum_{i=1}^n F(\vec{r}, \vec{X}_i) = \text{max}$, i.e. we select such $\vec{r} = \vec{r}_0$.
 (Say, to start with \vec{r}_0 is not a "wall" of its parameter space)

$\exp(F) = \text{probly.}$ So \vec{r}_0 express maximizes "likely hood" \equiv PC.
 There is a "True" value of $\vec{r} = \vec{r}_M$: that generated the data probabilisticly.

Or another way to look at it: A generator n probabilities, P_i : We montecarlo, use \vec{r}_M to generate a binary string of n bits, $= \vec{X}_M$.

We then test $P(\vec{r}, i) |_{i=1}^n$ function of $i=1 \dots n$ and try to find $\vec{r} \rightarrow$
 10 $\prod P(\vec{r}, i, X_i^i)$ is max i is an input config to P X_i^i is i 's component of \vec{X}_M
 $P(\vec{r}, i, X_i^i)$ is i 's probly assigned to X_i^i w. "input" $\leftarrow (to P)$
 12 So this has a probly value for each i . How much does it differ from

I'm not doing well in this problem! Define it is put it on "Stack" (IID 6.11.38)

Discussion starts on 192.20, goes to 194.12.

T. idea starts with Akaike factor $(86.16 - .35; 187.20)$ $G^2 \rightarrow G^2 \frac{n+k}{n-k}$
 18 186.16 is a vague outline of a proof, probably it would be well to make a careful proof of it before going on to the previous problem, which is \rightarrow We have an n bit seq \vec{X} generated by p.d. M .

We have a function of the previous prefix, that gives a best PC for each bit. This function has been optimized over with those k params, so the PC of the corpus is max.

21 What is the expected value of β_i (in P of i 's next symbol)? "Is it an unbiased estimate?"
 \rightarrow do we want $P \ln P + \beta + \ln P$? $= (P \ln P + (1-P) \ln(1-P))$ function of β alone.

\rightarrow A (new) poss. approach: Consider a continuous (Alpha) case is divided into d discrete levels. See how the "fuzzy" version works, then \rightarrow 2 levels ($\equiv 0,1$).

earlier work on this problem 152.24 - 153.19; 156.32 - to 157.00 - to 158.00 - to 159.00 ft.

193.37 N.B. 193.37 is probly wrong: $\sigma_{obs} \rightarrow \sigma_{obs} \sqrt{\frac{n+k}{n-k}}$: parity of that symbol is $\frac{1}{\text{obs}}$
 $\sigma_{max} \approx \sigma_{obs} \left(\frac{n+k}{n-k} \right) \approx \sigma_{obs} \frac{n+k}{k}$

Backtracking "is perhaps a special case of Bayesian (trying) evaluating a new Model."

0: 194.09 152.00 - 12 is also a useful direct copy
 also: 150.07 ft looks like good copy! 151.00 ft, Trouble is, I have to original \equiv "talk as written",
 150.07 - 192.12 (it's discussion of 195.10-196.13 on "future probabilities".) - But they don't fit together!

150.07 - 26 looks v.g.! : So perhaps do this: for each section I've written, write a short summary of its ideas (or a "Listing" of them). Then see how I could fit them together or make a paper from these parts.

30:

SUMACS: ID 378.00 Derive a Sumac Schema, but felt it wouldn't work because eventually, the initially finite no. of O_i functions would go to 0 as the no. of QAs ↑. If we do "backtracking", this will not occur.

Note also, that "backtracking" a complex corpus need not be so difficult. We can always do "soft coding", in which we make a lossy code for the corpus & then add the corrections.

A main difficulty here is that this does not give codes in the "pc order" — which is of importance in using the "Lsrch" to look for induction codes (i.e. O_i's).

0:

On assigning pc's to New reactor, Nat defense, Genetic catastrophe, ^{space shuttle} etc. Is this related to the computability of UPD? While the reason for the real incomputability of .10 are not identical to those for the incomputability of UPD, (that all models have not been fully evaluated... many models are p.r.), .10 is not the same as UPD.

In .10 we devise heuristic models w/ pc's that we may or may not have enough time to evaluate.

~~But~~ these seem to be different from the kind of things being Models for UPD. In UPD, we are trying to get good models (i.e. codes) for a ~~known~~ known (i.e. slightly unknown) corpus.

In .10, we take "known" (i.e. by pc.) scientific laws & try to use them for extrapolation.

18

One way to think about .10: In .10 we have no empirical data, so we are working only on the a priori. This will be to a priori on sequences of RW events that may or may not lead to the failure whose pc we want to evaluate. In the case of "Reactor failure", we do have empirical pc's of various event types that help evaluate our needed a priori. These empirical pc's may have been obtained from this or from other reactors.

40

That it may be necessary to expand to detail in .10, I think that it is essentially what we do when we use UPD for prodn. In the R.W., ~~the~~ there will be many possible

26

codes to evolve. The set of eligible codes is not altogether clearly defined.

or lack of sci knowledge!
e.g. Bad Sleep Cycles on People Monitoring Reactor

Sometimes a Reactor "eligible code" are not evaluated, for practical reasons. The actual failures of reactors have been by far the most have been by "ineligible codes" (codes not considered).

A code can be ineligible because (a) Very low expected a priori. (b) Politically impossible.

(c) takes too long to evaluate (d) Simply wasn't part of (related to) ~~the~~ the gen set different from (e).

Common cause of unexpected failures: Correlation of events not considered to be correlated.

10

So, it seems clear that .10 is a case of UPD, but that the "instruction set" or the rules for creating codes & assigning pc's to them may not be clearly defined — also, (26)

There is another way to evaluate pc of reactor failures: We have this large corpus of reactor histories — not many failures, but many malfunctions & partial failures, etc. from this corpus

NIPs.....

20 : 195.40 : An interesting Q: Int. problems of 195.10, The incomputability of the UPD is very relevant. —
Yat it is not in most production! — Why not? How do these situations differ?

Well, in 195.10; we either know no previous corpus on which to test our models or y. corpus is too small, so SSZ error is $> \epsilon$. upper bound on permissible failure probability

Hrrr, f. forgg'z abt doesn't consider P_M approaching $P_M \in P_M$ or $P_M \rightarrow M$. at all!

Maybe we consider $P_M^T \rightarrow P_M$, and that P_M is "as good as we can hope to get".

Actually, in these failure cases, the problem is always to show that $\|P_C\| \leq \epsilon$, and this is necessarily impossible. We hope to show that there are no codes w. pc's of $\leq \epsilon$.

These ones we find must be (usually) $\ll \epsilon$.

This is opposite to usual problem of prediction, in which we want to find codes of $\leq \epsilon$ pc as possible. (Well's not same thing in 195.10, but want to show this "Max $\leq \epsilon$ "

13 is $\leq \epsilon$.)

That discussion of 195.10 - 195.13 is good & useful, I still feel that I don't have complete understanding of what's going on. T. predn. for ϵ -items of 195.10 seem to be quite different from

Normal predn using UPD.

20 : 194.90; T. part es of new 193.26 SSZ uncertainty \approx 3.31 Model uncertainty: how it arises

So: read & prepare notes carefully, from within what I want to write down (i.e.) write it!

③ How this "incomputability" is different from inc. of $\sqrt{2}$ (Th. Pythagoras)

④ Errors of unk. size due to invisible regis in corpus (If we have a register SSZ, we do have a bound on P. error — But the bound will not be w. SSZ.)

⑤ complete \leftrightarrow incomputable

⑥ Always examples: always a better model very beyond.

⑦ Certain Scientists: have confidence in low error estimates for system failures.

150.07 (1) incomputability of UPD usually irrelevant to application

(2) what incomputability means in terms of unpredictable models. Uncertainty of "error size"

a.g. Pseudo Random say uncertainty that we have best model — & how much better our prediction is

(3) Prob. this difficulty arises independly of using UPD as approx for predn. — it is true in all

151.00 at any pt. we will need to use stochastic models, \therefore uncertainty arises from incomputable models.

NIPS

197.40: From the foregoing, it would seem that the incompleteness of UPD would have no practical impact on our use of it or any other model for prediction.

Not quite true! I will mention two situations in which the incompleteness of UPD is of immediate importance.

One is in intelligent systems ~~with~~ human machine. We have several possibilities in different domains that must be evaluated in a certain total time, T. If we were able to know how close to optimum we were ~~at~~ ^{each} our probability approximations, we could spend more time on evaluations that were furthest from optimum. The incompleteness of UPD assures us that we can't do this.

Intelligently allocate time to each problem so as to get most or an optimum

Another problem is the evaluation of the failure probability of very large ~~and~~ complex way,

systems, such as a nuclear reactor, a national defense system, the launching of a space shuttle, or rules designed to solve ~~the~~ ^{a set of} ~~problems~~ ^{problems} ~~concerning~~.

In either case many systems of this kind, the ~~analysis~~ ^{built} system whose failure ~~probability~~ ^{probability} we are to determine, has not yet been ~~built~~, or it has been tested in too few ways for too small a sample size for this data to be used directly ^{in the useful} ~~for~~ ^{extremely} small probabilities.

Instead, the probabilities are evaluated "a priori" using ~~test~~ ^{data} on parts of the complex system. Various models ~~to follow are proposed and~~ ~~sequences~~ ^{sequences} of events that might lead to failure are proposed and evaluated. Strictly speaking, there can be no "incompleteness" ~~problem~~, since ~~(2.2)~~

In the R.W. system like .11-.12 the ^{seq. of} events to occur have to take a limited amt. of time — So there can be no "Halting problem" — ~~There can be~~ "infinite loops".

Let (21) each ~~event~~ ^{event} in sequence can only take a known, finite time — yet for practical purposes, the number of possibilities and the logic time needed to evaluate some of them, make the computation ~~essentially~~ ^{practically} "incompleteness" from a practical standpoint

U.S.P.S

(Spec 199.22) The ~~the~~ incompleteness of the UPD is usually not relevant to ~~the~~ problems in practical prediction, it is of interest in the philosophy of science.

~~Some~~ Scientists are repeatedly disturbed by the need to revise ~~our~~ ^{their} understanding of ~~immutability~~ ^{their} sciences. They look forward to ~~the~~ ^{the} "Final Progress" that will put an end to all revisions. ~~Final~~ ^{revisions} ~~and the~~ ^{cannot ever} ~~incompleteness~~. ~~But~~ ^{we} ~~the~~ ^{cannot} ~~incompleteness~~ of the UPD assures us that this ~~never~~ ^{can} happens with any amount of data and finite computing time, we ~~can~~ ^{can} never know that ~~we~~ ^{we} have the best, ~~possible~~ ^{found} ~~the~~ ^{Final} ~~Final~~ ^{Theory}.

For many of us, this is not a cause ~~of~~ ^{of} ~~concern~~ ^{concern} but find it to be a source of ~~continued~~ ^{never-ending} joy in discovery.

For many of us, this incompleteness assures us that science will continue to be a never ending source of joy in discovery.

Big but a feature!

~~Let us return to~~ ~~the~~ ~~ancient~~ ~~Greeks~~ ~~who~~ ~~discovered~~ ~~that~~ ~~the~~ ~~square~~ ~~root~~ ~~of~~ ~~two~~ ~~could~~ ~~not~~ ~~be~~ ~~expressed~~ ~~as~~ ~~the~~ ~~ratio~~ ~~of~~ ~~two~~ ~~integers~~. ~~It~~ ~~is~~ ~~his~~ ~~number~~ ~~was~~ ~~"incompletable"~~. ~~None~~ ~~of~~ ~~the~~ ~~best~~ ~~approximations~~ ~~were~~ ~~used~~ ~~and~~ ~~eventually~~ ~~it~~ ~~was~~ ~~understood~~ ~~that~~ ~~the~~ ~~approximation~~ ~~could~~ ~~be~~ ~~arbitrarily~~ ~~close~~ ~~to~~ ~~the~~ ~~"true"~~ ~~(in~~ ~~incompletable) value~~. It took the mathematical community several centuries to get a good understanding of this problem, but well before that approximations were made and used. None of the approximations in a ~~any~~ ~~real~~ ~~sense~~, the square root of two was ~~incompletable~~ ^{was} actually $\sqrt{2}$, but they got arbitrarily close.

BIG Q: Does using purely Apri Model give unbiased estimate on testing?

It would seem that if σ^2 was unbiased then usually σ^2 would be not unbiased!

Now, if one uses rms error, $\frac{1}{n} \sum (x_i - \bar{x})^2$ given ~~the~~ ^{normal} ~~poor~~ ^{poor} distribution, one may get unbiased error for σ^2 (or σ) — For any power-law error $e^{-\frac{|x-\bar{x}|}{\sigma}}$ is ok for p.d. — T.L. value of σ is ~~invariant~~ ^{invariant} to "fitting". The ~~the~~ ^{the} true expected error will be ~~the~~ ^{the} wtd mean of ~~the~~ ^{the} something like $(X - \bar{X})^k$, if it may be that this will be "unbiased".

What does "Unbiased" mean? — One poss: that if one uses true σ for ~~some~~ ^{some} source of ~~a~~ ^a bunch of ~~copi~~ ^{large} ~~copi~~ — Test ~~to~~ ^{to} average ~~error~~ ^{error} will be ~~the~~ ^{the} actual error.

This is apparently trivially true. It is true of ~~any~~ ^{any} error ~~or~~ ^{or} estimation method.

NO! In 32 I was thinking of observed error rather than predicted error ~~predicted~~ ^{predicted} variable. However, in f. case of ~~linear~~ ^{linear} ~~predn~~ ^{predn} of ~~a~~ ^a ~~val~~ ^{val}, it may be that ~~if~~ ^{if} we sum over all poss. "parameter vectors" ~~in~~ ⁱⁿ space we will get some kind of ~~unbiased~~ ^{unbiased} error.