

20 : In the interest of "diversity" we may make ^{with} some power (≤ 1) of plain PC's start ^{with} QA copies
 ("Diversity" means access to substrates of a ^{great} variety of conds.)
 I could try $4+325.31-32$: i.e. with no. of QA's would w. "acceptable PC". T. "acceptable threshold"
 for each problem is a part of the problem, dem. — in which case S-induction is very much like induction.

HVR. While Zivenc is a very general universal P.P., I'd like methods of S-ind. soln.
 for various other common forms of S-functs: like Scuzed from induction (using
arbitrary context sets for induction... is other simple S-funct-forms, \rightarrow 3.00

10: **INEPT**
SN On Inept Monte Carlo Search is the St. Petersburg paradox.

Suppose the browser knows the soln (to a INV problem) has $PC = P_0$ & $CC = C_0$.
 Then if we do mt. Carlo trials w. diff. probab. of a trial CC its PC is unknown to estimate
 its PC , so we can spend $CC_i = k PC_i$ in trial i . We slowly \uparrow k as i such progresses.
 My impressn is that one spends an inordinate amt. of time on ~~the~~ by PC conds that
 didn't converge. > But still, what is expected time to soln?

If k is some function of the "trial no." i , ($f(i)$)
 To get to soln. $k PC_0 \geq C_0$ so $k \geq \frac{C_0}{PC_0}$. also, probab. of soln at each time is PC_0 .
 So we expect $\frac{1}{PC_0}$ trials after $k \geq \frac{C_0}{PC_0}$.

For a given value of k , a trial takes expected $CC = \sum (PC_i \cdot k) = k \sum PC_i$.
 $\sum PC_i^2$ is usually $\ll 1$ but can be ~ 1 . So after $k \geq \frac{C_0}{PC_0}$ (if its large enuf)
 we need $\sim \frac{1}{PC_0}$ trial costing $\frac{C_0}{PC_0}$ each = $\frac{1}{PC_0} \left(\frac{C_0}{PC_0} \right)$ which is $\frac{1}{PC_0}$ times CJS , $\frac{C_0}{PC_0}$
 $\frac{C_0}{PC_0}$ would be $\frac{1}{PC_0}$ needed for regular LSuch.

Now that's the time needed after we get k large enuf. How fast $\left(\frac{C_0}{PC_0} \right)$ should ~~be~~
 k increase w. i ? The i problem may not be easy to solve, for whatever ~~is~~ strategy is used,
 then kind of LSuch is $\sim \frac{1}{PC_0}$ slower than regular LSuch.
 If there are few type pc conds $\sum PC_i^2$ can be $\ll 1$.

30 : Actually $\frac{1}{PC_0}$ trials costing $\frac{C_0}{PC_0} \cdot \sum (PC_i^2)$ each or $\frac{C_0}{PC_0} \cdot \frac{\sum PC_i^2}{PC_0}$.
 $\sum PC_i^2$ can be quite small: Say $\tilde{PC} >$ the largest PC, & it's quite small (but $\gg PC_0$).
 $\sum PC_i^2$ maybe $\sim \frac{1}{PC_0^2}$, and $\frac{\tilde{PC}^2}{PC_0}$ could be < 1 !
 $\sum PC_i^2 = \frac{\tilde{PC}^2}{PC_0} = \tilde{PC}^2$
 $\sum PC_i^2 = \frac{\tilde{PC}^2}{PC_0} = \tilde{PC}^2$ which is probab. $\gg 1$.

But even if the time spent after $k \geq \frac{C_0}{PC_0}$ is not very large, the problem of how k
 would have to be solved.
 $k \rightarrow$ funcn of f , ~~which~~ be solved.

574

so

"BOOSTING" = Say one has a large source of not very good problem solvers: ^{See 2.1} for imp. choice of understanding

These problems have numerical or predicate solns. They may have predicate solns over a set of possible string replies. If x_i are the problems & y_i the solns, then a PSM will map x_i into v_i we want $(y_i - v_i)$ difference (error) to be as small as possible for a linear comb. of the PSMs

The Boosting idea is a successive approach method to get good wts. for a large no. of PSMs that are individually weak, but very good in a properly wtd. comb.

The method used gives more wt. to the problems for which the past solns haven't been so good. So at all times we have a wtd sum of PSMs: & we introduce a new PSM that is supposed to be interested mainly in the problems that the old ones didn't work so well on.

It reminds me of my "Problem Pool" method in which the hardest problems first but I'm not sure there is any real relationship

The papers in Boosting are in D:/PS: Dated 1.4.05

Google to Boosting Algorithms to get set of papers.

One long JPLP (PS) paper is on "Text categorization" dated 2000: (Name is "Boosting Low Rate Text Classifier Good Bib. ps")

Text cat. work is of int. to me: I want to see what kinds of (abstractions & concepts) they use to build "learners" (= simple PSMs) - so I can use their bibliography.

As for the "Boosting" method itself - it's a way to solve a certain kind of optimization problem.

But I may find useful. I think it's not very economical w.r.t. size, because it doesn't re-optimize entire set of PSMs each time.

Boosting Density Estimation (PS 1.4.05): I am mainly interested in density estimation, so this may be a useful paper.

What is unclear to me: They have this New PSM & it is supposed to try to solve a certain set of (the most difficult) problems. How do they modify the PSM to solve a certain subset of problems? It must be modifiable!

I guess it's a "learning system" that has some adjustable params: so it can "fit" a certain subset of a certain (perhaps small) set of examples. So each PSM is a little "learning system"

It's a "micro learner". - first, one would end up w. "little learners" that could approximate for a case each. The total bc of the system would be very large! So poor general capability... but the sum of these systems & total boost seems to be bounded! How does that come about?? I will need to read those papers more!

Trouble is they tend to not discuss the "micro learners" in much detail.

An imp point is that "Boosting" is not just a method of getting optimum wts for a set of fixed methods. It seems to involve choice of "weak learners": So I really don't yet understand it.

See: PS Density Estimation by Boosting 2002. Pdf for perhaps clearer discussion of some of the stuff about (see 2.0)

STM

2 days ago 1 April; Prof Barnes

W

20:11:08 : Perhaps the jumping arena set of ~~the~~ O_i 's ~~is~~ work ~~in~~ problems, to O_i 's that work only (TM) problems is like an "Escape" in PPM! (I discussed this sometime ago! ... (320.31-321.03))

Anyway! T. main diff. probs for "Strong AI."

- 1) writing good TSO: a. for Lsrch or "ProblemPool" & its variants.
- 2) Having no In Lschr phase of TM say "phases": Having found several O_i 's for various corpus sizes: How to efficiently srch for O_i for larger corpus. How does the new srch have access to regions that were visited in earlier srchs? How do we avoid testing entire past corpus, when we make ^{for our new cond} ~~what appears to be~~ a smaller changing/past successful O_i (of smaller corpus). (How can we know ~~what~~ what part of past corpus is affected by a modification of O_i ?)

13 3) What are some good ways to represent s-functs? TM should be able to decide which to use in any particular problem. Also how is information on effectiveness of various s-funct forms transmitted to srchs for different s-funct forms?

→ (6.20)

17 8 Φ What are some ^{fast} efficient ways to update the P.d. that Guides Lsrch (or any other srch)? How can we ^{find} good concepts (2 subareas ~~for sets of subareas~~) for induction into PPM.

Can discuss! Suppose we have an "Adequate TSO" for Lsrch on the QA problems.

By definition, this means that Lsrch would work. Hvr, the CJS can be normally very large, if it is actually PC:cc., where cc is cc needed to test entire corpus. By using the Recog. functs, we can reduce corpus needed for testing. T. analysis, (by the trainer) of the TSO must include discovery of R functs. (Remember that T. per is a major part of time spent on cond. ... we may reject cond's in much less than just five. Hvr, if there are some cond's of high pc that do not converge, all of T. per will be spent on such cond's.

20 Still, the R functs could considerably reduce pc ~~to~~ for the cond that the trainer has in mind, and it is this ~~is~~ that is the "CJS". (The ~~idea~~ by no means, "Solves the problem" - let's go on to another point) For each QA, the trainer has a threshold of pc to be assigned to the "victorious war". Hvr, there is the possibility that the "true" generator, M , of the ~~observed~~ observed A_i 's will occasionally generate an A_i of very low pc ... much less than the trainer anticipates!

On the other hand, remember that ~~the~~ induction is an O_i prob ... we just try to get as good a O_i as we can in the available time. This would seem to mean "no acceptance threshold" for any particular A_i (so it's not really much like a FNU problem!). We want to maximize the pc of the corpus within the available time.

A Reminder! That the idea of a TSO is that, after solving the early problems, TM has within it various concs. that would vastly reduce the srch time for finding solns. to subsequent problems. These concs can be in the form of definitions also in ~~observing~~ collecting of data on

STM

10:34:00... t.p.b. of ~~tokens~~ "Tokens" that follow various kinds of "contexts"
 Seems to me that T. TSO approach to TM has to work! (Expresses) forces
 solves no problems (in terms of ~~costs~~ costs) may be significantly different from way humans
 do it. This may give a large (but I hope manageable) "constant factor" better, homogeneous
 TMs, in solving problems.

To start off, we can let trainer derive TSO's that don't need R tokens to be computed by TM
 These will be separate TSO's for each R - ... a different subject.
 They will each have a different (set of) O^j functions, but. They will be tested on
 different corpora, but the search data used for searching for E 's will, to a large extent,
 be shared... Also in searching, TM will know which O^j 's in i 's corpora of
 solutions were used for which R 's (i.e. O^j will be "indexed"/ "indexed").

So: to start, we may ignore R problem, & i 's problem of testing either corpora on
 to entire QA corpora... because that corpora will be small.

Then: We will be penetrating O^j 's and testing them on QA corpora; How to do i 's testing
 economically: This will depend much on forms of i -functions. ~~As a result~~
 One way to think about it: we want test values of $O^j(A_i | Q_i)$ (\approx pc.).
 Expressing O^j 's as \exists some world is seem to be good in terms of cc, but need high
 forms of Universality. If we make O^j an arbitrary (perhaps so)
 function of A_i & Q_i , then we have to normalize, otherwise its maximizing (let to maximize i 's product).
 To normalize, we have to consider many possi- values of A_i ... which sounds expensive
 and ~~very~~ very inefficient.

23; 385,24 - A0

26

A pseudo 3-issue way would be to have several O^j 's: each O^j corresponds to
 a subset ~~subset~~ R values in pub of i -issue. Say we have 4 or 2 such functions
 They are sum for all QA problems... (which is a true approx of \exists issue).
 We then use "bit distance" to test one of i 's outputs is closest to target A_i .

I discussed .23-.26 before, but I seem to forget the details - here (there are no
 complex details): .23-.26 is a correct approx. The main deviation of (.23-.26) from
the original 3-issue in the pub is that in this last case, R can take values that have different
 wts. -- i.e. not all R values are of \approx wt. (as in .23-.26). In (.23-.26) the only way to get more
wt. is to have duplicate or approximate duplicate outputs enough to make the 2^k outputs
 Actually, it is possi to make an exact "approx" of \exists issue with (.23-.26) by suitable exact
duplication: say $k=4$: then if an R in \exists issue had $\frac{1}{2}$ wt. of particular A output W ($W \in R$), we would
 have $\frac{W}{2^4} = \frac{W}{16}$ duplicates of that A output.

Here, .23-.26 slows down the testing procedure by a factor of $\approx 2^k$. TM could learn
 to find the right functions of i 's 2^k to get \approx each A_i output - which will signifi-
reduce cc's by Quick Abort(?).

Mr

20:4440: In line w. 9.2.3-26: Perhaps have very many O_j 's ^{/each of which} ~~are~~ fairly good at solving Q's as poss. For a new QA - to do this problem, simply mt. from all ~~the~~ no. of O_j 's that will ~~do~~ ^{work} on it. Perhaps only count ~~of~~ ^{of exact} hits on Q's.
 — Maybe more w. for O_j 's that had by scores on "closely related" Q's.
 This needs much theoretical analysis!

Suppose as .00 ff we have soft Recog. functions and many O_j func's & we have R func's ^{2 wt} (No ~~no~~ to each O_j for each poss QA. We would then get products for all QA's & we could get the products, which we want to Max. — But what is the pc of the dcou? Many O_j 's are ill so we can't simply add their β costs.

In the past of EDSIA re-part. cu Rec func's, Peter went + bunch of Rec, func's, together were able to decide which O_j set, should be used on a given Q. These elementary R func's were 2 func's a' they could be used to form a subset on the input Q's.
 Note that it's much better to have 1 R func's because we can f. couple up (determinately) ~~are~~ ^{sharp} ~~are~~ ^{are} able to save c.c. in testing w. certainty.

In using \sum R func's, we may be able to use ~~the~~ ^{prob} problem testing on randomly chosen ~~QA's~~ Q 's of f. corpus. (Batterman → chose R. ~~to~~ ^{trials in} ~~to~~ ^{strict} "diff" order).
 .00 ff is closely related to a recent "problem pool" method, in which I would try to find O_j 's that would do many of the prob's poss... picking various subsets of f. Q 's, to "specialize" on. There, the idea was to use the partial solus (tho ~~it~~ ^{is} ~~not~~ ^{correct}) as a source of "parts" to eventually find an O_j that solved all of f. Q 's satisfactorily. (tho ^{is} ~~not~~ ^{usually} ~~not~~ ^{able} to do so.)
 In .00-17, we would be satisfied w. a set of O_j 's — (no individual ones that were able to work all problems) and an R func to assign Q 's to those O_j 's (or O_j 's to those Q 's).
 Say we only had 2 O_j 's & an assoc R func to assign pc's to Q 's — as to which of f. 2 O_j 's to use.

We would then code the corpus of n Q 's in 2^n different ways, corresponding to different choices of O_j 's for each problem. If, for a particular Q , we assign P_1 to O_1 & P_2 to O_2 , & O_1 gets P_1 for the answer, then the system assigns $P_1 P_1 + P_2 P_2$ to the correct answer.
 For the entire corpus, we multiply all of these sums for individual Q 's to get final pc for f. corpus, & But we also know to multiply pc of f. & func & the pc's of ~~the~~ O_j 's of O^2 .

It would seem that if O_1 usually gets small P_1 for various Q 's, then we should not completely "change" its entire pc of dcou? ^{against ?}
 A (False) assumption ^{against ?} (27-31) is that as the weights of O_j (or its P_j 's) get smaller & smaller, we have a discontinuity at P_i (for all Q 's) $\rightarrow 0$, since the pc of O_j 's dcou. becomes suddenly incorrect. A counter argt. to (27-31) is that we always have all codes in which O_j is not used at all, so they don't have the penalty of the pc of O_j 's dcou cost.

For ~~the~~ say we have many O_j 's: for every subset of them, we know pc of f. corpus w/ from mult by their pc's of creation! Mult by pc of creation of assoc. R function that assigns O_j 's to Q 's!

STM

00: In general, this best part of the R func may differ for each subset of U^i - \rightarrow will it's pc of creation.
-05 If there are n O^i 's there will be 2^n possl codes (including the identity code).

02 (5.32-33) Because of the redundancy adding of pc's (as 5.29), an additional O^i may not be such a big expense (pc of creation - wise). Also, say O^1 was used for all but one Q , in which case O^2 was used. T. pc of O^2 creation plus (the) opt. Recog. funct. would, indeed, be needed to justify using O^2 on just Q .

04 T. move (is: It is much better if we can do it w. only one O^i , but if not, then use $>$ than O^i .
Also, t. O^i 's need not be very expensive, if they share many tokens in "contexts".

So perhaps 5.18-6.09 may be a viable way to use many distinct O^i func's; There remain Q 's,
1) What is form of s- func's? How to generate a "very complete Univ. set" (Univ. t.)
2) What is good way to do such for such functions? Any good way to avoid ~~the~~ verifying ~~the~~ entire corpus, whenever a new O^i is found?

Re 1) Such! A not bad way is in "Prob. Pool" Mode! to try to find O^i 's that work on many Q 's well.

I wrote a fair amt. on this approach (4+325, 30 - STM. of is relevant, but: 320-31) Base Idea
One difficulty was how to wt. Counts across that solved varying nos of problems w. varying pc's, so one could properly (in token term - say via PPM)

211.30-10
314.20 - 1 imp
315.15
316.09
5TMS.184

2.3.16 1.7.05 On representation of S-functions! Some problems/goals: 1) we want it to be "easy" (low cc) to get pc of an object: This is to testify Cond. functions 2) we want it to be easy to generate objects of the pc ... say in a pc order: This is for prediction: it can be much slower than Q , hrs, since we do 1) much more frequently (in fact, def. STM) 3) T. method of generating t. s- func's: It should be universal & it should be easy to ascertain pc of each s- func. 4) T. representation system should be match w. "Recog. detection scheme" (e.g. by PPM) ... so it's schema can easily find reg's in func's of that "language"

2.2.20 EN On representation of functions (Both S & f): Instead of sum (A, B, C) where A is some express: keep A as a token; so all operations are in a simple form of finding ~~tokens~~ tokens & how they are to be combined. So all reg's are of the form A, B, C, where A falls low B & C are to be combined. This is supposed to make it easy for PPM to recognize reg's. So a ppm would be a series of definitions of "Non-Terminals". This would normally be written in a number & it could be easily compiled into a tree (but maybe not for recursive func's).
Hrs, I'm afraid t. reg's. would not work as is: It's really accy that t. sub trees be recognized - i. Pr. reg's system doesn't do this in a "reasonable way".
There is, hrs, t. passy that if we use a system that would not recognize reg's of a certain

STM

S - funcB .17

oo: sort. That 1. system would like to "eat a long" w.o. using flex regexs — by perhaps creating functions that don't "use" them (?? ... unclear on this pt... what's the press or likely!)
don't "read"

[SN] ^{As sub} Generators of "Context": A context is a set of suffixes that reduces entropy of the corpus following it.

En past we had considered only 8. immediate "info" of a context. Now we consider any prefix that compresses the following corpus. If a certain context makes it likely that the next symbol to follow would be "t", then this is a ^{potentially} context.

T. This is a sub-generator of my Global generator ^{local} context... as the entire previous corpus is known knowledge of its influence on the production (≡ compression ≡ entropy) of the future. (As opposed to "MacroContext")

ON INU problems: Use of L search for soln. of INU probs seems to work the probs by direct search — but humans usually convert them to OZ problems & solve them as OZ probs. But a good regular INU search would be able to find the best way to solve the problem... If it would convert to OZ, it would do that. Furthermore, we can design TSO's that encourage enable that kind of problem soln. — So L search can solve INU probs "optimally".

A trouble is: in phase I the OZ solns aren't very good. The OZ methods it causes aren't v.g. In phase II we do much better.

S - First Representation: For D.P.'s on rules, a very common way is to have 1 output be the use of a Gaussian or other D.P. (say Gamma D.P. Since these kinds of funds are the main kind needed for phase 2, this means we only need to do rules and introduction in phase 1 (if we need phase 1 at all!). Actually, we only need rules as output.

~~Rules~~ Rules can be optimally fitted to the data — as in usual time series regression

[1.22.05] For S - func's w. string output: S-CFG's are not bad. We can modify/generate them

in various ways (e.g. transformational grammars). Actually, many different types of grammars have been proposed for English. We can build from these devices a grammar that can generate all known grammars; plus any others. We can generate a universal grammar, grammar — but the main idea is to find a grammar grammar that will easily generate grammars that are known to be useful.

S - context sensitive grammars are perhaps universal & quite powerful.

"Universal" in the sense of being able to simulate a Universal Turing. They are also nice because they are an "add on" to S-CFG's so we can have FSM learn CFG's then probably have for CSG would not be very different (?!),

The easy way is know to discover S-grammars is BU — i.e. hill climbing — starting from a seed (very expensive) grammar for the corpus & successively improving it... usually greedy, but may be not always. (8.00)

10:7.40: Another, perhaps very likely way to generate s-strings is by a relatively short seq. of operations — logically defined, but w. some uncertainty in some of the operations.

This would seem to be a way we answer questions that require string replies ...

Q's themselves can have strings w/o numbers in them

Distributions of the kind can arise from "case" state models becoming not 100% accurate — having PDS for each case, rather than deterministic results.

Would like s-strings to have non-zero pc for all string outputs ... but we may avoid this

Need to in a fully way:

Given a corpus, we fit a s-function to it so all of A's are chosen, in Param we find params to give a whole model max pc. This model can then predict.

It can give ~~pc for~~ pc for certain futures even if, however. If it does, actually do so, we will have to completely discard it as a kind of new s-model.

Simplest example is Bern seq. w. h symbols: we know a simple k component vector that determines choices.

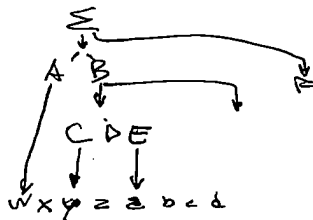
Next is PPM & Various Algs to decide on what context to use to get a k-vector for choosing next symbol. The Algs can do all probs, so we add Bern (in a wtd. way).

Next, we do operators in the past to get pd. vectors for future symbols. The operators on the past give something like contexts that are useful.

SN in CFG's when we generate a string: Let us look at "derivation"

I think each terminal symbol in the final output has a "linear" history of choices.

So perhaps we can usefully use a PPM to predict the symbols — using that kind of "regularity". Let's look at the derivation tree of a sentence:



The letter 'z' has context $\{S B C\}$ or $S B C$.

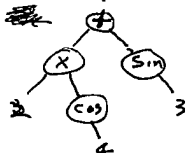
w has context $S A$

so a & b have common context, S .

y & c have common context $S B$.

We can do some kinds of contexts in construction of functions (which have a CFG) Unless they are use contexts?

order simple function trees:



4 has context $+, x, \cos$.

2 " " $+, x$

3 " " $+, \sin$

These are not, however, the contexts our gets in a Polish or RPN expression of the function.

$+(x(2, \cos(4)), \sin(3))$

456 P SMC
Stat P on line 25min, 1hr.

20 : Is there any good (cheap) way(s) to implement this kind of context?

I want to include my method of evaluating the utility of a particular improvement of search algm. I had a nice simple equation for the efficacy of a method, but I haven't been able to find it! The eq. appears about line 33 on left side of page.

It may be possible to create it: Maybe this: Say I have several such methods of various speeds & strengths. They have different ΔC profiles w. cc spent. Each time we use a search method, we get a ΔC which if our corpus size is quality, so each method helps itself: ... But the ~~very~~ analysis ended up w. some surprising factor that was very reasonable! Try to find that analysis!

I think I made refs to it several times in TM notes.

Was it w. the timer $\frac{G}{G_{max} - G}$ was destroyed?

21 SN on K. Szlotowicz & J. Schmidt. Productive Incompleteness (PIPE) 1997 PS: 7/26/07

FRAM (81.17 memo common.

185.00-03 memo common comment.

To generate prog, they use fixed, standard trees: k types of functions, 2 kinds of literals so each branch pt has k+2 possys. $\frac{1}{2}$ is to max out of all functions.

There is a "common" complete, infinite ~~tree~~ w. n branchings at each node. This tree has simple coords of each node. When a function is created, each of its nodes maps onto ² certain one of the ~~nodes~~ nodes of the "common" tree. This enables us to compare/nodes of diff. functions — so we may learn.

Main kind of trap: we probabilistically construct a population of candidates using probs at the node of the common tree. We take the Best cand of each generation, and move probs at nodes to make the Best of Generation, more likely. } for Generation based trap.

Also we occasionally do the ~~same~~ thing using the "Best Cand of all generations Run for" — "Elitest" trap. Also they use a certain mutation rate ^{of probs in "common tree"} to get trap. also; but mutation algm. is very strange! G.2 at Document

I certainly haven't gone into this paper very deeply: But my impression is that the method used to learn, is not very good. However, associating terms probabilities w. each node of the "common tree" is amusing — it might actually lead to good functions! That they only learn from the best cand. in each generation, seems wasteful, however.

In General, it might be worthwhile to go over this paper more carefully.

The probs at the nodes of the common tree correspond to the probs in "fixed length genome" of PBL (Population based Incompleteness) (Baluja & Caruana, 1995) — they also use copy to learning methods of PBL.

Grammar Model-based Pgm. Evolu. Dec 2003 - Guy in Australia - mostly Gubarna (!).

This seems about the same as what I had in mind! The use of a SCFG for summarizing a population of function trees, is very good. ^{useful/common} Sub-trees are recognized in the Grammar discovery process.

Herl PT vol II, 127: They're doing this close to ordinary GP. The way they select population to keep, etc., is not v.g.

My ^Q ~~Q~~ _{Genes - M_g} would probably be much better. They just use the "top 1%" of the population to create a SCFG; so Q is unlikely to be optimum.

Q by Q is whether to use SCFG or some kind of enhanced PPM. This last makes it possible to optimize Grammar very frequently, which does ~~not~~ "do it". Also, depends on how good SCFG's they can get in "available time". At any rate, it would seem that SCFG's would be able to discover any ^{common} sub-trees in the set of candidates.

So my improvements ^{on} ~~to~~ PEP paper are:

- 1) (07)
- 2) PPM may be better than SCFG... but I would stick with it. If they have a usable SCFG discovery program - try to get it from them.
- 3) They don't work with nice problems. They just use problems that others have worked on to show their methods superior. They don't do TSP.
- 4) They don't give system ability to "watch ~~the~~ fitness evolve".
- 5) Phase II in which the constructive searching of s-models is reproduced as a problem for a system to solve: it has a universal model accessible to it.

Think about simultaneous L search & BV search, using S-CFG for PL guiding. Doing PC order is V. diff. later.

24 **HA!** We can do Simult L search & BV search S-CFG's! (or S-CCG's);

It is "easy" to list Cands in a PC order: Use ordinary L search method of waiting all nodes until pc gets below threshold = T_h ; then do the $\frac{1}{2}$.

My main idea was that it was somewhat like "A river today".
Main idea: first given a PC threshold, it is always possible to list all Cands in $PC > T_h$.
It may be like Paris: one explores the trees of decisions until either a cand is completed or $PC < T_h$ - at which point one backtracks, (perhaps increasing PC), then continues. Alternatives backtracking (which is time efficient, but messy in practice?), start over at "root" whenever we complete a cand or get $PC < T_h$.

What the eqn does, is simply an exhaustive depth first search of the space, in which the termination of a path is either a leaf or $PC > T_h$. It keeps track of the PC level of each node so it can update T_h properly after a "backtrack". I think a stack can be used to keep track of where one is in the tree, & what the PC's of the nodes are. A possible serious intention of this algo is that it can spend much time going out on branches that it PC, so it becomes unlikely that it will have time to complete the search. - Herl, if one is doing L search, one of the ~~problems~~ -

termination criterion is not $PC < T_h$ but $CC > T \cdot PC$. This last results in

In discussing Advantages of ALP (→ Stoch Comp) over MDL - I said that the diversity of Models was useful in dealing with new depts. problems - so we tended to use A.P. ~~rather than MDLM~~ A.T.

Another area of use is in G.P.: If we use single grammar to Model population, this is MDL: However, if we use several grammars, we have much less likelihood of getting stuck at a local optimum - which is called premature convergence in "GA talk".

If we are using S.C.F. grammars to Model our population, we can get diversity in our models, if we obtain them by ~~"early splitting"~~ "early splitting".
What this means: In obtaining a good grammar, we start out with a ~~grammar~~ grammar that fits the data, but ~~it~~ gives low probability to the corpus. To obtain a better model, we try modifications of the grammar... Checking out many are two common types of modification. Often, there will be several possible modifications that will all give better fits to the data. If we follow these branches out and continue improving them, we will tend to end up with several grammars that are quite different - giving the desired diversity of trials.

Possl. Inadequacies, Bottlenecks in IAP Model(s) of TM.

1) A MTM - QATM could be v.g. using PPM, (sub)string context oriented PPM, S-CFG, S-CSG
~~Real Bottleneck~~ (a) Diff't to write TSG (b) It could do self-improvement of its ~~induction~~
induction such a thing, since this is an s-induction problem. The system s. can't

ever be very smart (unless I search long I've found good stuff) → Nov. 2000 13.00!!

Also TM would probly have to learn ~~the~~ Methods (eg. c's how to use it (?))

The this MTM could be a "small TM" project. Could use useful experience coming from PPM, SPPM & CFG in such.

2) I'm not sure CFG's can recognize sub ~~strings~~ (sub)functions occurring with frequency. BTW (sub)functions occur very unusually

frequency ← BTW I think ASCII and NT's always represent the units only, It may be that a S-CFG isn't any (or much) better than PPM | Perhaps each word w. smaller SSZ | Also reps detected by SFG could be associated with reps recognized by PPM. → 12.00

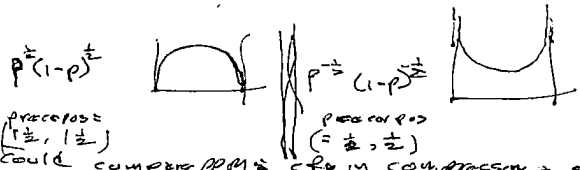
3) If 6 long for expressing Functions (≡ conds) has deductive facilities, then

PPM can be v.g. for noticing unusual freqs & unint freqs w. context. Still, even after a function has been "defined", it will not recognize F (α , β), where "don't care" is of any size. 37

4) Perhaps best search method would be sub-tree context PPM; but implementing it would be non-trivial & possibly slow. Normally speed & efficiency can be traded off - but certain kinds of reps may not succumb to that sort of "Boogey". → 12.20

57. DGA F (n , α) ! "n" means "Don't care"., so f (n , α) is an unint. So the source CFG's can recognize some units, may be that can recognize f (n , α).

Note that "n" can be a legit N.O.T. i.e. to set of all finite strings; except to define! May be includes own + "free" definitions.



30
SPEC
11.29

We could compare PPM to CFG in compressing a set of unordered strings. Note that PPM is normally not used this way, so we have to be careful in comparison!

PPM can be used for unordered strings: I did devise a vgg. very fast way to do it. We just use a special "initial" symbol, (Δ , say) before strings, & fix the whole thing like a sequential corpus. We can get p's for first symbol in any little string. It has context Δ only: somewhat like in fact any context string is allowed to go past Δ . (I may have already written much on this!) I think this is easy to arrange when we put all p's of corpus in lexical order. For the initial symbol in a string, we may have ~~no~~ no context - just frequency.

For the final symbol, we ~~have~~ usually have much context, and for the Δ following, (which we ~~do~~ do need a PC for) has much context.

For these unordered strings, we can also use PPM backwards: see if we get same compression.

→ Hur. Note 11.37! CFG's may be able to recognize $F(\alpha, \beta)$

The "AND" of CFG sets $F(\alpha, \beta)$ AND $F(\beta, \alpha)$ is $F(\beta, \alpha)$.

Hur. in CFG's we do OR (union) of sets, but not AND (intersection) of sets. There is some O'Leary proof of the "n" symbol, Hur.

30 : 11.35 : 5) For QA induction, I considered using $\prod_{i=1}^n P(Q_i|A_i) \geq z$ (partial) GRC, because it was easy to devise stats for $Q_i|A_i$ or $P(A_i|Q_i)$. We can do this for all forms of QA... (string of numbers). This makes it poss. to start w. poor grammar (i.e. A-Grammar) & still climb by checking & fixing & prob. other operations. T. troubles, P is for a CFG is not universal. T. CSG is universal, but I am not sure as to whether it's an easy way to write P's in.

Ordinary ones just stop after they print an output. all define s-grammars, but they are not practical ways to define a P.D.

30 One Poss. way: Consider a pgm with a naive "case" decision on list of integers (to n; its output is $\{\alpha_i\}_{i=1}^n$ α_i are strings say. This pgm has a certain length to decide unordered set, $\{\alpha_i\}$ (including the set of choices) How can we mod. by the pgm to find the total PC? ~~the total~~ But find steps to find PC? The only prohibitive part is the 1 to n integers: we can lower that number & try to find PC of a set $\{\alpha_i\}_{i=1}^n$: Unstartly, this would not easily translate to a new algorithm.

Another reasonable posy is 13.18: perhaps close to how humans do it?

00:12:40 6) Would it be poss. to get Phase 1 to work on improving induction s rel for d-induction?

01 The corpus pairs would be $\left[\left(\begin{matrix} \text{sequence of accepted } Q_i^j \text{'s upto } t: \text{ is problem} \\ \text{if } Q_i^j \text{ for the } t+1 \text{ problem} \end{matrix} \right) \right]_{i=1}^{T-1}$

02 It may vary per $\left(\begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \right)_{i=1}^n$ to be interpreted,

It would seem that in theory, this is a subset of a s-induction problem. To work it, TM would have to have an adequate TSO! So Phase 1 is capable of SI.

How, in Phase II TM has lots of other OE problems, so the training situation is, presumably, much better. ALSO (N.B.) Phase 2 UNDERSTANDS optimization Much better than Phase 1 does.
Also Note Phase 1 can improve its s-induction algms in way w/o need for d-induction to re-iterate 01-02 TSO QAs of corpus etc

$Q_i^j := (O_1^j, O_2^j, \dots, O_n^j) : A_i^j = O_{i+1}^j$; T, problem: Given Q_{n+1} to get a pd for A_{n+1} .
TM is given $\{Q_i^j, A_i^j\}_{i=1}^n$

So, to start, we will use (perhaps augmented) PPM or S-CAG to do this. But in solving this s-induction problem, TM is able to use a universal instruction set, so this is str. & diff. of (2.26 ("G")).

7) A poss. soln. to a pd on an ordered set of strings (i.e. the QA (help strings) problem).
Every AI problems (Simon & Newell) were of f. form: Given initial state of systems' set of sequences/

x-params $[F_i]$, x-lem Q_0 into A_0 using some metric on the set $[F_i]$, and there is some metric of how good the seq. of x-params is: The problem of $A_0 = \text{Min}$ function of Q_0 , using a set of (universal) re-define insts, γ 's of Q_i 's etc. Their "GPS" was a set of heuristics for this problem.

My take on this was the OP-OB algebra: Have to trainee learn to "chunk" various ops & various obs. to make the problems easier. A combination of OB followed by an OP is "feed back". (Action-observation-action = Experiment),

At any time, there will be a set of problems that f. trainees "Goodat" & Other problems will take longer, or get low scores, or be insoluble.

The technique of learning may be simply "chunking" & possibly "heuristics". Newell & associates had a "SAIL" system but was forgotten/ & obscure (System for AI) — it was rather complex.

T. Big Break Rev in GPS was the idea of "Difference ~~operator~~" is a assoc w/ each OPS to that difference. Instead of optimizing a single Gene, Ray had a vector Gene & had different ways to deal w. each Gene component.

∴ $\sigma_{obs}^2 \cdot \frac{n+1}{n-1}$ for the Normal distribution: 1 multiplicative coeff.

* $y_i = a x_i + \text{Gaussian noise}$

∴ $\sum_{i=1}^n (x_i - y_i)^2 \rightarrow \min \frac{dS}{da} = \sum 2 x_i (2 x_i - y_i) \quad \frac{dS^2}{dS} = 2 \sum x_i^2$

$\Phi(a) = \sum (x_i - y_i) = 0?$

$\sum_j \left(\frac{\sum x_i^2}{\sum x_i y_i} x_j - y_j \right) \rightarrow \frac{\sum x_i^2}{\sum x_i y_i} \sum y_j = \sum y_j?$ if so then $\frac{\sum x_i^2}{\sum x_i y_i} = \frac{\sum y_i}{\sum x_i}$

$S = \sum x_i^2 - 2 \sum x_i y_i + \sum y_i^2$

$\frac{dS}{da} = 2 \sum x_i^2 - 2 \sum x_i y_i + 2 \sum y_i^2$

$\frac{dS}{da} = 2 \sum x_i^2 - 2 \sum x_i y_i + 2 \sum y_i^2$

$S = \left(\frac{1}{2} - 2 \right) \sum x_i^2 + \sum y_i^2$

Best try it empirically for various $a \rightarrow N^2$ values.

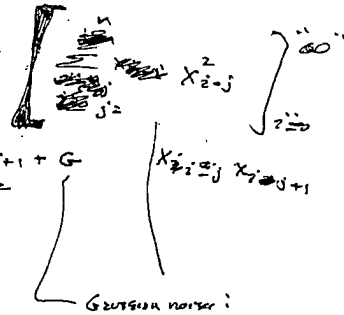
$A = : n =$

generate seq; keep / track of

$\alpha \quad Axx = \sum_{j=0}^n x_{i-j}^2$

$Axy = \sum_{j=0}^n x_{i-j} x_{i-j+1}$

$AL = \frac{Axx}{Axy} \rightarrow x_{i+2} = A \cdot x_{i+1} + G$
 $\text{sig} = (AL \cdot x_{i+1} - x_{i+2})^2$



$z=0$: ~~some~~ $z=0 \Rightarrow$

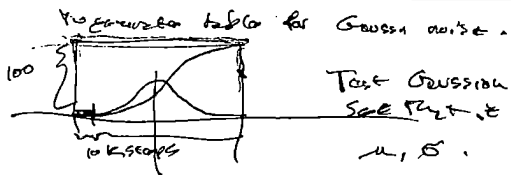
$z=2+1$; ~~some~~ $z=2+1$

$z=2+1$: ~~some~~ . If $z < 1000$ then good

$Sig = \text{sum} / n \cdot (n+1) / (n-1)$

$Sig =$

Initialize: $x_1 = 0$ for $i=2$ to 1000
 $x_i = x_{i-1} \cdot A + G$ Gaussian Noise



Test Gaussian fit to see that it works
 μ, σ

For $k > 1$, the trust of σ^2 could depend on out. coeffs : whether or roots of ϕ . And polynomials (and all inside unit \odot) could be very complex! I could just choose coeffs at random, then reject to seq. if it's not too far from ϕ - i.e. if σ 100 as like this divergence.

Best Drop Plus for time being!
 for All Notes: Probability $\frac{\sum x_i}{n}$ say that I'm not sure of formulas at present time.

STM

MEJ Paradox

'MEJ'

00:

③ **A BIG Q:** on SIAP 48, 02 > 20 + formula $\frac{\sigma}{\sigma_{max} - \bar{G}}$ is Question!
 It is indep of threshold if +. Given simply: order of merit of cand! So + formula depends much on \bar{G} func. w.r.t. monotonic functional - like only direct proportional selection in GA.

It may be poss. to use a functional of \bar{G} \rightarrow a suitable threshold so as to maximize $\frac{\sigma}{\sigma_{max} - \bar{G}}$. Simultaneous opt of threshold in functional form.

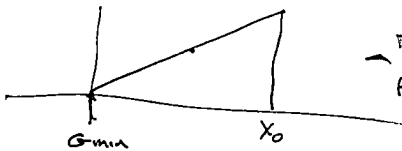
If would mean that the MEJ formula would keep to work if \bar{G} was a measured Obj/ly.

How, if only useful outputs MEJ is X_0 , + threshold, (or any other way to determine \bar{G}).
 Inputs used to determine the search language. But the output (i.e. corpus or language) depends (perhaps much) on functional form of \bar{G} + i.e. any measure X form will change output!

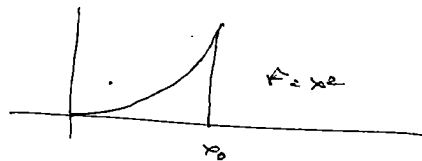
The PDF (or other probability) used to get a PD out: chosen lang; would like PDF M is σ of the PD is the same as Statist. population body model, if was use of MEJ eq. for lang selection.

look at picture at SIAP 48, 204: Muz suggests MEJ into word for that Density

Profile: Consider $M(x) = \frac{\sigma}{\bar{G} - G_{min}}$



$f(x) = x$; Area as in probability distribution



$$\frac{1}{X_0} \int_0^{X_0} x = \frac{x^2}{2} \Big|_0^{X_0} = \frac{X_0^2}{2}$$

This is first moment of $X=1$ from 0 to X_0

$$\frac{1}{X_0} \int_0^{X_0} x^2 = \frac{x^3}{3} \Big|_0^{X_0} = \frac{X_0^3}{3} = \mu$$

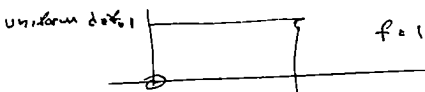
first moment of X (to X_0).

$$\text{mean} = \frac{1}{X_0}$$

$$\text{second moment} = \frac{1}{X_0} \int_0^{X_0} x^3 = \frac{x^4}{4} \Big|_0^{X_0} = \frac{X_0^4}{4}$$

$$\sigma^2 = \text{second moment} - (\text{first moment})^2$$

mean of square



zero moment $\int_0^{X_0} 1 = X_0$

$$\text{mean value} = \frac{1}{X_0} \int_0^{X_0} x = \frac{x^2}{2} \Big|_0^{X_0} = \frac{X_0^2}{2} = \mu$$

$$\text{first moment} = \int_0^{X_0} \frac{1}{\mu} \cdot x = \int_0^{X_0} x = \frac{X_0^2}{2}$$

$$\text{nth moment} = \int_0^{X_0} x^n = \frac{X_0^{n+1}}{n+1}$$

$$\text{mean} = \frac{M_1}{M_0}; \quad \sigma^2 = \frac{M_2}{M_0} - \left(\frac{M_1}{M_0}\right)^2$$

$$\frac{X_0^3}{3 X_0} - \left(\frac{X_0^2}{2 X_0}\right)^2$$

$$\sigma^2 = \frac{X_0^2}{3} - \frac{X_0^2}{4} = \frac{X_0^2}{12} \quad \left\{ \begin{array}{l} \mu = \frac{X_0^2}{2 X_0} = \frac{X_0}{2} \\ \sigma = \frac{X_0}{\sqrt{12}} \end{array} \right.$$

$$\sigma = \frac{X_0}{\sqrt{12}} \quad \mu = \frac{X_0}{2}$$

$$\text{so } \frac{\sigma}{\mu} = \frac{2}{\sqrt{12}} = \text{indep of } X_0$$

$$f(x) = x; \quad \text{zeroth moment} = \int_0^{X_0} x = \frac{x^2}{2} \Big|_0^{X_0} = \frac{X_0^2}{2}$$

$$\text{nth moment} = \frac{X_0^{n+2}}{n+2}$$

$$\text{mean} = \frac{M_1}{M_0} = \frac{X_0^3}{3} \cdot \frac{2}{X_0^2} = \frac{2}{3} X_0$$

$$\text{second moment} = \frac{M_2}{M_0} = \frac{X_0^4}{4} \cdot \frac{2}{X_0^2} = \frac{1}{2} X_0^2$$

$$\sigma^2 = X_0^2 \left(\frac{1}{2} - \left(\frac{2}{3}\right)^2 \right) = \frac{1}{18} X_0^2 \quad \left\{ \begin{array}{l} \mu = \frac{2}{3} X_0 \\ \sigma = X_0 \cdot \frac{1}{\sqrt{18}} \end{array} \right.$$

$$\text{so } \frac{\sigma}{\mu} = \frac{1}{\sqrt{18}} \cdot \frac{3}{2} = \frac{3}{2\sqrt{18}} = \frac{1}{\sqrt{12}} \text{ is constant.}$$

$$f(x) = x^2 \quad M_0 = \frac{X_0^3}{3} \quad M_1 = \frac{X_0^4}{4} \quad M_2 = \frac{X_0^5}{5}$$

$$\frac{M_2}{M_0} = \frac{5}{3} X_0^2 \quad \frac{M_1}{M_0} = \frac{X_0^4}{4} \cdot \frac{3}{M_0^2} = X_0 \cdot \frac{3}{4}$$

so $\frac{\sigma}{\mu}$ is again constant.

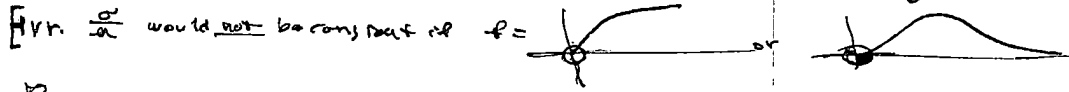
16.00

STM

"Linearity of G" .06

00:15:40

For models, if f is a power power of x , then $\frac{df}{dx}$ is constant. But any other (non-power) relation of f to x would probably give a non-constant ratio. So perhaps the constancy of $\frac{df}{dx}$ is a characteristic of x^n only.



Perhaps what characterizes $f(x) = x^n$ is that $\frac{f}{x} - \left(\frac{f}{x}\right)'$ is constant. $f = S f$.

Re: fitness func: I have 2 defs: 1) G is linear if G is equiv to $G \circ P \circ G$

2) G is linear if the utility utility of G to $G \circ P$ is constant (n. probly) & only a non-const of G . \rightarrow This is "utility is const"

I think these 2 defs are v. same, but no proof! Perhaps all electrical and differences between 2 states

0:13:40 8) On "Phase 2": I had considered starting out on phase 2 (instead of doing Phase 1 first)

This does seem like an alternative Path to TM & could be use fully pursued in //.

deciding on how many more bits is a constant size, per year in forward by prediction

SN

King & Cover (Entropy of printed English w/ 978 $\frac{bits}{letter}$). Using 12 "experiments" to predict next

Symbol in printed English: I expect to obtain better than others, but mixing 12 predicts $w_i = wts$. Got better than best expert!

So Q is, what wts to use? (As C shows that wts are obtained plus for are best expert but say pct of entire corp is divided by n & no. of experts.

A good way to find wts. Choose wts \rightarrow if they were used to predict each symbol of entire corpus, the product of probs would be max. Getting costs is a N.L. optm problem, solvable by Lev, Marg method. - possibly use small sample of corpus to get good approx. wts, then use entire corpus to "fine tune"

This is essentially a problem in linear prediction! - we have to optimize K costs, but if K is too big or just very large optimizing costs becomes very difficult. Perhaps the no of costs (or experts) could vary w. SSZ. (18.00)

SN

In devising a usable S-funct for string output: $P(A|Q)$: A is a string;

Somehow Q is made to create a grammar from which A is to be created. (In general,

finding such a grammar is not diff, because in Y. TSCQ, the A's will be of the

PC, so Ray will be easy to find by search (trying generate strings in PC order)

T. elements of it. Grammar can be a set of functions to be combined.

So O' is a function on Q 's that produces a set of strings that define a grammar

plus a set of params for that Grammar. The problem was: how to get these params? They have to be a function of Q . - I wrote about this diffy recently... at 24

Spec 17.00



o: 16.40: objection (of ideas): A poss. implement. TM categorizes Q_{10} (say) into a set of Q 's of f. Part. (Recognition functions of any part of ID's input). That set of Q 's all had to ~~same~~ same set of parameters & param. — words, so they would all be essentially f. same Q w. a D.F. over A^* . Not very General! (actually not bad! — see (10))

biology: All Q 's use f. some set of functions to be combined, but each Q generates a diff. set of p.c.'s for using process functions in the grammar.
Doesn't sound much good!

A better way to do (w. f. problem of s-functs w. string output & try some actual cases in which I have a ^{stochastic} soln. I want to implement.

One common case is in Baum Seq. Same Q , but limits no. of A 's. More often need a continuum of Real A 's.

For discrete case, I can postulate f. Grammar model & post hole, find param. but param. p.c. to comp. S. Say Q is a context (in a PPM): we have partial matches for various partial strings of Q : How to mix & prod. of various context lengths?

Perhaps find a way to do f. w. of (16.20-29) "Optimum" w. assignment.

What about a more interesting (complex) stock string output?

Consider text prediction: PPM's one way, but reasoning about what's going on, is another

What about predicting f. course of a story?

would Boolean Boltz Net (BBN) be useful? They are a natural problem I/O transducer. → 18.03

SN on the "Gram Mod Based Prog. evol." Stan McIlroy, Porter et al.

P7 eq (10):
$$\frac{P(\theta)}{P(\theta)P(\theta_n)} \cdot e^{\frac{c}{2}} \cdot \left(\frac{1}{12}\right)^{\frac{c}{2}}$$

$$P(\theta) = \frac{e^{-n} c^{-1}}{\prod_{i=1}^c \theta_i}$$

$$P(\theta) = \frac{1}{P_c(\alpha_1 \dots \alpha_c)} \prod_{i=1}^c \theta_i^{\alpha_i - 1}$$

$$\exp \left(\text{Eq (2)} \right) = \left(\prod_{i=1}^c \theta_i^{\alpha_i + n_i - \frac{1}{2}} \right)^{-1} \cdot P_c(\alpha_1 \dots \alpha_c) \cdot e^{\frac{c}{2}} \cdot \left(\frac{1}{12}\right)^{\frac{c}{2}}$$

Factor A & B look like Stirling Approx to Product of factorials divided by factors of sum.

$$\frac{\left(\sum n_i \right)^{\frac{1}{2}(c-1)} \prod_{i=1}^c \theta_i^{n_i + \frac{1}{2}}}{\prod_{i=1}^c \left(\frac{n_i}{n} \right)^{n_i} \cdot \prod_{i=1}^c \left(\frac{n_i}{n} \right)^{\frac{1}{2}}} = \frac{\left(\frac{n}{\sqrt{n}} \right)^{c-1}}{\left(\frac{\prod n_i}{n} \right)^{\frac{1}{2}}} \cdot \frac{\left(\frac{n}{\sqrt{n}} \right)^{c-1}}{n!}$$

$$\theta_i = \frac{n_i}{n}$$

$$\theta_i = \frac{n_i + \alpha_i}{n + \sum \alpha_i}$$

so should we put π instead of n here? so \sum & \prod formulas may be excise for constant factors

exp of MML is reciprocal of PC

so proper factor $\frac{1}{\left(\frac{n}{\sqrt{n}} \right)^{c-1}}$

for f. formula; if α_i all = 1 then formula is $\frac{(k-1)! \prod n_i!}{(k-1 + \sum n_i)!}$ if some $\alpha_i \neq 1$ then $\frac{\prod (\alpha_i + n_i)!}{\prod (\alpha_i + n_i)! \cdot \left(\sum (\alpha_i + n_i) - 1 \right)!}$

Work work out out!

STM

BBN.03 SM.12 Z.141,20 SVM.30

Spec
 30: 16.29: One random dirty approach: pick k best codes plus var. specify in k dimension of best code, will have a certain pc. Want to select $k \rightarrow$ the best pc of k -vector (with k in n space) times k multiplica of pc obtained, is max. So we try various k to see how much k mpc would

03: 17.19 **BBN**: One way to do BBN's: A particular Net of this sort will have a certain set of adjustable params. finding good values to codes particular corpus can be very difficult. A possible way: say k params. in net: pick random pts in k -space: for each, compute pc of corpus. Try many random k vectors, get P_k : keep 100 best k vectors in memory. After many trials, make predictions by using simple network product of 100 best k vectors.
 Will this work? How many random trials is needed before it begins to get good?
 Sounds like a "Wisdom of Crowds" problem!

12: **SM** 2.03 A unit would be a set of security, could be used to realize Cover's scheme: Or perhaps Cover's scheme w. finite (core exponential) memory. We write two forms of system by selecting k "best" of n different stocks in style of \dots !

0: Well, I wrote IAPSS which does $\dots \prod_{i=1}^n p_i^{a_i} d_i \dots$
 Lower base bias $p_i \geq 0, \sum p_i = 1$
 I got $\frac{(c-1)! T(n;!) }{(c-1 + \sum n_i)!}$ - Not far from what I expected.

It is similar to formula used in **STM FEB 2000** in Letter to Wolff, I also had an additional factor of $\frac{1}{c+1}$: ~~which was wrong~~ But they had to do with some other coding costs... I don't remember what, but $p_i = \frac{1}{c+1}$
 I'm uneasy about this: Consider "polynomial stream" expand $(p_1 p_2 \dots p_c)^{\sum n_i}$: p_i terms
 in value $\frac{(\sum n_i)!}{n_i!}$ if we integrate over one $\sum p_i = 1$ hyperplane, we get +20% of net by per plane, which is a function of c , but not of $\sum n_i$: If we look at the integrals divided terms, $\prod n_i! / c^{n_i}$, $(c - \sum n_i)!$ only partly cancels: $(c-1)!$ doesn't cancel, - if they are area of hyperplane, but c factors cancel

30 **SVM** They have try a separation using hyperplanes: If that doesn't work well, they use a "Randomish" non-linear transform into a dimension $>$ order of c or more hyperplane trial. Then they try hyperplane separation in larger dimension space. If it seems to hyperplane space is random, it costs nothing. If we do a random trial I select best one, it may be O.A., but we have SOY distribution. Taking average of several random n.l. features would give no. SOY;

ungstar also what kind of "Universality" it is.
 $\prod n_i! / c^{n_i}$, $(c - \sum n_i)!$ only partly cancels: $(c-1)!$ doesn't cancel, - if they are area of hyperplane, but c factors cancel
 $\prod n_i! / c^{n_i}$ vanishes; may shouldn't!

Rissanen: $\log_2^* n$ 100

Q: $2N \log_2^* n$ is approx for integers.

Actually using $\frac{1}{n}$ (which diverges) may be ok, and much simpler. We see only integers for den part of PC's, & this may be good enough!

T. Main Problems T. 3rd IAP Lecture deals with Alpha (Phases) & 2.
Also a ZEP system that supposed to $\frac{2 \text{ CPU}}{\text{sec}} \rightarrow$ Strong AI.

What are differences in these systems?

1) For Alpha: Preliminary exp on MTM type problems (MTR). Use of Lisp-like lang. & Lsearch. Using PPM or ~~subtree~~ augmented PPM, or S-CFG discovery for search. How good may rule be, is unclear. We may be able to "soup up" such alg so it would have some interesting results.

Perhaps to take it logic & train it to use logic in problem solving. (20)

2) Getting to S-induction in Phase 1 - is not so easy! I can do induction ~~from strings to nos.~~ from strings to nos. - a not bad, but not universal way. The output is a no. & we try to get how m.b. error. Ideally, if output would be an array p.d., but I haven't figured out a practical way to do it yet. (33)
The 3 input func is universal, but I have no good way to fit it to a comp. (33)

T. Problems HOW to Express/define/formalize in elements S-Functions

g: (11) For induction: Using Lsearch seems v.g., but perhaps use "partial ~~math~~" from Set of Problems & cand solves. It's certainly wasteful not to use PCs in it? We have been considering something like it for "problem pool" training (arts) in which we get partial P.B.

(By wt. in corp) auto basis of (how many problems, & how type) was obtained. In both d-induction & s-induction, a Exp Q was how to allocate "wt in corp" on f. basis of how many probs solved & which ones solved (in case of s-induction, how bad was failure & pc = 0 if a possy for certain of f. induction problems.)

One way to deal in P.B. is by finding ways of doing it that involve (a) of Exp PC (b) effective in the past. E.g. in s-induction, if a ~~prob~~ cand assigns ϕ to a particular correct A_i , we can ~~use~~ give it a default pc, or a pc that depends somehow, on ϕ or Q_i or some other "correlate".

33: (15) It would probably not be difficult to define a modality & a f.d. of all possl. P.P.'s on any range e.g. (0,1), (-oo,0), (1,oo). Thus, normally, if it's a unimodal d.f., there are only a few possy skewness skewness ~~and~~ & they are good enough (Gaussian, & Gamma). Modality could be (1) letter pairs (2) skewness (3) multimodality. So we could get "all possl. functional forms" by a lowly array (normal) functions - but actually, only a few forms will get much use & \therefore have very by arrays. So don't worry about universal from universality, for now - this looks good enough for now (99cc. 20.00)

0: (space) (19.40) : Can I genz. this idea to strings (discrete)? Our way! we have several distance

NB Implementing this distance could done by the "R" input of the 3 IUM model of induction. This also applies to approximation.

- measures better. strings that can measure "error" of 2 particular AD production.
- 12. ① Information for fan: Universal, correct, ^{very} expensive ② Same as ① but using 2 non-universes (best of X MAS (as a universal test w. heavy wt. on an important subset of insts) in the style of New-Simon GFS: TM learns how to do it. learn ~~how~~ having least best in available insts)
- ③ Edit distance (I don't know how expensive this is out. even so) ④ Affinity distance (cheap but not r.p.).

06

② (0.25) could include any "invented" distance measures appropriate to domain

Question Changes.

- 3) continued from 19.19. We have what seem to be kinds of Goves for such 1
- a) Lsach b) BU such. In Lsach, d-induction is simple (usually) in general, but in PC order - The guiding p.d. is modified by past success. Hw, in invar. 19.20, we may want to vary to suit of w. given to previous "successful" d. counts - - depending on how many probs. they successfully solved. The most serious problem is updating of the guiding p.d.
- I expect to use either S. CFG's or PPM. Modified to detect subtrees as context.
- Since I'm not ~~using~~ really using Lsach (I'm not doing inv. probs & I'm not doing t. Lsach version for OZ probs), I don't expect a very good induction (until I get to phase 2).

19

A poss. way to use BU such for d-induction: Use "problem pool". BU such tries to find a sym. subproblem as a subproblem in pool as poss., using 19.20 as "feedback" info. → 20

→ SN **SN** Mite it not be poss. (practical) to use Lsach for OZ probs (in particular for searching for good induction hiddls)? The search is over option techniques, & TM starts out w. a grammar that can readily construct ^{good} option techniques. (predefined set of ~~the~~ option tqs.)

set OZT = option technique

This would (perhaps) be very similar to starting TM work on "Phase 2 first"

SN I usually think of ~~the~~ OZT's as being independently tested - but if OZT's work by winning trials, then they are not at all independent, since they can use each other's trial info. Hw, a v.g. induction system could locate a bunch of trials & suggest a trial w. very high expected G. & so in this way, a v.g. induction system could do much v.g. "Self-Improvement"

0: (19.5)

On Phase 1 s-induction: We can use 2 part codes, both for numeric and for string output (0.00-06). In "Lsach" ^{we do} trials in order of PC of the discrete "first part" of the codes. For numeric outputs, ~~we do~~ for each "first part" trial, we have to optimize parameters & get midth of optimum. If there are several good optima, we have to get value & width for each & "sum" them. It will be better ^{under certain condns...} w/rid with work time using while to do random params... (but I vaguely remember this not working!). And the trick is to use a small SSZ & find broad peak, then SSZ & narrow down peak. There is probably an optimum SSZ ratio to SSZ.

PCA(a) v.s. PCA(a).

For induction w. ~~the~~ string output, → SSZ 21.00 → 21.00

20:20:40: S-induction w. strong input: For "Lsueh": 2 part codes: do first part in pc order, ~~then~~ which will have a strong output: If it fails is not for degraded fit, use 20.00.06 to get to correct fit. This can be done consuming a wa into for using analogs off: 400000 output of 20.30:40.

T. for: were born for "Lsueh": For ISU such, we start in 1 or more cards that give ~~pe > d~~ pe > d. to corpus, but are usually rather poor. We then use a kind of "filling in" by induction on corpus of ~~equations~~ ~~and~~ their Govc. values (\equiv pc's). There are many many other techniques of this sort in the ~~Math.~~ / Europe Community, & a sizable subset in the Mach. Ling. Community. I am most familiar w. feed-forward ANN's & GP/GA. — But there are other impl. approaches: SVM, BBN, Decision trees, ...

$P(A|Q)$ v.s.
 $P(A,Q)$.

10 (See Dada, Hart, Stone for more ways). While recurrent ANN's are universal in Turing sense, feed-forward (i-type) are familiar w. do not — So I'm doing GP: I think BBN is about "Turing Universal". Decision trees probably isn't. SVM's I don't know, but they would seem to be universal (only in sense that feed forward ANN's are "universal").

"Inductive Logic Programming" may be ~~the~~ Turing Universal. } Are there any other T. universal systems?
GP/GA
Recurrent ANN

20 SN: 19.03.05 Here I ~~have~~ ^{have} a bunch of parts & I'm trying to put them together to make a TM, & I'm seeing if the parts don't work right, or if they don't fit ~~together~~ together, or if I can fix them or get different parts. So far, seems to be going along ok: Also really ~~bad~~ ^(problems) ~~bad~~ (troubles) found!

19.03.05 is the beginning of a parts list: Perhaps a quick ~~runover~~ runover what I think a complete component is: "parts" ~~list~~ ^{list} ~~list~~.

1) Problems for QATSR: (2 approaches) for expression of QA prob & solns. Use of "Lisp": forms $P(Q, A)$ or $P(A|Q)$. both possi. for MTM $\ominus A = F_3(Q_3)$ is available.

30 2) For d-induction on QA probes / MTM: Use Lsueh: Guiding P.O. PPM, PPM + substructure discovery ~~for~~ ~~discovery~~, Other things ~~are~~ ~~CSG~~ CSG discovery.

3.2k mi ~ $\frac{60}{50} = 500$
~~500~~ $5 \frac{2}{3} \frac{3}{12}$

∴ A possl. Intro to talks! We have this program that can look at all of its past experience & make predictions/judgements based on ~~that~~ experience in the best possible way. It could bring all of this experience to bear on any problem given to it.

We have what appears to be a ~~simple~~ way to ~~construct~~ write such a program. We have a program that could, indeed, learn optimally from its past in the best possible way — but there is a catch — the program takes too much time and memory.

We will describe a sequence of developments that ~~enable us to~~ ~~overcome~~ overcome this difficulty.

~~We start out with a program designed to be prediction only.~~

First a description of the system, prediction system that is very accurate, very economical of data — but is much too expensive to run.

Next explain how ALP works & what its problems are.

At a certain point in the ~~ALP~~ talk, discuss universality; feed forward ANNs are "universal" in a certain sense: Describe just how they are universal.

Then explain how they are not "fully universal" but there is a very narrow category of things they can learn & a very broad category of things they cannot learn.

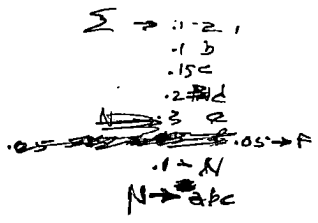
Intended to find examples of things that a universal feed forward ANN cannot learn.

Well, we can regard it as a decn. (arith.) problem. To use an ANN ~~to solve a certain~~ ~~kind of~~ ~~problem~~ ~~is~~ ~~to~~ ~~do~~ ~~decn.~~ ~~functional~~ ~~form~~. Using a universal machine, we need more ~~than~~ bits to decn. functional form. It has already occurred to us that partially occurred before. A sine wave can be coded by 2 samples by a sine — no matter how many cycles are given. If, ANN needs a no. of bits or ~~no. of~~ ~~cycles~~.

2.8.05
STM

David
Lindsay Room 131 (Unit 123)
david@cs.rnu.ac.uk
or D.G.Lindsay " " " "

CFG diffy! In analysing Tree Grammars, I got into a serious problem, which I considered that there were very many different ways to parse a single string or set of strings in a simple grammar:



i.e. Enumeration If N occurred n times in corpus, there were 2^n different ways one could parse that corpus. —
assigning either (N or directly abc) — Thus if particular words had a much higher probability any one parsing would give to it.

11:52 P Δ3
2:56 A Δ3
6:31 A Δ3±
12:18 A
4:12 4 to
6:14 2

2.20.05

22.00 off is inf. rta direction for Invo to first talki But I may want to First Define "Strong A.I.", then How could we achieve this? (Defn. Exptly of ALL)
What are the preconditions for such an approach?

1) Defn. of S.A.I: Ability to synthesize S.I. on hardware. This plus "Moore Law" gives extremely capable A.I. in finite ("short") time. In

D.K. on way to
Toward History-based Grammars:
E. Black
P. DeLima
J. Lafferty
et al.
≡ "HGB"

One test of Turing Universality is Existence of an "Undecidability" or "Halting Problem".
A better positive Criterion! That any Universal Modeling System can Simulate any computer in finite time with only a finite amount of info (logically).

on convergence of $\sum PCU$:
Say Naiman constant = $A > 1$. $A = \frac{1}{\prod(1-u_i)}$

$$\ln A = -\sum \ln(1-u_i)$$

$$\ln(1-x) = -x - \frac{x^2}{2} - \frac{x^3}{3} - \dots$$

$$\ln' = \frac{-x}{1-x} \quad \ln'' = \frac{-x}{(1-x)^2} \quad \ln''' = \frac{-2}{(1-x)^3}$$

$$\ln(1-x) = -x + x^2 - 2x^3 + 6x^4 \dots$$

$$\sum \ln(1-x_i) = -\sum x_i + \sum x_i^2 - \dots$$

$$\sum \ln(1-x_i) = -\ln A = -\sum x_i + \sum x_i^2 - \sum x_i^3$$

$$\sum x_i^2 = \ln A + \sum x_i^3 \quad \sum x_i^3 < \text{say } \ln A$$

$$\text{So } \sum x_i < \approx 2 \ln A$$

I need to see just what I want in a scheme to detect proper copies
in functions (Lisp & /o forth & /o machine lang). CFG's is augmented with a set of systems
I need to see how Good Reason. This seems like a logic problem to me, and to be understood!
→ 24.00 Spec.

30

32

10

on comparing prodn models that have different nos. of ^{continuous} params: n_0, n_1 ; $n_0 < n_1$

One way: Increase n : no. of params in t : smaller param. model ^{with n_1 params}

Then, to compare ≥ 2 models integrate each over t : n_i dim space.

This usually means that $\int |dim = 1$ for each dimension. We end up w. a coding cost that is \int integral over all of t : n_i dim space — for $i = 0$ or 1 .

So this is nothing new: We have to get an a priori density function on param. space.

~~We can't assume that~~ How to do this is still unclear: t : different params could be

of diff. "dimension" & \therefore not comparable. In ordinary linear regression n is

sort of comparable, and " $n \pm 1.0$ ". When n is not true, we could do it!

Unclear as to what to do.

IMPT! If n params are really different, the only way to get a reasonable a priori for t is to look at "similar" prodn probs in t : past

" n to sub-observations in past" may be via complicated functions, not — or rather not!

This has to be analysed in greater detail: ~~It's~~ very impt!

20

(SN) It may be easy to store t : equiv. of several (1) CFG's in a W.B. xline (PRM).

Very simple & fast } A given context could occur in a low PC grammar, but have very high PC within Prod grammar. \therefore product of these 2 PCs could be $>$ Prod t : corresponding sp's in 2 / more likely grammar? Story of Models of comp.

30

On Normalized of Stoc. Comp-like AEP model. We could use unnormalized form. PC ratios would all approach correct values, ~~but~~ I do do a new way much slower, but.

2-11-05
STM

MAJOR HEURISTIC .06-.09

Quick About: New View? .11-.17

α1-α14 follows
STM 25

25

Then STM 26 follows.

So α1 → (STM) 25α1
= 2 → 25α2..

0: 23.40 : 044. Q of 23.32 : Ideally, we want TM to be able to look at a bunch of function trees, find ~~some~~ a set of common sub-trees and then decide which definitions would be best: (considering overlaps of sub-trees). Criteria for (actual code lengths) will have to be defined/calculated/evaluated.

36 2-17-05 Hm, A MAJOR (DEFAULT) HEURISTIC: After solving a problem (w. Diffy):
To try to find heuristic rule(s) that soln. would be faster. Occasionally,
→ 1 heur can be gleaned from a single soln. How to put this Heur in more usable form --- UNCLEAR. 2 KINDS of Heur: (a) ↑ of PC of soln. (use "standard" methods for finding Regps in corpus) (b) Quick About techniques. — I think I had generalized this to cover a large range of such heurs. One (not so general, but perhaps very useful) — as a cand. is being generated, develop ~~some~~ obs to predict ↑ PC of failure (Unclear that this is not covered by ~~normal~~ assignment of low app to cand. to be "Aborted"). — See 26.06 for analysis of "Quick About" that doesn't seem to be covered.
17 → ?? Could it be that the PC method of .15-.16 covers most (if not all) Quick About methods?
In which cases, finding good productive assigning good PC's to cand. (i.e. parts of cand.) is all that one has to do in Phase 1.

SN Trouble w. Bernoulli Seq. formula! say 2-symbols only: If no φ's & n₁ 1's have occurred.
Prob (PC of 1/PC of 0) (for next symbol) is $\frac{n_1+1}{n_0+1}$ ← this is ratio of probs for n₁, n₀.
P_{1,3} → n₁ → n₁+1; n₀ → n₀+0. These probs were obtained by direct integration!
 $\int_0^1 P^{n_0} (1-P)^{n_1} dP = \frac{n_0! n_1!}{(n_0+n_1+1)!}$ Oh! its ok! ☺

2-11-05
- 5 IAP
- 5 TM
②

STM

25 of 1

This paper follows STM 25

Re: Paul Vitányi: I've known him for many years: Since he and Peng Li

began to think about writing "An Introduction to Kolmogorov Complexity and its Applications". I've come to
ask to me, ~~as~~ ^{as} ~~one~~ ^{of} the founders of the field of ~~the~~ program-complexity.
- Preparatory to their writing "An Introduction to Kolmogorov Complexity and its Applications".

His influence

His main influence on me, has been to force me to express in more exact form, ideas that I had been a bit vague about.

His influence on this field has ~~been~~ ^{been} enormous. At the beginning it was Peter's application of ~~ideas~~ ^{ideas} about ~~randomness~~ ^{randomness} to a variety of proofs in Mathematics. Later, his influence was largely due to a flood of his book on Kolmogorov Complexity — it reviews much of the work in this field and ~~is~~ ^{is} continues to be the only ~~text~~ ^{text} available textbook in this area. In addition the book has been an important ~~stimulus~~ ^{stimulus} for continued research by others.

As a person, I have been much impressed by his ~~curiosity~~ ^{curiosity} and erudition in a great variety of fields — in ~~mathematics~~ ^{mathematics}, Science, and other Arts. He is a prodigious world traveler. Likes to meet new and interesting people ... taste new and interesting foods.

See those cheap very fast SONY graphic CPU's - used in PS2 or PS3

Very Imp't Idea! Try to see what kinds of functional Registry

Occur in Machine Code w/o forth or perhaps a specially restricted superputer
Machine. I did have this idea (long ago) of a Machine w. many addresses being sent
around. A few (perhaps 4) function arguments. Addresses and ~~data~~ ^{data} used to
get processed. Perhaps work out General problem of Registry in Computer
Machines. 2.12.05 5 IAP STM 1.00 per decm. of T. Address Machine.

(SN) while SW. patents would slow innovation, it would probably not slow approach to superlativity (much),

we will have Lisp 30y, plus a complete set of Rice insts. So gradually TM would learn to use + Rice insts.

One trick: when one wants to declare a variable, a special symbol follows to "data" symbol, ^{enables} ~~enables~~ special insts that make var decls poss.

2.12.05
5AIP
STM } mix
STM
⑧

T. ADDRESS Machine (TAM)

25x2

T. address machine (TAM): Each address holds name of the sub-function to be evaluated. If keeps track of where we are in the sequence, & knows when we are done. T. addresses contain (or point to) each sub-function evaln. When (few) addresses need a reverse in function evaln., duplicate evaluators are constructed by a hierarchy system.

11 - 2-2
2495
5855 3/6

In (Lecture) ANL: T. notation: $(I \oplus \text{Then} - \text{Push}) + \text{Add}$. would be useful. means "I + push" but "I, ob" is just a word w. "Add" + "Op". In general, we may want to modify the language syntax so say "I + Then" is one word, so push is ~~push~~ (push) push & a prefix of I, Add. ~~push~~ + I + Then, automatically push next symbol. It may be that (I + Then push) should be a primitive. So I + push + I, Add means I + Then pop & push Add. So we may want to modify syntax so as to make PPM a better way of being more visible to PPM. \rightarrow a 4.00 spec.

M. glickstein@ucrl.ac.uk
17R
Actually, there are at least 2 ways to do it. Grammar Corpus in learning of Grammar. \rightarrow x3.00

SN

Discuss "Grammar Grammar" in Derbug how to assign a prior to it.

Grammar: t. discussion. I have may be O.G. ! But (what?) - E. Grammar

Grammar can end up w. better (years) for the entire corpus. better for starting out - to "small" increments of corpus.

SN

Intro Lecture!

what is my goal in S. AI? (Terry test as Psych. info about man)

Requisit University - (Must be Pliers in Barrel)

- 1) ~~Easy~~ stuff poor understanding of Learning probability, induction
- 2) ~~Use of~~ ~~incomplete~~ ~~unlabeled~~ ~~induction~~. ~~Not a full~~ ~~backwards~~ - Law of
- 3) poor T3Q5
- 4)

1985 votes & more evidence.
① NIPS comparison
② AI Sub-culture on Uncertainty in AI.

00

0

16

17

20

20

STM

Expected Value.

Mean of $PCC \frac{1}{n}$

if $p = \Delta \alpha^{-n}$ then prob of extra NT is always α & usually $\frac{1}{2}$.

10: & 2:17 : In deriving the set of NT's, we need to prep of an integer n ; $\frac{1}{n}$ is actually not bad, because we (probably) don't have to normalize. It does, however, give lots of trouble.

to large n values. ~~$\frac{1}{n} = \alpha$~~ : very fast divergence. $\frac{1}{2^{log_2 n}}$ also has $\frac{1}{n}$ divergence rapidly. I imp'tly say $\frac{P(n)}{P(n)} = \frac{n}{n+1}$ close to 1

In f. case of $P(n) = \Delta \alpha^n$ ($\Delta = \text{arbitrary const}$) $\frac{P(n)}{P(n)} = \alpha^{-1}$ which is usually close to 1 (I have computation for $E(n)$ in terms of α within last wk. or so. of STM notes)

So usually (probably if $\alpha > \frac{1}{2}$, say) $\frac{P(n)}{P(n)}$ being close to 1, will give us essentially no central bias for another NT! Here an additional prob is that of deriving the NT

SN

Could it be relevant to linear regression (MLC) problem? For this not. Just saying we need an extra bit, isn't where the PE lies - its

in the narrowness of the coefficient D.F. - which .00-.08 doesn't help with.

3: .08

Another way to see it is in terms of the no. of terminals + non terminals. $\frac{1}{n+1} \approx \frac{1}{e}$ for k references to symbols. This connects to ~~...~~ #28

This "small \downarrow in pe " again occurs when we add a "production", or when we \downarrow no. of productions (Merge). This is because tallying the no. of productions for each NT requires statement of an integer. So in both cases \uparrow or \downarrow by 1 of no. of NT's or

\uparrow or \downarrow by 1 of no. of productions for $\geq pc$, will be multi. of p by a no. "close to 1".

The second "punctuated symbol" is the one that marks the end of a production string. I think ~~...~~ again "dpc - 1".

To further kind is "end of list of production strings", a P's also has a "dpc" ratio of close to 1.

I think the moral is, we can use $\frac{1}{n}$ or 1 or α for any of these quantities. Saves a lot of bookkeeping & calcu.!

30

It looks like the most complicated part (by far) of calculating pe is a grammar mod. is involved w. ambiguity - multiple parse.

Probably it would be best to (keep) track the sequence of production strings as a single string w. punct. symbols betw. production strings. - It would seem probable ~~...~~ symbol (E. or N.T.) appear w. by freq. in one production, it would have by freq. in all productions. I could look at this α experimentally

- I really haven't ever tried any of this on real data.

STM

SAAR-Broden
SARB

20 : α, β — In related contexts At SAARB in ANL # using stack: P20 trees
 were generated by a CFL: So $\langle \text{if} \dots \rangle$ was chosen, we got
 IF α Then β else γ : $\langle \text{if} \dots \rangle$ was a certain type, which would restrict them
 considerably. Also, if I had a pd of α, β, γ based on past successful
 functions, P20 would be very useful. So $\langle \text{if}, \text{Then}, \text{else} \rangle$ would have its own
 ppm" BW xfm. We'd like to context to go further if poss. — e.g. 1.
 context of "if then else" could be considered.

To what extent can I represent all that I want by a (Σ) CFG?

I.e. If a certain Context occurs w. enuf freq., then I declare it ~~is~~
 a technique frequency.

One trouble w. "defining" is that one tends to parse only 1 way, so one
 loses other (possibly better) solns. By referring to "Ambiguity" of a CFG,
 one could get around that.

13-15 gets around an objection that I had to "defining" i.e. "Freeze"
 to parsing. In fact, when ever a bunch of new data. have been
 obtained, one should periodically "reparse + past" (like Wolff) \rightarrow

Actually, humans don't reparse (involunt new concs) very often: — Mainly
 because its time consuming, diff., is often unpleasant. Tho it can be very
 pleasant, indeed, when a much better parse (\equiv formulation \equiv code) for
 the past is found.

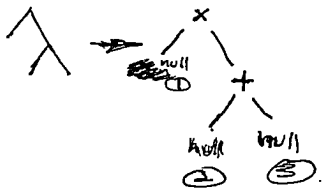
Some effects of reparsing: Reassignment of (perhaps many) different
 pc's to various concs, ~~the data~~ IRTPC of entire corpus is much f,

— we can expect much better prodn. in future!

On deriving simple (non-looping) Grammar, using search
 sub-trees only: It is easy to define a sub-tree: we just give
 its top node, then for each branch, we have usually a finite N.T. or T.
 Here, we add a.s a terminal, + null — which makes ~~the~~ best point
 an argument of the resultant subtree. See 2.5.01 \rightarrow

SNL on used on normal
 re for QA lang;
 This will be forward.
 forests that are total
 recursive.

5TM



This gives the subtree $x +$, w. 3 arguments:

$$x(arg_1, + (arg_2, arg_3))$$

One can define a subtree, it becomes a N.T. as can occur as arg to any function. Htr, we may want to keep track of ~~the~~ trees of errors (N.T.'s & T's) at each "context" of 1 or more nodes.

This "context" idea is certainly impl. — but unclear as to how to implement it!

Anyway, this "finding of common sub-expressions" is just "dynamic programming" — no magic. The grammar itself is loop free. Because of its non-occurrence (\equiv no loops) it is very large. Mindful of PPM in that respect.

I did write about how (in theory, but not really practically) how to use ~~the~~ sub-expressions as contexts. If I could, indeed, implement that ~~is~~ w. ~~the~~ keeping sub-expressions as distinct objects, I would get a ~~very~~ fairly powerful function predictor.

Perhaps if I did define sub-expressions, using distinct sub-expressions as contexts in an ordinary PPM, it would work O.K.!

$$\int_0^{\infty} z^n e^{-\alpha z} dz = \frac{n!}{\alpha^{n+1}} \quad e^{-\alpha z} = \sum_{j=0}^{\infty} \frac{(-\alpha z)^j}{j!} \quad \text{I that this finite sum will be}$$

useful for summing effects of all parses for data of α in defining α as a set of primitive symbols, but the formula isn't quite so simple. I think, htr, that there is a corresponding approx. that does give a well-known function (or may be close to it).

$$\frac{\prod (n_i - r)!}{\left(\sum (1 + n_i - r) - 1 \right)!} \cdot \frac{\left(\sum (1 + n_i - r) - r \right)!}{\prod (n_i - r)!} \quad \text{But this is cubitarily! — 1.2 factors probably don't exactly cancel.}$$

Also, this r assumes that all of the primitive symbols are used in α (i.e. newly defined N.T.) once.

The exact no. of symbols of each type is the no. of each symbol in the original corpus determining the range of summation of "r".

Perhaps easiest to do this as an integral — ... Some kind of modification of the found Beta function.

SN

Re: Discussion of Apri's "Subjectivity" in Lecture 1. I may want to expand

Re: I did write how the a priori person changes during his life, but I didn't make it clear how this was relevant to anything!

My Expos more exactly: T. Apriori at any particular point in a person's life, is based on the sum of all of his experiences plus the info built into him at birth.

Say this means: info at birth plus ...

T. Apri itself, is a kind of statistical summary of all of his info. Needless to say, this apri is a different, the data on which it is based, will differ much from person to person.

For those that feel science is ~~of necessity~~ must, of necessity be purely objective — ^{would seem to be} — The Re is, indeed, a lot more of his a priori that just about all humans have in common, there is much that is highly individualized.

How ^{can we} deal with this difficulty? — I will divert to a closely related problem: How statisticians deal with a priori information? .?

A common way is to select a set of hypotheses to consider, then find which one best fits the new (non-a priori) data. For example, in linear regression, we

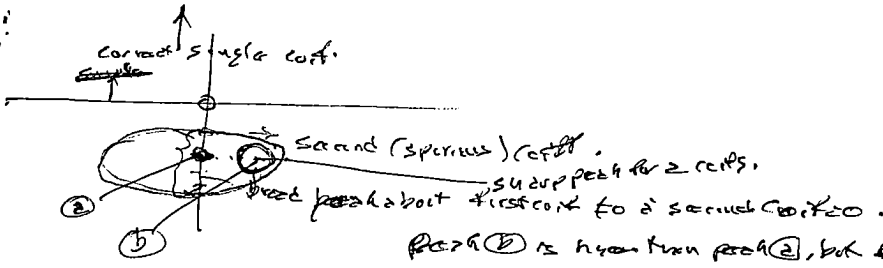
decide a priori, to use ^{linear} ~~use~~ β coefficients to predict the next element of the ~~series~~ ^{series}. The coefficients can be assumed to have a uniform a priori distribution. (We will gloss over certain non-orthogonal difficulties with this approach).

The statistician says "If we make these assumptions about the models for the data, then we will obtain predictions such and such predictions with such and such variance. These are rather objective statements."

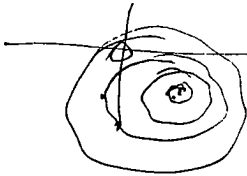
i.e. apri do not converge
 $\sum_{k=1}^{\infty} 1/k^2 = \pi^2/6$
- but they may not converge to the same value
Re is uncorrelated
no relationship across
of parameters so we
can get rel. w/ty of
various hypotheses.

SN AH! Consider β coeffs for linear regression! We get a sharp distribn of coeffs then in 1. subspace of 5 coeffs (w. 1 coeff = 0), we get a broader dist. w. much more total var. but its post predns about as good.

It may be that if we averaged over the entire β dim space, in the generation of β seq. was only 5 coeffs, then we would still get terrible for the sets of coeffs. It may be that to peak for 6 coeffs it will be narrower but if the 6th coeff zero, then peak is very ^{broad} ~~narrow~~. Essential visualization of β was actually 6 coeffs but we consider a space of 5 coeffs.



What would the ms error contours look like?



It would be worth while to do a computerized study of (the ~~relative~~ coil on the generator) i.e. considerations of i.e. 2 dim space of 2 coils. What do the contours look like?

Actually, the different coils could be "equivalent" in part in each direction. If we were doing linear regression on a 2 channel system. But if each linear coil is entirely different (copy of different dimensions), this is not true, it fits unclear as to how to

compare with the first different coils - in terms of info content. One way would be to look at the past & see what relative weights were of what relative importance.

In the case of two series (linear or non-linear) we can look at the past of the same timeseries to see what relative weights should be used.

In specifying coils of the most general kind: One kind of info to specify is "center of distribution" - another kind to specify distance from the center.

The first kind is needed for noise prediction. The second kind holds most of the info for (noise) as by precision product. It's usually the second kind (always SSE) that statisticians are most concerned w.

A unitary matrix method for coils of different dimensions is the partial derivatives of the ~~error~~ (mse) w.r.t each coil. In general we have a Hessian matrix. If off-diag terms are $\neq 0$, then the coils are not indep. i.e. "to use first" one. If info in a particular coil is independent!

So how to compare Hessians of different nos. of dimensions? Since Hessian is symmetric, it is diagonalizable by rotation: so \exists linear combs of coils that lead to indep. "coils" w. distance $\frac{2 \text{ trace}}{2 \times 2}$ for each "coil". So one inserts info in each coil by amount necessary to ~~minimize~~ \downarrow ms error by same amount. After the coils have each gotten to the pt. of a contribution to error, then each coil will transfer to that same no. of bits of info.

So, initially (as perhaps part of the "discrete" form of the model) one gets the coils down to processors, & each error in coil contributes equally to error. Then however - some processors in each direction contribute of (bits of precision)

00

0

20

30

31

spec
 01 $\alpha 7.31$: but not in terms of simple error in each "coeff" : This gets us out of $\ln n$ in \ln "Cost"
 form of info. I guess assumed cost $\propto \frac{1}{\epsilon}$, which seems like $\frac{1}{\epsilon}$. By problem is to max Cost not $\propto \ln n$: Th. discre
 So $\alpha 7.31$ may be great break thru ! STM 36.01 Cost not $\propto \ln n$: Th. discre
 part.

Hvr, to start off, the ~~the~~ determinant of the Hessian matrix gives r. info
 ($\rho \epsilon^{-1}$) of the "peak". Further info is \propto to no. of dimensions, m & space.

In general, the amt. of precision needed to specify the peak will depend on SSZ .

For large SSZ , the info needed is $\propto \ln SSZ$: - in cost for SSZ \propto to no. of dim. Hvr, for small SSZ info needed is more complex, it depends on much app info.

2 (AM)
 ~ 3z
 4:30 AM
 7:47 3:17

To compare k coeffs w. $k+1$ coeffs: Compare best of worst
 determinants of Hessians of $k \geq k+1$ coeffs. This will give values $\rho \epsilon^{-1}$ to start off
 for high precision, equal amt. of info will be in used in each dimension.
 So it may well be that $\rho \epsilon$ at peak divided by Determinant of Hessian is
invariant (or maybe $\sqrt{|Hess|}$?) is measure of $\rho \epsilon$ at best peak.

N.B. In adding up $\rho \epsilon$ of various pts in coeff space, I was using hypercubes
 w. same edgesize in each direction. $\rho \epsilon$ is average. Use as edges,
then as eigen vectors of the Hessian is more or less small hyper parallelepipeds
 w. edges \propto (the basis vectors) (or perhaps the eigen vectors). This automatically
 makes all directions in space equivalent.

I still don't have a good inductive grasp of what's going on, hvr.!

I'll have to keep going over $\rho \epsilon$ in various ways until it becomes inductive
 clear. (or false)

Still, the specification of the Hessian as part of the discrete deriv of ϵ is correct "... ?! @? "

Say one changes dimensions (much to fit in one off coeffs) - Hessian is
 most by 12. (or 12^2 ?)

Majorant was (each each dim. as being deriv (or derivatively), the
 dimension is just a "small" addition (discrete part of deriv) part. - yet this
 discrete additive part is very imp. for small SSZ . Hessian (linearization) is easy to compute!

$$S = \sum_{j=1}^n \text{sq. err} = \sum_{j=1}^n \left(X_{j+1} - \sum_{i=0}^{k-1} X_{j-i} \cdot 2^i \right)^2 : \frac{\partial S}{\partial a_k}$$

(SAS) ① Use of $\frac{1}{2}$ APL for four. (Good Arrows functions; but sure
 version is pass.) but be sure recursion is poss.

② for $\frac{1}{2}$ use \pm very easy (Agnos. (\pm is random))

$$\frac{\partial S}{\partial a_k} = 2 \sum_{j=1}^n \left(X_{j+1} - \sum_{i=0}^{k-1} X_{j-i} \cdot 2^i \right) \cdot (-2^k)$$

Diags elements "not", but "is". I got my result some time ago when working on "G" $\propto \frac{n+m}{n-m}$
 They differ much when k was close to \log or $k \approx \frac{1}{2} n$.

HA! It may be that "Fisher Info" is relevant: It has to do w. derivative of ln of something...
 which seems in direction of what I want.

STW

4:32 min.
back 5:00p

$$\sqrt{x^2 \cdot y^2 \cdot z^2} = xyz$$

Perhaps a Good Way to think about α, β, γ off; Think of $SSZ=1$; $N=1$; $R=3$
 Entries + Hessian is the discrete part of the denom. The rest of the denom is just "normalizing down to distribute."

SN As an example to test various/linear prediction methods: Consider the time series with 2n

of obs: T. char. eq. is $\prod_{i=0}^n (x - a_i)$ maybe $i=1$ to ∞ $\Rightarrow |a_i| < 1$ \rightarrow Root a_i are random: once chosen, they are fixed.
 a_i is random no. chosen uniformly, value $\neq -1 \pm 1$. λ is like to expected no. of coins.

T. roots of char eq. are all in "unit circle".
 Say $\lambda = 3$: for small n ($\approx SSZ$) we will only be able to resolve a few coins, because a smaller λ \Rightarrow 1st coin will be very small if $\lambda \gg 3$. As n increases to resolve smaller & smaller coins.

Another possibility $a_i = \frac{1}{\lambda} \cdot (-1)^i \cdot \omega$ $a_i = \frac{1}{\lambda} \omega^i$
 Here, if $\lambda = 0$ $a_i = 1$ which is not so good \rightarrow $\lambda = 0.4$ will be bad for $\omega = 0$ (usually).
 \Rightarrow I just pick $(\frac{1}{\lambda})^2$ in as a factor.

Anyway! Re: to "Hessian" form! for linear regression, it looks like I may end up w. something close to $\frac{n+k}{n-k} \sigma^2$ obs. Just use asymptotic for $k \ll n$; (finding mean is tedious, messy).
 Then do a dimensional analog, via s. diagonalized Hessian.

So, to compare num. of info in 6 v.s. 5 coin models: for 5 coins: we have constant containing $\frac{\text{Denom of model}}{\text{num of model}}$ + shape of Hessian that $\frac{5}{2} \log_2 n$ bits to desc ω pt. of interest + $n \log \sigma$
 for 6 coin $\frac{\text{Denom of model}}{\alpha^1}$ + shape of Hessian $\frac{6}{2} \log_2 n$ + $n \log \sigma^6$
 $\alpha \approx \alpha'$ maybe \approx I don't know about β v.s. β' : They may differ by < 1 bit(?)
 Also, both models have to be built by pc. of corpus wrt that model. Since there are just n points, it will be $\propto \sigma^{-k}$ v.s. $\propto (\sigma^6)^{-n}$ in 1. 2 cases.

For logon to 2 info forms $\approx \frac{5}{2} \log_2 n + n \log \sigma$ v.s. $\frac{6}{2} \log_2 n + n \log \sigma^6$.

for k coins $\frac{k}{2} \ln n + n \ln \sigma_k$ v.s. $\frac{k'}{2} \ln n + n \ln \sigma_{k'}$

or $k \ln n + n \ln \sigma_k^2$ v.s. $k' \ln n + n \ln \sigma_{k'}^2$. T. comparison of 2 Hessians should also be made: $\frac{\text{Denom}}{\text{num}}$ how to do this. or $n^k \cdot (\sigma^2)^{-n}$ or $n^{\frac{k}{2}} \cdot \sigma^k$

Look at $\left(\frac{P(\vec{w}_0)}{H(\vec{w}_0)} \right)^{\frac{1}{2}}$ as width of distrib; $P(\vec{w}_0)$ is pc. of corpus wrt. that point in \vec{w}_0 space.

$|Hessian| = 2k \cdot \sigma^k$

$P(\vec{w}_0) \propto \sigma^{-k}$

$P^{\frac{1}{2}} \cdot H^{-\frac{1}{2}} = \frac{\sigma^{\frac{3}{2}k}}{2 \sqrt{k} \cdot \sigma^k} = \frac{\sigma^{\frac{3}{2}k-1}}{2 \sqrt{k}}$

$\frac{|H| \propto (\sigma^2)^k}{P^k} = \int \frac{\sigma^{-kn}}{\sigma^{2k}} = \int \sigma^{-k(n+2)}$

$= \int \sigma^{-k(n+2)} = \sigma^{-k(n+2)}$

$\rightarrow \frac{\sigma^{3n-2}}{4k}$

T. stuff about Hessian in last P of Lecture 5 is wrong.

00

(H) $\neq \sigma^{2k}$. \leftarrow This would be a loss of a total ~~sq.~~ error (perhaps)

but what if I'm interested in the f. function of σ σ^{-n} - f.c. of the corpus with f. model (?)

I may know $\frac{\partial^2 \sigma^{2k}}{\partial (\text{any param})^2} \approx n \sigma^{2k}$ $\Rightarrow \frac{\partial^2 \sigma^2}{2(\text{param})^2} \approx \sigma^2$

If $\frac{\partial \sigma^2}{2x^2} = \alpha$ then $\frac{\partial^2 \sigma^4}{2x^2} \Rightarrow \frac{d\sigma^2}{dx} = 2\sigma \frac{d\sigma}{dx}$ $\frac{d\sigma^4}{dx} = 4\sigma^3 \frac{d\sigma}{dx}$
 $\frac{\partial^2 \sigma^2}{4x^2} = \frac{d}{dx} \left(\frac{d\sigma^2}{dx} \right) = \frac{d}{dx} \left(2\sigma \frac{d\sigma}{dx} \right) = 2 \left(\frac{d\sigma}{dx} \right)^2 + 2\sigma \frac{d^2\sigma}{dx^2}$ $\frac{\partial^2 \sigma^4}{dx^2} = n(n-1)\sigma^{n-2} \left(\frac{d\sigma}{dx} \right)^2 + n\sigma^{n-1} \frac{d^2\sigma}{dx^2}$

This may work out ok - but drop it for f. value: Check to work on 4. Lectures (43)

I think the correct formula for width σ is

$\left(\frac{p^k}{|H|} \right)^{\frac{1}{2}} \cdot p = |H|^{-\frac{1}{2}} \cdot p^{\frac{k+1}{2}}$: T. value of |H|, has to be computed

Use coeff $\frac{\partial^2 \sigma^{-n}}{\partial x^2} = -n \sigma^{-n-1} \frac{\partial \sigma}{\partial x}$
 $\frac{\partial^2 \sigma^{-n}}{\partial x^2} = -n \sigma^{-n-1} \frac{\partial^2 \sigma}{\partial x^2} + n(n+1) \left(\frac{\partial \sigma}{\partial x} \right)^2$
 $= \frac{-n}{\sigma^{n+1}} \frac{\partial^2 \sigma}{\partial x^2}$ say $\frac{\partial^2 \sigma}{\partial x^2} \approx 2\sigma_{\text{observed}}^2 \approx 2\sigma_0^2$

so $|H| = \frac{p^k}{\sigma_0^{-n}}$

$\left(\frac{-n}{\sigma_0^{n+1}} \cdot 2\sigma_0^2 \right)^k = \left(\frac{-2n}{\sigma_0^{n+1}} \right)^k$

$\sqrt{\frac{p^k}{|H|}} \cdot p = \frac{\sigma_0^{-n}}{(-2)^k \cdot n^k} \cdot \frac{1}{\sigma_0^n} = \frac{\sigma_0^{-(n+k)}}{(-2)^k \cdot n^k}$

I think $(-2)^k$ is irrelevant. Also mult by 2^k for actual full width.

so $\frac{\sigma_0^{-(n+k)}}{n^k} \ln \rightarrow -(n+k) \cdot \ln \sigma_0^2 + k \ln n$

$p = \sigma_0^{-n}$

$\frac{p^k}{|H|} = \frac{\sigma_0^{-nk} \cdot \sigma_0^{(n-1)k}}{(-2n)^k}$
 $= \frac{\sigma_0^{-nk} \cdot \sigma_0^{nk-nk}}{(-2)^k \cdot n^k}$
 $= \frac{\sigma_0^{-k}}{(-2)^k \cdot n^k}$

2.19.05
57M

$$\pi \cdot e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \max$$

The gamma distribu. is a summary of data.

if $\sum x_i = n\mu$ $\frac{1}{\sigma} e^{-\frac{x}{2\sigma^2}} = \max$

$\sigma e^{-\frac{x}{2\sigma^2}} = \min$

$$e^{-\frac{x}{2\sigma^2}} + \sigma \cdot \left(-\frac{x}{\sigma^2}\right) e^{-\frac{x}{2\sigma^2}} = 0$$

$\frac{x}{\sigma^2} = 1$ so $\sigma = \sqrt{x}$.

$$\prod_{i=1}^n \left(\frac{A^\alpha}{\Gamma(\alpha)} x_i^\alpha e^{-\frac{x_i}{\beta}} \right) = \frac{(\prod x_i)^\alpha}{A^\alpha} e^{-\frac{\sum x_i}{\beta}} = \max$$

$A^\alpha e^{-\frac{B}{\beta}} = \max$

We need normal constants, i.e.

$$\int_0^\infty x^\alpha e^{-\frac{x}{\beta}} dx = \frac{\Gamma(\alpha+1) \beta^{\alpha+1}}{\Gamma(\alpha+1)}$$

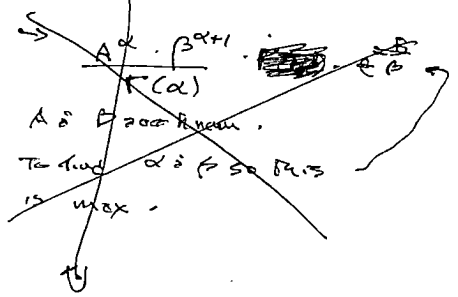
$\frac{x}{\beta} = y \Rightarrow \frac{dx}{\beta} = dy$
 $x = \beta y \Rightarrow dx = \beta dy$

$$\int \beta^\alpha y^\alpha e^{-y} dy \cdot \beta$$

$$= \beta^{\alpha+1} \int y^\alpha e^{-y} dy = \beta^{\alpha+1} \Gamma(\alpha+1)$$

No! $\alpha!$

$$\int_0^\infty x^\alpha e^{-\frac{x}{\beta}} dx = \Gamma(\alpha+1) \beta^{\alpha+1}$$



$\Gamma(x+1) = x!$

$\int_0^\infty x^\alpha e^{-x} dx = \Gamma(\alpha+1)$

$\int_0^\infty e^{-x} dx = 1 = \Gamma(1)!$

Normal for each is $\frac{\Gamma(\alpha+1) \beta^{\alpha+1}}{\alpha!}$

So we want max $(\Gamma(\alpha) \beta^{\alpha+1})^{-1} \cdot A^\alpha \cdot e^{-\frac{B}{\beta}}$

A & B over known. we want α & β so as to max this

T. values of α & β are functions of A & B i.e. $\frac{A}{B}$ i.e. $\frac{A}{B}$ i.e. $\frac{A}{B}$ i.e. $\frac{A}{B}$

Tables. Maybe $\beta \propto B = \sum x_i$?

Get rid of π : $n \cdot \text{root } \Gamma(\alpha) \beta^{\alpha+1} \left(\frac{A}{\beta}\right)^\alpha e^{-\frac{B}{\beta}}$

so $A \rightarrow \frac{A}{\beta}$; $B \rightarrow \frac{B}{\beta}$ would fix it so we only have to max wrt A & B .

replace A & B by $\frac{A}{\beta}$ & $\frac{B}{\beta}$ i.e. $\frac{A}{B}$ means.

so max $(\Gamma(\alpha) \beta^{\alpha+1})^{-1} A^\alpha e^{-\frac{B}{\beta}}$

max $\frac{A^\alpha e^{-\frac{B}{\beta}}}{\Gamma(\alpha) \beta^{\alpha+1}} = \frac{A^\alpha}{\Gamma(\alpha) \beta^{\alpha+1} e^{\frac{B}{\beta}}}$

$-\ln \Gamma(\alpha) - (\alpha+1) \ln \beta + \alpha \ln A - \frac{B}{\beta} = 0$

$\frac{d}{d\alpha} = \left(\frac{d}{d\alpha} \ln \Gamma(\alpha) \right) - \ln \beta + \ln A = 0$

$-\frac{(\alpha+1)}{\beta} + \frac{B}{\beta^2} = 0 \Rightarrow \alpha+1 = \frac{B}{\beta}$ so $\beta = \frac{B}{\alpha+1}$ - correct!

or $\beta = \beta(\alpha+1)$ or $\alpha+1 = \frac{B}{\beta}$ $\beta = \frac{B}{\alpha+1}$

so max $(\Gamma(\alpha) \left(\frac{B}{\alpha+1}\right)^{\alpha+1})^{-1} A^\alpha e^{-\frac{B}{\beta}} = \Gamma(\alpha)^{-1} \cdot \frac{(\alpha+1)^{\alpha+1}}{B^{\alpha+1}} \cdot A^\alpha e^{-\frac{B}{\beta}}$

max $\frac{A^{\alpha+1}}{B^{\alpha+1}} \cdot \left(\frac{\alpha+1}{B}\right)^{\alpha+1} \cdot \frac{1}{\Gamma(\alpha)^{\alpha+1}}$

max $\approx \frac{(\alpha+1)^\alpha \cdot A^\alpha}{B^{\alpha+1}} = \frac{1}{\Gamma(\alpha+1)} \cdot \frac{(\alpha+1)! \cdot A^\alpha}{\Gamma(\alpha) \sqrt{(\alpha+1) \pi}}$

$C = \left(\frac{A}{\beta}\right)^{\alpha+1} \cdot \frac{(\alpha+1)^{\alpha+1}}{\Gamma(\alpha)^{\alpha+1}}$ so given C find α

$\frac{(\alpha+1)^{\alpha+1}}{e^{\alpha+1}} \approx \sqrt{\alpha+1} \cdot 2\pi \approx (\alpha+1)! = \Gamma(\alpha+2)$

$\left(\frac{C(\alpha+1)}{\Gamma(\alpha)}\right)^{\alpha+1} \approx \max$

2.19.05
 α. STAY

25 α.12

Max_α $\sqrt{\alpha+1} \cdot A^\alpha \rightarrow \frac{1}{2} \ln(\alpha+1) + \alpha \ln A$

$\frac{1}{2} \frac{1}{\alpha+1} + \frac{1}{\alpha} + \ln A = 0$ $\frac{1}{\alpha+1} + \frac{2}{\alpha} + 2 \ln A = 0$

$\frac{\alpha + 2\alpha + 2}{(\alpha+1)\alpha} = -2 \ln A$

$\frac{3\alpha + 2}{(\alpha+1)\alpha} = -2 \ln A$ $\frac{3\alpha' - 1}{\alpha'(\alpha'-1)} = -2 \ln A$

no! $\int_0^\infty x^n e^{-x} dx = n!$ not $\Gamma(x)$

$\alpha' = \alpha + 1$ $\frac{1}{\alpha'} + \frac{2}{\alpha'-1} = -2 \ln A$

$\alpha'' = \alpha + \frac{1}{2}$

$\frac{1}{\alpha'' + \frac{1}{2}} + \frac{2}{\alpha'' - \frac{1}{2}} = -2 \ln A$

$\frac{\alpha'' - \frac{1}{2} + \alpha'' + 1}{\alpha''^2 - \frac{1}{4}} = \frac{2\alpha'' + \frac{1}{2}}{\alpha''^2 - \frac{1}{4}}$

12:10 $2 \frac{1}{4}$
 2:26 2
 4:45 $4 \frac{1}{2}$
 9:12

So it's $\sqrt{\alpha+1} \cdot A^\alpha \rightarrow \frac{1}{2} \ln(\alpha+1) + \alpha \ln A$

$\frac{1}{2} \frac{1}{\alpha+1} + \ln A = 0$ $\frac{1}{\alpha+1} = -2 \ln A$ $\alpha+1 = \frac{-1}{2 \ln A}$

A = Gamma function. so $\frac{1}{\alpha+1} < 0$!

I must have gotten Recy mixed up somewhere!

$\alpha+1 < 0$; $\alpha < -1$
 unlikely!

Max $(\alpha! \beta^{\alpha+1})^{-1} A^\alpha e^{-\frac{A}{\beta}}$: $\beta = \frac{A}{\alpha+1}$: $\frac{A}{\beta} = \alpha+1$

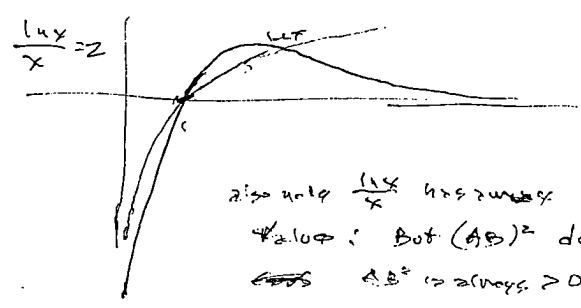
Max $\frac{1}{\alpha!} \frac{A^{\alpha+1}}{\alpha^{\alpha+1}} e^{-\alpha+1} A^\alpha = \frac{(\alpha+1)!}{\alpha^{\alpha+1}} \cdot \frac{A^{\alpha+1}}{\alpha^{\alpha+1}} \cdot \frac{1}{\alpha+1} A^\alpha$

$\frac{1}{\alpha} \sqrt{\alpha+1} \cdot \frac{A^{\alpha+1}}{\alpha^{\alpha+1}} \cdot \frac{1}{\alpha+1} \cdot \frac{1}{\alpha+1} A^\alpha$

$\frac{1}{\alpha} \sqrt{\frac{2\pi}{\alpha+1}} \cdot \frac{A^{\alpha+1}}{\alpha^{\alpha+1}} = \max (\alpha+1)^{\frac{1}{2}} \left(\frac{A}{\alpha}\right)^{\alpha+1}$

$(\alpha+1)^2 = x$ $\frac{1}{x} \cdot (AB)^{\frac{x}{2}} = \max$ $-\frac{1}{2} \ln\left(\frac{x+1}{x}\right) + \ln(AB) \left(\frac{x}{2}\right)$

See if I made a mistake in α.11.22 pp other than $\Gamma(\alpha)$ should be $\alpha!$



for $z > 0$,
 there are two solutions.
 One probably for max, another for min!

$\ln x = 2x \ln(AB) = x \ln(AB)^2$
 $x = (AB)^{2x}$
 $x = (AB)^x$

$\frac{\ln x}{x} = z$

$x = e^{xz}$

also note $\frac{\ln x}{x}$ has a maximum value.
 value: But $(AB)^2$ does not have a max value > 0 .
 AB^2 is always > 0 : could be made as large as you want

$$\int_0^{\infty} x^{\alpha} e^{-\frac{x}{\beta}} dx = \alpha! \cdot \beta^{\alpha+1} ?$$

$$\frac{x}{\beta} = y \quad \frac{dx}{\beta} = dy \quad dx = \beta dy$$

$$x = \beta y$$

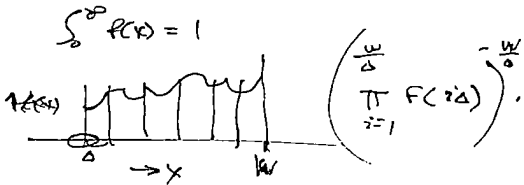
$$\int \beta^{\alpha} y^{\alpha} e^{-y} \cdot \beta dy = \beta^{\alpha+1} \int_0^{\infty} y^{\alpha} e^{-y} dy = \beta^{\alpha+1} \alpha!$$

so d.f. = $\frac{1}{\beta^{\alpha+1} \alpha!} x^{\alpha} e^{-\frac{x}{\beta}}$; first moment =

$$\frac{1}{\beta^{\alpha+1} \alpha!} \int x^{\alpha+1} e^{-\frac{x}{\beta}} dx = \frac{\beta^{\alpha+2} (\alpha+1)!}{\beta^{\alpha+1} \alpha!} = (\alpha+1)\beta \quad \text{mean} =$$

Geom mean of d.f. $f(x)$

β = actual mean.
 A = Geom mean.



$$= \pi(f(a))^{1/w} \quad (\alpha+1)\beta = \beta$$

$$= \exp \left(\frac{1}{w} \sum_{i=1}^w \ln f(i \cdot a) \right)$$

$$= \exp \int_0^w \ln f(x) dx$$

$\exp \left(\int_0^{\infty} \alpha \ln x \cdot \frac{x^{\alpha}}{\beta} dx \right)$; To geom mean of any $x^{\alpha} e^{-\frac{x}{\beta}}$ & would be zero!

Max $\frac{A^{\alpha}}{\beta^{\alpha+1}} \frac{\alpha!}{(\alpha+1)!} \rightarrow \left(\frac{A}{\beta}\right)^{\alpha+1} \frac{1}{\alpha+1}$

Max $= \left(\frac{A}{\beta}\right)^{\alpha+1} \cdot \frac{1}{\alpha+1}$

$$= \left(\left(\frac{A}{\beta}\right)^2\right)^{\alpha+1} \cdot (\alpha+1) = \text{max}$$

$$x = \alpha + 1 \quad \left[\left(2 \ln \frac{A}{\beta}\right) \cdot x \right] \frac{d}{dx} = \text{max}$$

$$2 \ln \frac{A}{\beta} + \frac{1}{x} = 0$$

$1, \frac{1}{2} \rightarrow .7$ param \leq avr.

A = Geom, β = arith mean.

$$\frac{A}{\beta} < 1 \quad \alpha + 1 = \frac{1}{2 \ln \frac{B}{A}}$$

$$\alpha + 1 = \frac{1}{2 \ln \frac{B}{A}}$$

$$\alpha = \frac{1}{2 \ln \frac{B}{A}} - 1$$

B = Arith mean
 A = Geom mean.

$$\alpha + 1 = \frac{1}{-2 \ln \frac{A}{B}}$$

$$\alpha + 1 = \frac{1}{2 \ln \frac{B}{A}}$$

This assumes Sterling Approx is exact.

So the value that we end up with is $\alpha + 1 = \frac{1}{2 \ln \frac{B}{A}}$

Substitute for goodness of fit,

In a PC sense it's $\frac{1}{2 \ln \frac{B}{A}}$ maybe (divided) by A is to take RR power.

So the f. pc of f. Gamma D.f. (for best fit of data) is not used to raise.

β is used for f. Data squares all $>$ ($<$) a certain val,

For RR's calcn. we need Arith & Geom means of data.

β Gauss d.f. " " Arith mean is mean of squares of data.

2.23.05
2.21.05
STM

25 α/4

This is the end of the 25 α/1 to 25 α/4 secy.
STM 26 comes next.

on fitting data to $P(1-p)^{n_1} (1-p)^{n_2} / (n_1+n_2)! p^{n_1} (1-p)^{n_2}$ as $p = \frac{n_0}{n_0+n_1}$

T. set of pfs on $p \in [0,1]$ for n data pfs, ~~probability~~

say X_2 data: $\bar{x} = \frac{1}{n} \sum x_i$ probly $n_0 = \bar{x} \cdot \alpha$
 $n_1 = (1-\bar{x}) \cdot \alpha$
 α should be $\approx n$ usually

No!

so, what is α ?

we want max $\left(\frac{\bar{x}^{\alpha \bar{x}} (1-\bar{x})^{\alpha(1-\bar{x})}}{\alpha \bar{x}! (1-\bar{x})!} \right)^\alpha \cdot (n+1)! \leftarrow \text{constant}$

$\frac{e^{\alpha \bar{x}} e^{\alpha(1-\bar{x})}}{\alpha \bar{x}! (1-\bar{x})!} = \frac{1}{\sqrt{2\pi \alpha \bar{x}} \sqrt{2\pi \alpha (1-\bar{x})}}$

$\frac{e^\alpha}{\alpha \bar{x}! (1-\bar{x})!} = \frac{1}{\sqrt{2\pi \alpha \bar{x}} \sqrt{2\pi \alpha (1-\bar{x})}}$

$\frac{e^\alpha}{\alpha \bar{x}! (1-\bar{x})!} = \frac{1}{\alpha \bar{x}! (1-\bar{x})!} \cdot \frac{1}{\sqrt{2\pi \alpha \bar{x}} \sqrt{2\pi \alpha (1-\bar{x})}}$

$= \frac{e^\alpha}{\alpha \bar{x}! (1-\bar{x})!} = \frac{e^\alpha}{\alpha \bar{x}! (1-\bar{x})!} \cdot \frac{1}{\sqrt{2\pi \alpha \bar{x}} \sqrt{2\pi \alpha (1-\bar{x})}}$

No!

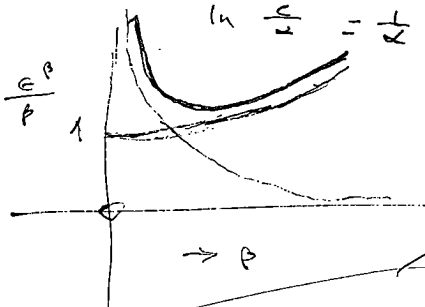
nb! I want $\frac{\prod x_i^{\bar{x}} (1-x_i)^{1-\bar{x}}}{\bar{x}! (1-\bar{x})!} = \max$

$= \left(\prod x_i^{\bar{x}} \cdot \prod (1-x_i)^{1-\bar{x}} \right) = c^\alpha$

$\frac{e^\alpha}{\alpha \bar{x}! (1-\bar{x})!} = \frac{c^\alpha}{\alpha \bar{x}! (1-\bar{x})!}$

take ln -21 $\alpha \ln c - \alpha \ln \alpha - \alpha \ln \bar{x}! - \alpha \ln (1-\bar{x})! = 0$

$\frac{\partial}{\partial \alpha} = \ln c - \ln \alpha - \frac{1}{\alpha} = 0$



$\ln c = \frac{1}{\alpha} + \ln \alpha$ $c = \alpha e^{\frac{1}{\alpha}}$ $\beta = \frac{1}{\alpha}$; $c = \frac{1}{\beta} e^\beta$

$\frac{1}{\beta} e^\beta$ has a min at $\beta = 1$

$-\ln \beta + \beta = 0$ $-\frac{1}{\beta} + 1 = 0$ $\beta = 1$

So c must be ≥ 1 .

$c = \left(\prod x_i \right)^{\bar{x}} \left(\prod (1-x_i) \right)^{1-\bar{x}}$

Also since $\frac{d}{d\alpha} = 0$ at $\alpha = 1$, we assume $c = 1$ is max, not other Min.

so c is \bar{x} completely determined by \bar{x} . or \bar{x}_i from mean \bar{x}

\bar{x} is geom mean of $(1-x_i)$ $\therefore \geq \text{geom mean}$

Re: Lectures 3: What are the major apparent obstacles to final TM?

In "Phase 1"

1) In d-induction phase: Use of Lisp-like language:

a) Method of implementing recursive is unclear ^{INTSO}

b) Method of efficiently searching over corpus is unclear: I'd like a quick way to test

new cases: putting problems in order of empirical dirtiness is not bad. - its a "quick and dirty" technique.

c) related to (b) The "recognizer function" idea in the IDSR report, seems good,

but needs much more work! Go back over it carefully! -> see 28.22 for note.

2) In s-induction phase:

a) The 3 input μ mc (3IU) is universal, I don't yet see a practical way to implement it. (Re: "2 input" part is "under control")

b) I have 2 way to get 3 kinds of probabilistic output of rules: 1. Beta 2. tools

say from 0 to 1 use Beta distrib. 2. from 0 to ∞ use Gamma P.D. 3. from ∞ to ∞ use Gaussian D.F.

These are all monometal d.f.'s. - But these methods are not universal.

It may be that by using a uniform D.F. from 0 to 1 as ~~one~~ ^{one} ~~input~~ ^{input} problem

down ~~the~~ ^{the} ~~input~~ ^{input} to a universal (log (Lisp, say) I can get any output

D.F. ~~this idea may~~ be able to solve/clarify problems of ~~spread~~ ^{spread} in

say, "How Many ~~bits~~ ^{bits} ~~are~~ ^{are} linear & non-linear")

c) I don't have a good way to get probabilistic D.F.'s on ~~strings~~ ^{strings} outputs.

Some ways: 1) Output is a single string. We use some kind of distance function to other strings, for e.g. D.F. on other (non-control) strings. Initial output string

can be center of a cluster but not rest of function (?). The distance function can be: Edit dist; Hamming dist; An ~~other~~ ^{other} ~~into~~ ^{into} distance" using instruction codes.

This last can be approximate (not use full instr. set is a "Edit dist") or fuller set, &

"heuristic using" (Approx-Sim. GPS) to create desired corpus strings (expensive). "T. edit distance" may be essentially identical to the "3 input μ mc"

3) In ~~the~~ ^{the} Induction of Phases:

a) I have 3 approaches to using rays of just:

1) PPM; May work a bit; but would be much better ~~if~~ ^{if} ~~sub-trees~~ ^{sub-trees} could be recognized as "contexts": I have a way that would pretty work, but it seems very slow.

2) S-CFG discovery! I don't know how good this is - to what extent does

it recognize sub-trees as contexts or as "chunks". Also, I don't know whether

it's really very well, or I suspect it's slow. I could use ~~some~~ ^{some} S-context Sensitive Grammar w. context being in function of sub-trees.

5 TM

20

A big poss. trouble: Both the need of that ratio may not be monotonic in x_0 (if threshold or output).
 a/o be a bump, irrat. function of " x_0 " (outputs).

③ "G.P." for "range of part" — use of S-city for summarizing corpus & coding good. (see 26.25 for some diff's w. PLS):
 Use of $\frac{\sigma_g}{G_{max} - H_g} \Rightarrow$ criterion for threshold selection — seems good, ~~but~~ at first glance, but is very dependent on Gorc odd. Here Gorc is PC, so its unambiguous.
 T. Q is: will the formula work with that Gorc? — It may, because ~~the~~ this Gorc may be "linearized". \rightarrow 28.10 discusses Phase I D. Induction TSO's.

Phase II

1) The induction for phase II involves split to real distribns. They are usually monomodal, so I can use either Gauss or Gamma D.F. for output: \Rightarrow means f. output need be but easy to read, ^{usually a combination of D.F.} I can then find params of f. D.F. to give best fit.
 [Output is p.d of G for application of PSM; to problem, ~~for~~ a certain fixed time.]
 IDSA has an algorithm of this. T. system assigns a PC to a given PSM, depending on how good it was based in productivity & resources used. (Ideally to assign loads of (PSM_i, Prob of t) ^{in allowable} gives P.D. as output.

I haven't much look into practical aspects of realizing it.

2) I want to have a grammar to generate PSM's. Seems like a very abstract problem in itself. It could be a CFG — am not sure it need be a ~~CFG~~ Loopless Grammar may be adequate. Intact, ~~a~~ Loopless Grammar may be adequate. ??

3) It may be poss. to do Phase II first w/o doing phase I ever! I haven't look into this much ... I have no idea what it diff's are. I originally had the idea that a fairly good induction system would be needed. How expensive in many kinds of problems — E.g. Reason in Phase I. Say a well developed "Phase I" induction system ... but this may mean many kinds of problems — How expensive in many kinds of problems — E.g. Reason in Phase I.
 not generally — But 20 does have to work eventually. If I did in now (first) 2.23 found out to be feasible, I'd save a lot of time!
 4) I want Phase II to be able to do arbitrarily good self-improvement, which is the main idea of Phase II: whether the system of 23 could do this, is unclear; well, it is able to improve search form for estimates of Expected value of G for (PSM_i, Prob of t, T)

Logical Reasoning: Much used in "Heuristic PSM", Math, etc: I haven't

gone in to it at all, but ... it's probly an imp. part of "intelligence". T. plan was to have PM logic as a mathematical study, then learn to correct it w. heur. substitution that would arrive in a few kinds of problems.
 It seems likely that Logical reasoning would be imp. in solving Phase II's problem of allocating PSMs to probs. Also creating new PSMs.

33.20

26.00 - 27.40 \Rightarrow a kind of outline of "state of TM": I would do well to "flashout" the outline by biblirefs. In particular, it would ~~be~~ make a good way to organize my notes! Make special indexing for (notes not in the outline).

SN NB Look at Abstract of IDSI A report. This is an earlier version of the present "Phase I, Phase II" system. Using the tricks of PIM; PSG Drey; Enhanced GP: would it also be workable?

Re: Phase I duffies: Writing is learning of TSO!

Poss. initial TSOs: ^{from} ① Elementary Algebra Book: Notation: Function Eval; (Soln. of linear, (then N-linear eqns), of greater complexity, Differentiation, ^{open diffy} integration.)
② Learning of Notation of Maple, Mathematics - etc. \rightarrow Soln. elements, integration. ^{difficult, inefficient, numerical!}

Will PIM' (augmented?), PSE (augmented?), G.P. (augmented) be adequate?
Will ~~the~~ contemplated TSO's be adequate - are there many more "intermediate concepts" needed in learning. How we substitute?

It may, indeed, be that we cannot get far even in Phase I to get to 2 used phase II! ^{No to that Phase II needs ϵ -induction w. numerical output. 26.12-20 discusses how to implement this, but suggests no TSO's: SM & HR would seem to be reasonable TSO's.}

~~Not mentioned~~ ^{26.07 does mention} (in 26.00-27.40) \Rightarrow the "Recognition function" (= "Cases ~~of~~ function") of

Early part of IDSI A report. These Recognition problem types $\hat{=}$ allocate an appropriate set of \hat{O} functions to them. This markedly reduces amt. of time needed to check each \hat{O} modification past cases: Only cases of appropriate R category need be tested. This Accords w. much "Human Scrutinizing".

SN ON "A-LIFE": A-life is interesting & diverse because w/ ^{almost} immediate get into a great diversity of problem & solution types, but about eq. to diversity of life forms itself. It is competitive and like competition \Rightarrow strong & diverse, & usually of some level of bit by as to competing forms. The life forms solving problems over themselves available (via parasites) as parts of solving to new problems.

5TM

Karl Newett : } T. Sci Community as
+ 1 co-author } an Intelligent System }
MIT Report

! Source revisited keep 100 best models; but for search, use
Evolver or PPM to explore search, so one doesn't lose
diversity as TSO continues. Try to keep diversity at minimum
level: at least default. a tall time.

Could .05 - .07 be "trick soln" to TM? See .05 ff for objectives suggesting to overcome objectives.

There may be several different ways to "↑ Diversity": If one doesn't ↑ diversity,
One saves best 100 solns. to induction problem. As TSO continues, fewer
fewer of the models fit the data "acceptably" — until none do.

One way to ↑ diversity (at least fast and to compensate by loss of diversity
viz .05-.07) is to make stock Models of good cards much for (viz S-CFG
or PPM). This gives a large soft (probable) language of cards over

can be very "diverse" — but not really "clever". Could I formulate the creation
of a set of "clever" diverse cards as an OZ problem?

The amount of CC spent on S-CFG discovery w/o Augmenting PPM
could be regulated by the amount of diversity needed.

NB Increasing amt. of time spent "updating" PPM's S-PSG discovery, will ↑ THIS MUST BE ANALYSED
accuracy of predictions — is this necessarily tied in w. "Diversity ↑"? → .30

Another trick: I was a bit uncertain about how to model a Game (≡ "Githess dnd")
in ways that would ① ↑ speed of comp. a Game, or test approach ② Enable crossover
of the Game to find Cards of by G. In both of these applications, it would be well if TM
could "watch the Game being computed" — or be given to Alg. that computes the Game, —
which are a "No-No" in Genetic/evolutionary prob. Solns, because its much not like Gp. Evolu.

It would be well to attempt modeling a Game in "reversible = invertible" funcs, so that we
could more easily implement search for cards of by Game.

30: (20) One method ↑ diversity was to store best 100 or best 1000 cards at all times. In S-PSG discovery
the cards not branch away from each other, early in attempts Hill climbing used in
will tend to be diverse. This is also true of cards obtained by branching early in TSO
Hvr. if one discovers one hasn't enuf diversity, it's probly not practical to backtrack back to
to beginning! One could backtrack back a few critical parts, hvr., but one could
retain more early branches in S-PSG discovery (say, retain all early branches until the
no. of stored nodes was 100 or 1000 or whatever. One could store 10000 early branches
then slowly winnow them down to only 1000 or 100 as we work toward an accepted
set of cards.

This is
domain
to Sci
community

S. TM

10:28.40 : ALife! - cont: One way to get it to solve more practical problems: Make Environment have features that one can learn to take advantage of. E.g. T. rules of reproduction are always, slowly changing. Any other rules about getting food, consuming other organisms, etc, are subject to slow change (perhaps "randomly" ... unclear as to how random, is "random in what way")

Part of this is ~~image~~ analogous to peoples in varying climate developing more rapidly than those in constant (benign) climate.

While it's easy to see a line: still, how do we get the population to solve what are to the designer, "useful problems"?

Perhaps by experiments, the designer could try various "rules" (innovations / changes), to see what the community is able to do w. effectively.

Might start by giving equations to solve. For a soln, individual gets more food/reproduction ~~image~~ rights or whatever. We would want

individuals to be able to help one another: possibly complete as well. This "equality" reward - would want to dit. w. other "normal" kinds of competition within: ALife "Community"; Perhaps have a "patent system" so a person gets credit every time

Someone uses his ideas (s). - But have reasonable limit on duration of patent ... so it soon becomes "Public Domain".

~~image~~ could .13-.16 be adapted by itself? T. fraction of "reproductive rewards" obtainable by "solving problems" is an input param of t. system.

I have to work out mechanics of presentation of problems: realization by population that they get rewarded for soln, etc. Mechanics of TSQ's, etc.

0
13
16
20
30

2.28.05
J. TM

Criticism of "A proposal for Strong A.I." 1961
adequacy for Strong A.I. of man

In Phase 1, "An adequate TSO" is a computer "regy such as (em)" means only that the subfunctions needed for a soln. to the new problem have been detected (or found). So the soln. has a not too small app. This can, for many large problems, be still very small.

A cleverer machine will somehow be able to consider ~~the~~ exact nature of the problem & assign very large ~~pc's~~ pc's to certain functions; also by "logical analysis", come upon a soln. w. pc = 1 (or very high). It is .00-.03 that characterizes a "successful" phase 1. Its about as far as phase 1 can get. It assumes all significant "sub-procs" (= subfuncts) can be found by the "regy such" [Regy such is art. corpus of past successful finds.]

There are, then, 2 ~~types~~ types in phase 1: (1) Regy such over past successful finds — to find usable sub-functs. (2) Search over combinations of detected functs to find solns to problems. ||| T. Q was: In which (combination) of these was Phase II to be applied? (2) can use contexts of various kinds. (1) could be augmented by "context" — contexts of ~~4~~ subfuncts that have been found.

Both 1.9 & 2.0 would speed up TM considerably. Since "context" in its most general form, involves the entire previous corpus & ~~is~~ suby regys in that corpus — that "context" can be quite powerful. (The of course, contexts of progressively greater power/complexity become progressively more difficult to find.)



If it would seem that the 2 types of (.16-.19) would be rather well defined, it is not hard to do in a near optimum way — Regy are rather simple "rote" suches. (The it may be poss. to speed up detection of sub-functs). S. C. G. discy 2/0.

Augmented PPM might be adequate. I guess I was thinking that phase II could detect more regy types ~~using~~ using general contexts (.19-.20; ~~21-24~~ ^{general context}).

Any context detected by "regy such" should be translatable into a way of searching (= such) more efficiently — so the 2 types of (.16-.19) must of necessity be detectable. Essentially the "regy such" is essentially a stochastic grammar.

Context Sensitive Grammars (special) : 32

00:31.40 : to generate subsequent Cands. A "perfect" Phase I would use full ALP to extrapolate past set of Cands (i.e. Inductive) to propose future Cands.

The PSM would be "unsorted", Phase I still doesn't understand what "optimization" means! — So clearly Phase II has much to offer!

A "perfect" Phase I would use ALP to directly extrapolate past PSM's (Ques, Ans).

I guess I'm a bit vague on Ph II's all over operation! I did want to construct a "complete" grammar over all the ~~known methods~~ PSM's

know of. — But the method of ~~deciding~~ deciding which PSM was appropriate for each problem, is a (perhaps closely linked) problem of generating a new PSM appropriate to a new problem Phase I haven't much worked on.

One kind of "perfection" for Phase II would be a $\forall j, T$ PSM's for $(PSM_j, Prob_j, T)$. E.g. even so, we have normally, not such a good set of PSM's is knowing the problem of finding PSM's that have good "fig. of merit" for T . So we need ~~some~~ ways to generate good PSM's for a given prob j , a way to find T v.g. PSM's for that prob j .

3.1.04



T. problem of finding common sub-trees int. "Successful cand" corpus, would seem to be identical (or a most identical) to the problem of discovering S-CFG's; except in S-CFG discy, we also look for "Merges", which can give loops. So maybe S-CFG d. is uniformly better than simply looking for common sub-trees.

Common sub tree prob is also done as a (serious) enhancement of PPM.

Is there any way to merge CFG & PPM to get a system uniformly better than each?

I.e. ~~perhaps~~ CFG doesn't seem to use context (enough or at all) in choosing

substitutions. Context sensitive Grammars would help much; But I have to make sure the proper kinds of contexts (i.e. sub-trees) are used,

What needs to be done for Phase I!

- 1) Write PPMs for PPM, s-cfg disc.
- 2) write TSO's { First for D. Induction
Next for S. Induction
- 3) Augment PPM w. sub-tree Contexts:
Augment s-cfg to s-cfg dirty using sub-trees of Contexts.

Both input-kinds of Augmentations
 PPM aug. has been specified above:
 T. S. c.f.g. → S. c.f.g. has to be sub-tree
 much more

4) Try improving GP & letting it to work ~~with~~ w. interproblem lang,

Recall TSO's, then various improvements of GP →
 As is Phase I can only have numeric output for s-ind (future has string output for d-induction)
 5) Work on Phase II. $\left. \begin{matrix} \text{qual. (in time) may make string output possible.} \end{matrix} \right\}$

I'm having second notion $\frac{G}{G}$ as very useful! It needs better mass understanding! It $\frac{G}{G}$ may be ought to work in $\frac{G}{G}$ (culture) ought to be linear.

6) T. "Core Cancer" problem is Very Imp. Any long-term

research that TM does — internally also in R.W. involves solving PPM problem.
 It may be easy to devise a TSO of "Core Cancer" — type problems of ↑ diffy: E.g. slowly ↑ horizon (horizon = not error double better finish). Or no of poss. expts; diversity of results can be "Adjusted" for t. TSO.

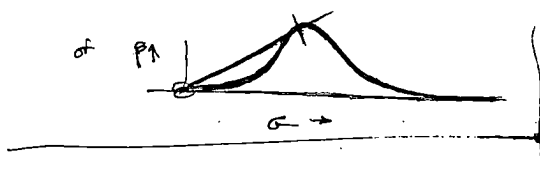
My Soln. of "Core Cancer" maybe a bit complex; I may be able to significantly simplify it by expressing it in Recursive form. Lots of work on "Approx Dynamic Ppm" started by Look up "Approx Dyn" (Ppm on Google).

7) Testing each OI cond. on t. entire previous Corpus is very inefficient. Several approaches:

- a) The Recogn. functions of each part of DPSM report
- b) Putting past problems in order of diffy. ...
- c) Problem Pool approach.

27.32. Phase II (cont): 5) It should be poss. for TM to look at a problem & generate

(or somehow get a pd on all poss. PSMs) one or more PSMs that would have by FOM (Fig. of Merit) for that problem j. i.e. $h_j \theta \rightarrow G$?



Actually, it may be O.K. to select PSM w. max h_j for $(PSM_j, Prob_j, T)$. — T. details of t. D.F. are usually not of interest!

Avr. to get t. h_j 's I may need to get t. entire G v.s. PC curve for each (j, T) .

This is so I can best fit t. curves using negative as well as positive data.
 It may be that I can get a simple approx formula for h_j , given t. set of G values for past success & set of G values for past failures.

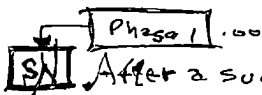
In case "T" is not given; I may want to work first on PSM Max θ !

Phase II: .10

3 Input UMC

Brasano? ??

2.5	33
3-1	14
	11
	258
	25 d.



10: After a successful cond search; - yielding 1 or more acceptable cards - Th. following "raggy search" examines these new cards - expecting to find novel (concs/ sub-functions), i.e. How do the new cards "differ" from the old "non-successful" cards? on correct problem "cards"?

- N. Zealand
- Aust: Max 2 places
- RHUL London
- Rupano?
- MIT
- Carif. - some how
- China
- Japan
- Singapore (Aot!)
- India
- S. Africa
- IBM Almeda?

ON PHASE II: Does Ph II really understand what Opten means?

10: How much Better is it than Phase I in its "understanding"? A good Optimizer will look at a problem, consider different ways to solve it, decide on one on the basis of max expected G in t. available time. This is what Phase II does. I guess what I felt was "mising" was: we give TM "G = 1 - x^2" to find x > G is as large as poss. (in available time). In the present case, a person would usually do $\frac{dG}{dx} = 0$ at $x=0$ as a cond. Then try it thru by $x=0.01$ to a few max. TM could have a general second order n-linear opten routine built in.

10: Hyp. I'd (also) like TM to be able to discover routines like that on its own

On 3 Input UMC

21: We can simply use 2 general Lisp function w. 3 inputs for pooling, "easy" problems; searching over to "R" input shouldn't take much time, since solus will be of type. A Barn seq could be possible by using (Hvr, do cond) + solut. Barn seq. obtained by interdy overall prod. factors: T. saved data distribution. (C-12) P. 11

On Reparsing in QATM

22: Very imp't thing to do! T. set of reparses ist. Cost. product of reparses for each QA problem soln. (well maybe not such a "prohibit" - T. "individual solns" approach for a subset of corpus, & the subsets "Nests"..... A different kind of "soln." for QA problem has many solns. to various subsets of the QA corpus

On Phase II

20: After TM has collected a lot of (parts of PSMs) & partially assembled PSMs, this is in a good position & data on FOM of various PSMs for various fields. W. this data, TM should be able to construct PSM's that are likely to have the FOM for a particular problem. This Expresses Ph II's goal as a problem, not as a Ph I problem!

34: Could Ph II's main problem be worked on by P&I's "raggy. Snek" algom(s)? Well, I think G.P. can work on it. G.P. being a type of BU snek. This isn't a real OE snek, hvr. It's a "Better than" snek. I could use a searchover a finite no. of PSM's to find one w. best FOM for a given problem. This is a step toward 31-34; It could be used for OE problem of 34 - so this would be a "round about" soln to Ph II - No a very "Round About" soln!

0 (34.24-23) ^{Spec}: 31UMCS (cont.): 34.21-23 seems like it might work! (so Universality for String Outputs).
For π and numeric output: T. ~~idea~~ output could be for mean & var. & t. R input would tell

how many s & 's away & desired "A" is. \leftarrow (Re Re is not clarifying mind)
 23. 34.24 I still don't see how to unify its treatment of t. Bern sep. w. t. Beta function formula Approach.

- 2-I 4, 6, 7
- 2-II 1, 9, 5
- 3-I 10
- 3-II 6
- 4-II 6

0 : 34.30 : T. Reparsing idea is related to t. tradeoff between making definitions using ~~the~~
 (loop constraints) (like PPM) for Proda. I probably want to do Diff ... but
 this somehow seems "inefficient"!

- 1) APL to 2014
- 2) ISRAE recursion
- 3) Fun. rec. v.s. Prim. Rec. Efficient choice of Post-recursive

0 : T. written by E. Frank : Data Mining - Practical Machine Learning tools - 2000
 Therasos & WEKA data Mining System: I may have bug w/ o. its out. WOB. David Lindsey used
 Messengers to implement his "Improving Reliability..." paper. \leftarrow See WEKA (Google).
 Also: "date mining" @ Google: first use. (Download w/ 120 PDF of "R" handbook.
 Also "MySQL")

0 : T. Languages I've Considered for PHI: Lisp, Fort, assembly, APL:

- 1) Is APL recursive? If not, can it be easily "recursified"?
- 2) Their zero focus of recursion but don't use a stack, but do much faster than usual ~~stack~~-systems
 Are they capable of full recursion? - can they express AK funct?

One way:

Rec Rec.funct(z): Do $x=f(x), y=g(y)$ until $y=z$ return Rec.funct = x	Modulo: we can count to recursion as (j = count) $\hat{=}$ $x=f(x, j), y=g(y, j)$. - or various modulos: $x=g(y, j); y=f(x, j)$.
--	--

25 $d(13)$ is found ~~with~~ by the TM 25 is TM 26.

20 "25 x 14.40" : On "Hessian" : T. Hessian is specified in the "discrete" part of the corpus Dera. This includes any "dimensionality" into that may be needed for each coil/param of f. Hessian.

12 Midway... 32.5 3 1/2 5-55 2 1/2 595 66.2 kg 3.6.05 595 - 115 for 65.6 no delay 1 4 4.32 lbs .6 kg for clothes = 1.32 lbs!

Partial recursive v.s. Prim. recursive. In fact, we never use part-rec functs: We use part-rec functs w. a "Computation Bound". This makes them Total rec - but are they really prim. rec (i.e. recursively Do loops?)

Look up Chris Moore's website! He has a very clear paper explaining prim. rec. v.s. part-rec. v.s. total rec. functs. My idea of prim. rec. involving Do loops is a bit unconvincing: I think this paper clarifies it. If I can't find it; look up "Chris Moore" papers in "PS" directory.

On "Phase II": Instead of the way I now have it planned! ^{Not such a good plan! - see 34.34!} Could we teach TM what option is, by example? At first I thought TM would simply try to "imitate" v. "example" is miss the essential point. (I'm not now clear on why this was so). I later found that I could modify c.c.'s of various insts. so as to ~~make~~ make "copying" not an option: In general, I think I have to review this.

Say I show TM various OZ prob. defns ($M_i(x), T_i$), in either simple or complex form - i.e. the value of x / or T_i actual prim. rec. bound. "optimum x ". ^{UNFORTUNATELY} From the error of Phase I (act z, TM would not be able to guess the "optimum" concept ~~correctly~~).

Say we give TM several OZ techniques for a given problem. All but one ~~have~~ are negative cases (I think we ^{can} get to Phase I to use neg data)

34.31 is about the best (treatment/understanding) of Phil Davis for, but it's still "incomplete" (t.d.t.t. 34.34)

33 [We can derive various OZ methods: (Avery) (zero no. in "Literature"). Perhaps use one of these (not really optimum) OZ methods on 34.34!

One kind of OZ is a so-called "Hill Climbing" problem: We try several PSM's on a PSM, & we get some PSM's. From this info we are able to construct new total PSM's, & their PSM's, and so on to "peak" in available time.

T. problem of: Given a bunch of PSM's, find PSM pairs, to guess the new PSM by foot, → 37.00

12:40	1:30
2:10	1:05
2:07	1:16
1:33	2:10
6:23	2:05
8:32	

Phase II (cont)
36.40 is a standard H.c. / G.A. sub problem. We can do it our fancy way; to approximate
 + FOM \Rightarrow PSM; Prob relationship by a test function 2/0 an invertible function, so
 we can get PSM's by FOM; This trial can be used in G.A. or any other H.c. Scheme.

Another Approach: T. max of a finite set of nos. is a clearly defined function; could we
 ask TM to extrapolate Phase to an ∞ of nos?
 T. problem seems to be "Recursive": if we had a 02 algm, we could use it (v. 34.31-34)
 to get a better 02 alg. So perhaps we write v. 34.31-34 as a 02 algm. "In power" —
 that could be applied many times to some initial 02 algm, to get better & better 02 algms,
 .09-.10 may be adequate! Go Over 34.31-34 & v. 09-10 to see if it really can work!

Eventually, we would discover, for long-time use, a v.g. 02 algm for 34.34 & this would solve
 f. problem. So (at present), I think .09-.10 (v. 34.31-34) may be adequate for Phase II.
 We note that the 02 algm. for 34.34 does have to significantly improve copying TM's info... since
 at first next algm. will be able to optz over only a finite set; then, eventually, it will have to
 optimize over an infinite set of PSM's (that it, itself generates). It is at this point —
 f. transition from finite set to ∞ ... that the "hand waving" occurs.

We would need a good TSG to get TM / to do this particular finite set \rightarrow inf set induction.
 37.05 36.40 — 37.005 is a view of Ph II that is a standard H.c. problem. Grammatized (induction (Augmented
 w/ context sensitivity...), PPM (augmented w. sub-structure contexts (products)) could be useful tricks
 for this, as well as G.P. ... or any other methods I'll be using for Ph I very such
 So, on Ph I very such isn't so hot, we are using it at a META level, to choose good PSM algms, so it is
 effectively Much Better

Note that both Ph I very such & Ph II (second order Ph I) such are capable of creating
 new object w. by expected Score (or "FOM") — but, presumably, f. 2nd order Ph II
 such is potentially better. It should be poss. to go up one (more) (levels) higher ... to get
 even better PSM's: but it may be too time-consuming for the amt. of "goodness" obtained.
 Note that all 3 methods contemplated for Ph I very such (PPM, CFG, G.P.) would work
 as well for Ph II PSM (creation/such) — So any improvement in Ph I very such is
 usually applicable to Ph II PSM such.

So: .09-.10 may be fine, in general, but using Ph I very such for Ph II PSM such/creation
 is a simple, pretty good soln.
 A big problem is (evaluation of) FOM for various PSM's that have been created!
 (most all PSM's we create will have zero SSZ) So how to get FOM? I pass by
 an inductive process, we have lots of (PSM's, Prob, T_k) \rightarrow FOM_{ijk} empirical cases,
 we just have to find a (algm) that extrapolates to new triplets

Aug 777

Send stuff
on cover business.

So 37.21 it is probably a more or less adequate Model for Ph II. 34.31-34, from
Perhaps 37.21ff. I suspect that it may be close to what I wrote in IDSIA Report: -
except for variations of using PPM, PSG-ding, G.P., for Phase I regy such as ^{also for} 34.34 such
of Ph II.

I do want to see how far I can go w. those 3 "weak inducers" for Ph I regy such
Look at Dem of Ph II in IDSIA: See what has to be added. (Corr Rem) → 30

The points of 33.00ff suggest that Ph I may be not far from theoretical realization!
perhaps TSO's are most critical, diff't. problem with now. Perhaps try to track
mentioned in Sol 30 of taking a diff't problem & working out a "concept net" to
solve it: This could be done for D-induction ("Math" problems).

RE: TSO's: At one pt. I wrote a bunch of things but I expected TM would need/want
to be preparatory to "Ingr Algebra".

An alternative approach would have TM Inv. y. meanings of various commands in Maple/M-headers, etc
& be able to execute those.

Re: Looking at IDSIA Report: Abstract: What was described Phase II only:
Ph I was not described in "Abstract":

Re: Abstract: Some up & down can follow each problem Solus but w. each/problem, we may
want to support it. PSM set. in the direction of that new problem.

Was there anything in the report that told how PSM's were "added/deducted"?

I think that a real understanding of Ph II would occur if I could write
a TSO for it that would really be able to "discover" good PSM's: (Phase Good
PSM's could be normal part of TSO - designed in by "trainer")
"Logical Analysis" must be imp't. part of this, but, I don't know how to tell either!

→ It is possible to "logic" as a kind of problem reasoning? → Perhaps!!

: SN In Math, Proofs a thing called "Non-standard Analysis". calculus & continuous concepts

using only discrete concepts: Edward Nelson devised a PC scheme for it.

1) Non-standard Analysis — Nelson; PC.

Look into N.S.A. It may bring a nice continuum below. Strong output & Num output of \mathbb{R} objects. T. models I'm using are not so hot!

2) How PH II is into PH I.

SN PH II. can be regarded as a slight general of PH I initially may!

PH I initially, solves only probs that are Max. of PC. It does this by looking at functions that have been good in past, making S-program. for these funcs, & using this S-program to generate func's!

func's in L such or BU mode.

We can do the same thing in PH II, but instead of PC's to max, we have to maximize

1. FOM (log of market) for $\sum PSM_i$, Prob, T_0 . Prob, T_0 are fixed, so it's a d.f. over

PSM_i . And PH II PSM_i should be a set of all poss. PSM_i , just as in PH I,

1. argt. is A, which is the set of all poss. feasible solutions. → 30

Some poss. differences: Th. goal for PH I may be linear; Th. Goal for PH II

~~is not necessarily linear~~. is not necessarily linear: it is, for CE problems, & expected value of

only Goal — is that Goal need not be linear. It may be actually mandatory that

we linearize it. Goal before trying to maximize it!

Any way, for this Phase II induction we do have a suitable corpus of $(G, PSM_i, Prob, T_0)$'s

(30 R) PH I problem, since PC's are assigned to these A_i, P_j pairs.

We still need something like expected value of G — which does mean we need it linearized

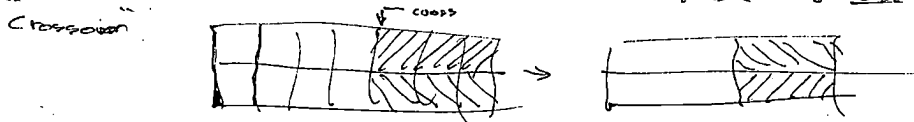
if we are to try to find optimum soln. to it. (PH II is discussed/mentioned in footnote in I & 2's report.

For Phase II INU problems (usually solved as DP probs), Th. Goal is automatically linear; we want to minimize time to solve.



Onishi, Ryan, Collins
: BFI (Boys for Ireland) & 3 Irish authors: Use of each their form to generate ~~random~~
DiPS Grammar Discovery / Grammatical Substitution

Searches (usually PSM or functions): Use in Genetic Algs: I think it's too slow, Poor



As to search, usually everything is & f. of crossover pt. is essentially variable, because it changes in numbers after "crossover" and ^{usually} for every different form the chances of making better crossover.

Occasionally a crossover will "preserve state", in which case we will get real "crossover".

Their use of "wrapping" is eqv. to having the finite chromosome repeated no times. So each chrom. is as long.

It would seem that Koza's "GP" would be uniformly better. → 40.00

and invest
Patterson ^N 97
Cultural Caching Alphas in C by GP

20: (50 sec 29.40) : ~~They~~ provide themselves on being closer to Natural Cognition - ~~but~~ P₂'s is good, only if you understand R. before your copying. It may or may not be appropriate.

In the set of papers in "Forum Evolv" folder:
or contains most info - 13 pp long! other papers w/ 6 pp long.
2. ps discusses "Pruna & Duplicate" 2 operators that make it default. from novel G-A.

In this paper Ray also discuss their "Crossover" method.

"Interacts" are bytes of chromo that occur after a function has been completely described.
By Pruning, Ray became better "inbred" - so "crossover" is more likely to do something.
(A crossover at a "inbred" produces no effect at all on individual).

By "prune" Ray uses "Byte" in chromosomes

That discussion of "Gene duplication" sounds insane ^{if my interpretation of what Ray}
over says, is correct. Any dup of a gene usually will produce a completely different effect than it did originally.

They say (7. ps) that their Caching algm is clearly better than Newsworld by GP. (!)

Sounds unlikely! - Unless I seriously misunderstood what Ray is doing!

5. ps "Under the hood" analyzes how it works (is opposed) to work.

I collected 6. ps is a WinZip file - both is 5. ps.

4. ps is very short nothing new to Paul.

6. ps is big file but only 7 pp. 2. ps is 11 pp - Most info.

20: 39.27 : Actually 39.30 is a bit deceptive: waits to compare for the needed induction (5, 0, 2).

What is PHI output is ... ~~is~~ 2 pd. on ~~PSM_i~~ { PSM_i, G }. The actual form of PSM output may not be useful. i.e. we may have to put in various trial PSM_i's and various G's to get PC's out - or may be not so useful forms.

What we'd like is either we put in a PSM_i in a set EG for output, or ideally a list of PSM_i in EG order - largest first.

This "Form of output" is certainly a major problem in all stochastic calculations,

Getting what we want ~~could~~ could involve a lot of (expensive) search! → 4.00

30 [SN] On nyms (or whatever) of unusually low empirical freq. I never really found a good way to use data of such nyms: An Alternative View. If a nym has "unusually low" freq. in a given context, then this is a hint to look for other nyms at that context having unusually high freq. - this will ↓ length of code

0: (40.29) ^{Spec} : Actually, that Big Search Cost May not be so bad!

It may well be that we ~~end~~ end w. only a small no. of PSM's, but we do constantly spend time looking for (new/better) ones.

[SN] Back to an old Q: Could we teach TM what Optim is by a suitable TSO (by a suitable set of examples)?

One kind of TSO: An OZ prob dec'd given G, T_0 . Paired w. it is a pair to solve it. We give Phil many of these pairs. Would it be able to induce what was needed in A_i ? (Q, A)

A trouble is: that the examples given are not nearly the best poss. ways to solve the particular

problem. } Perhaps Main Q: Can T. transfer Princ of a good recursive defn. of Optim? } CRITICAL IDEA!
If not, its unlikely that Tomo can devise TSO for "Optim".

[If I have a clear idea as to what I mean by "Optim", it must be poss. to define a recursive function to express it.]

→ Hint: That Phil doesn't understand Optim! It's soln. methods are indep of "Horizon", T_0 .

0: How Am I advanced in TM, since Sol 89?

1) ALL PSMs from "or a name of technique of Phil: found a unified set of soln. of DUV, OZ probs. ("Mixed ~~code~~ Theorem?") MCT?

2) Understanding of how Logic forms: (42.13)

3) Understanding and proof of convergence of "Big Corps", "QA Corps"

4) Use of QA form for Phil complicates TSO design

5) Understanding of what Allys might be suitable for probs w. string string v.s. Numeric Numeric output. This still needs much work. 3 forms of Numeric output (all have Model).

6) The use of PSM's PSB decy are but slight improvements in my understanding of how to search for solns in Phil. G.P. may be of interest (partic. body for BU search) if Star-M6 works. - But certainly I was aware in '87 that such for sub-firms dec imp.

7) Understanding of 2 modes of Search: LS v.s. BU Search! when each is appropriate

8) Listing of characteristics of STRONG AI. That only ALP seems to have (frustrating)

1) Accuracy for small size (as well as large size) 2) Diversity of small codes 3) Turning complex into partial dec. funcs

9) Understanding as to why part. dec. funcs may - They have hyper PC that corresponding prim. rec. funcs.

BIAS ASSOC. w. SEARCH ALG. is (O.K. / messy): (27)

- 0: (43.25) This will begin to work only when the O^j is quite large.
- 01 Another posy is to use Several "good" O^j 's ; construct a common Grammar or use them w. Augmented PPM, ~~but~~ as elements of an "ordered corpus" but give each O^j its proper wt. (i.e. \approx to PC it processes corpus). This will mean that the individual O^j 's don't have to be so large, but their total length should be large.

I guess that a smaller more attractive way using all previous O^j 's as corpus

Somewhat. Set of O^j 's that solve the "nat" seems more "natural" a corpus.

An apparent problem w. .01 is that, it may (comp. w. many || O^j 's) not present sets of useful sub trees.

A related idea is oops: Problem sets given in sequence. TM solves each set:

1. into sets of probs in each set are tightly coupled, but data betw. sets much less tightly coupled.

Also note in INV prob. we are looking for a single alg. fact can solve all INV. Prob. \leftarrow This is partly OK.

- 13 Writing a constructor of optimum $\left\{ \begin{array}{l} \text{we do want different coupling betw. different problems, but in OOPS, it's us fact} \\ \text{to couple betw. different prob-sets is too much - 2 is only of one kind!} \end{array} \right.$

A ollar (perhaps adaptor) view: T. corpus is series of pairs ~~(O^j, A^j)~~ Express that data corpus in more compact form;

Recall is much redundancy: $\{Q_i, A_i, O_i^j\}_{i=1}^n$ does express all of the info. — T. Q_i is:

does $\{Q_i, A_i\}_{i=1}^n$ O_n^j really have betw. info? — i.e. are all O_i^j ($1 \leq i \leq n$) uniformly worse than O_n^j ?

46.20

In discussing incompatibility, diversity, etc.: refer to total wt. eq. (3) factor

$\sum_i 2^{-|M_i|} \xi_i$ is BIG if poss. (Not exactly rite: we should also include ξ_i : No diffy!)

prediction... but this actually puts us into trouble: it enables us to bias f. search! (27)

So using .21 as a criterion would seem to get around the search bias diffy, but it means we "One shot (try)" ($z, \xi \geq$ shot (try is o.k., hvr).

Maybe var. z, ξ , but sum over all "next symbol" posys.

27 (22) AAAA! Break Qw! Any Bias int. search is Legit.! The bias is int. condn. of

to Apply it: search alg: T. combination of f. z, y results in final effective "sprid"

The user can bias "effective sprid" by manipulating theoretical z, y of search alg.

- There are 2 sources of bias are unavoidable & really desirable, since they never way to meet
- his apri info $\left\{ \begin{array}{l} \text{Hvr, use of a universal time for apri, does constrain sprid 2/05. The modn of f. apri due} \\ \text{to Search Method, is not so clear! — As example of extreme effect of search alg. —} \\ \text{Consider an algorithm such alg. that rejects a (1) today that gives 2/05 for the "1" continuation!} \\ \text{I shot (try).} \end{array} \right.$

34 (22) If we include Prads, we want to include/codes for the next symbols. In e. case of OSL, this would \rightarrow 47.30

mean that we write add one or more M_i 's that had definitions for the ngrams involved in poss. continu.

Micro Cosmic God - Fred Sturgeon

Assembling Set AP 41: PP 46-68

46 James Kidder.

Banner: Konant

48. Factory Made Food
AC Day

Liter Pump?

Syntax Geography:

Propolar at hach

U238 Power

49 President Konant at Bank.

501 formalized pe & reduced it to such

low terms that no know which experts to not try

502 Oxalic Acid \rightarrow Cutting Apart for Blood?
Elavir Made Sleep Under

511 speeding up life/metabolism x20

Generation in 8 days: life span 15 days

6 hrs gestation, eggs hatch in 3 hrs.

Sexual maturity in 4 days.

Female \rightarrow 9 eggs.

Males died 2, 3 hrs after hatching.

3" long. Much N₂ in atmosphere & Power
So couldn't live in Earth Atmosphere

512 4 races, inviolable diversity?

Neoterics

521 How Kidder Got People solve initial problems.

Wars betw. sections of Trib.

522 Making Story Alleg: —

"Rule of Reed".

53: ^{Kiddier} looks like picture of Calva Floors

54, T. Tealotype is a "social sociologist"

in chron to be by Neoterics

56.12 Power source, on xmtt

57 How it works! @

p52,21 Subject N's to Fed warfar of various extremes, to
problem to invents add phve devices, to survive.

So Basically, Each Threat is tot entire race & f. individuals
can cooperate to solve it.

The "Crusher" is a. As alloy solm. is a similar problem.

522: "rule of fear".

53: The teatotype w. Commandy: They had to develop

what was asked for — full time no other projects by anyone —

Or they would be penalized much ($\frac{1}{2}$ population killed).

3 stages of Education off N's! (Neoterics)

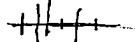
1) Organic Evoln: \rightarrow Diversification by Geneties,
followed by selection by Catastro plus.

or: slow catastrophe w. benevolent "teaching".

2) Impending Catastropha: Culture solves problem.

3) Culture can Read instructions (like "language")

If instructions are not followed — Catastropha follows.



20 1) has to get Culture to state where 2) is possi.

2) or 3) has to furnish N's w. language.

Lang. can either grow up in culture or (better)

be externally imposed — contracting much into —

Excuse for "God" to understand & give

orders in it!

Getting 1) Org Evoln to get lots of individuals to

cooperate and to form a culture & ^{adaptation}

(w. partners a.H. help) develops ^{adaptation} ~~language~~

language. ^{interdependent symbiotic} (symbiosis)

Getting some cooperation (like symbiotic species)

may not be so hard; but getting a group of individuals of

some species to cooperate on common projects,

seems diff!

A first step toward cooperation: "Mating":

by having "Mating" done by ≥ 2 individuals, a stronger

step toward "Community" (a "culture") is taken

Family raises children: Family is embedded when tells "mate".

1) I want individuals to (work on/solve) problems

2) I want ~~some~~ individuals to be able to cooperate in

solving problems

500 Neoterics is an "A-life" problem, not "G.A." — I'm not so sure! They

did have viruses, but it looks like G.A. rather than A-life — @.

These were not a lot of animals/plants, interestingly w. G. A.D.O.S.

Another approach: Have Pool of problems

accessible by pop(ulation): Each problem

has its own stated Prize for gain.

Have probs, but prizes. Prizes are

31 "Money". Money can be used to buy

"production" services. ^{Endowments} ~~tail~~

can cooperate to solve problems

in share prizes.

Money is "Carrot": T. "Stick" is

death of individuals in genetic kin?

if indiv. does Bad Things.

"Bad" is Not following orders.

The Bus system has "Buy" of all computer

parts: It does what you tell it to do,

not really what you want it to do! \rightarrow \$6.00

The "God" words "reward" is well-timed (1st)

STM

"Micro Cosmic God" (cont.) - Sturgeon

0: (45.40B): As is, the collection of "indivs" as pems. How do they "interact"? How do they reproduce? It may be that after the system has a sort of suitable communication, between individuals, it may be that little advancement is then "organism" most advanced cases thru development of technology in the system of individuals.

Tho, we could have "organism" continues - it's easy to continue to "improve" (i.e. pems that represent "individuals").

Indivs: can communicate w. each other by "messages" in some common lang. This is common lang. But "God" uses for no reason (Orders/assignments).

Note that we may want individuals to be similar enough so that 1 communication can occur
2 that observations by 1 individual can be useful to another indiv.

Could individuals be like objects on pems. "Objects"? - They would be a subclass of "Objects".

1: (45.3) Re: Money: Can also be used to buy things that might generate more money. - unclear how to do this, here. The advantage of Money is that individuals can have many different interests (→ "diversity") & be rewarded for them.

30: 44:13
SN

Problem Pool idea: [in OOPS! Consider the world of Hamlet! We are given a seq. of problems w. 1, 2, 3, 4 ... ~~discs~~ discs on the pikes to be mounted. J. considers the sequence of final pems that solved the problem {programs}^n, pems solves a disc problem. So his corpus is {pems}^n. The only way to address is freq. of use of each inst.

Two badnesses 1 This is too much a copy to the help person (at least use some "context")
2 Lsearch doesn't work that way anyhow. The ~~adoption~~ corpus that guides adaptivity of Lsearch is not {pems}^n, but the pems that found the pems's. J used to same pems to find all off. pems's - i.e. to "Lsearch"

"More Predictable Statistical Data Compression"

Steven J. Ross | Sep 2003

≡ D/PSI PPM Argumental Basis.PS 3/24/05

Seems Readable! He seems to get maybe 10% more compression than PPM.
 (Much more compression for LISP code).

$$\frac{1.934 \text{ kbits}}{2.160 \text{ PPM}} \approx .895 \text{ so } \sim 10\% \text{ better}$$

$$\frac{1.571 \text{ kbits}}{1.972 \text{ PPM}} \approx .796 \text{ so } \sim 20\% \text{ better}$$
 book = Carney/Carney
 LISP (page 2)

He did much worse for some EXCEL files here: - Part was only case he did worse on.

Blurring prior (pp 60-61): I could do this much better w. probably simpler results, by using a β distribn. β smoothing is w. a normal dist is a β distribn of constant vars, but w. a mean that corresponds to point being smoothed.

I suspect result is almost trivial! We end up w. a β distribn w/ some more as before but var increases (just as if we smoothed Gaussns. w. Gaussns is $\sigma^2 = \sigma_1^2 + \sigma_2^2$.)

Also this "blurring" stuff can probly be ~~done~~ exactly in ALP.

Anyway, his discussion of criticism of PPM is probably well worth reading: Put my mind in pool

— Ross w. 20 pp, but I don't have ~~time~~ time to read all of it! (state to improve PPM)

He also discusses M. Li's DSCA compress!

He used a binary alphabet — which would seem to eliminate "escape" problem. (No, it doesn't at all!)

Before reading this I was very carefully familiarizing myself w. my "solns" to PPM problems.

D/PS/ MADL-Search.PS
Jan Davidson

3 (24) 05
3 (24) 05
3 (24) 05

Deals "Gibbs sampling" & other such methods, with statistical models. (Stack complexity?)
 I'm not sure I understood his crit & my reply!

He seems to be using a ALP! Sub-over Models. (Stack complexity?)
 Perhaps some hierarchical version of Lecture I!

Re: Search alg ^{influence} on ALP & Adaptive Algorithms: We should at least use strong algorithms that \rightarrow ALP as cc \rightarrow oo.

But, we can still have an algorithm that's bad for $T < \text{Adaptive}(cc)$ seconds \rightarrow const. method for larger P.cc!
 so it \rightarrow ALP as cc \rightarrow oo. But, seen Alg. is chosen by User, so he will avoid such methods that seem to change as cc \rightarrow oo.

(17)
(97)

0: : **More Remarks on Ph II**. We can view Ph II in a narrow way: i.e.,
T. general improvement of Ph II's search — or, more generally, improvement of all/any
PSM's. Usually (?) the more general problem is easier; also corpus has more "diversity"

13: So cross-fertilization is possible before PSM improvements → (07)
[SN] Another thing I hadn't considered much in open problems: That TM could
watch GC() operate — so that it knows how G (i.e. optimization) works. Open PSM
is (poss./imposs.) in GAs, it is assumed to be impossible.

07: (03) → Having a large corpus of PSM's certainly makes it easier to construct a useful impl!
S- Grammar for Ph II.

0: That method of doing Ph II in IDSIA reports is "Not Bad", it's certainly only a single way
to ~~define~~ "define" open for TM. There should be a more general way: — but it eludes me!

13: So it would be good to have Ph II solving a great diversity of problem types
(border is while) it works on improving Ph I's such. (30)

→ G, Z, Skuitdhu: 2004 "Adaptive online allocation to search algms."
D/PS/-3/25/05: "My title "Adaptive CC allocation - Juergen-Doll".
I have 2 versions: Powerpoint presentation in PDF 38 slides
PS. 11 12 pp - actual paper.

20: This is very close to my Ph II problem. Their method could be much improved
(i.e. a set of PSM's is limited and for each run) but probably can. is useful ideas, perhaps.
They don't use exactly like Gore, etc, but --- Oh.

It can be used to "start off" on a generic Ph II.
(P. 68 of Juergen's "DOTA" 29 Sept 2004) such proposal to FNSMF ideas for help/suggestions for improvement.

30: (13) → Ph II: On improvement of Ph I's Meta search: If (i.e. long. we use to generate trial PSM's for
PS meta search) is universal, then its meta Meta algms don't have to be v.g. If they are
poor, it will just take a longer time to get a good such algm for Ph I, but we will get
it eventually. So we will get a "bootstrap" effect.
T. usual way Ph II works: we have a universal s.g. raw using to 2000's, a set of
good PSM's. One impl goal is to find a good PSM for Ph I's such problem.
We have one default/initial gsm for TM, say Aug PSM or Aug PSMing if we want to improve it. (Already,
Aug PSM & Aug PSMing are quite close to universal)
So we use (P) M generating Grammar to create seeds for Ph I such. We use (50.00)

On ALP is uncomputability, etc.

1) T. uses definition of "primitive recursion" is Rest of only uses composition
"primitive recursion" — (i.e. $h(\vec{x}, 0) = f(\vec{x})$ $h(\vec{x}, y+1) = g(\vec{x}, y, h(\vec{x}, y))$)
function
known
known

It does not use "μ recursion" : i.e. $h(\vec{x}) = \mu y f(\vec{x}, y) =$ smallest value of $y \rightarrow f(\vec{x}, y) = 0$

"total" in prim. rec. is also projection: $T_i^n(x_1, \dots, x_n) = x_i$ ($1 \leq i \leq n$)

So ~~it~~ it would seem that μ Ackermann funct would not use μ recursion. b/c prim. rec.

Since Ack funct. is a triple recursion, it's not immediately clear how we get it using μ .
Perhaps 03 is used!

It would seem that any sequence of functions trying to approach ALP, would each be computable & very probably prim. recursive. As soon as we found a funct that had no value (or seems to have no value) for 2^N args, we would drop that function.

So if we tried to prim. rec. functs in their recursively enumerable order, we should probably approach ALP.

No! — Because there are ^{total} recursive functs that are not prim. recursive. (Kleene, however —
e.g. Ackermann (0))¹. So ~~intuitively~~ probly it is "not bad".

If we tried all the prim. rec. functs to approximate ALP! eventually, ^{function} f_i decrs would get so large that we wouldnt have to consider them? So in theory ALP using prim. rec. functs is

well defined & "computable" ... but in practice, no more computable b/c as we used ^{partial} ~~primitive~~ rec. functs!

On my list distinguishes ALP from usual prim. rec. such: ALP will try args & functs that will occasionally be beyond CB, so it will not use those functs on next round of LSrch.
Hvr. using prim. rec. functs, people usually consider only functs that converge & they devote a lot of time to prove to enable them to converge. — They don't delete functs because they take too long.

E.g. on lists (~~primitive~~ recursive) all prim. rec. functs, I don't think they answer by pc ever takes much time? ... Maybe not: 2 prim. rec. funct of n can take n^2 operations to compute (e.g. $X!$). actually $> n^2$ operations as x & may be $x \times x$. x^x can take a lot of ~~time~~ ^{time} for x .

Basically ALP does μ recursion & ~~prim~~ prim form Rec. doesn't.

Since ALP only approximates μ recursion — Could approximate prim. rec. simulate it? — Also would it simulate it w. much \downarrow of ~~pc~~? (No theoretically tho into int. ~~approx~~ approx μ recursion, should be same as into in prim. rec. simulating μ recursion)

I really don't understand this!

STM

to evaluate these cond. Evaln. can be done in 2 ways:

1) Direct: we try to cond on PkI's corpus - (or a sample of that corpus)

2) We use a faster algm for approx. evaln. (This is from IPS's Report):

This has developed an algm that can look at any problem, PSM part is got a Gene for that part.

(This is faster than O(n^3) operating on full PkI corpus, but much less accurate @. We have to balance this inaccuracy with O(1) operating on samples of various sizes.

But 1) is perhaps easiest, & we can kill climb on it.

A big plus for 2) is fast to evaln Algms. works for a great Diversity of PSM's & problem types.

This Diversity is imp't in suggesting novel solns. to problems.

It means that PSM generated by Grammar gives a good diverse set of (inds for all kinds of probs (not just PkI's/ problem).

Also by mixing info from various problem types, it has a larger size than if no used info from only 1 type of problem to optimize PSM's for that problem. (comp. PkI problem).

10

20

CONTEXT in augmented PSM: Consider "External context" (like's problem in Math, Physics, Chem, SM, Sociology ...) we can assign a multivalued symbol to this info. How it fits in w/ other contexts in PSM is unclear. I did work out a way to work in sub-tree contexts. Perhaps it can be applied to α as well.

"Regular Expressions" are easy to recognize (I think). Are Reg expressions simple a Finite state Grammar? Can they help discover common sub-trees?

16

In expressions like $+ , A , B$ we can recognize that A is a sub-tree. Because to string "A starts w/ (and ends w/)" & parentheses are always "Balanced". Balancing is easy to do. We open w/ "(" and $x \rightarrow x+1$. every (increases x by 1, every) $\downarrow x$ by 1. If x ever gets < 0 , we have a syntax error. If $x = 0$ before end of string A, then A is not a proper variable. so (gives inc x ,) gives dec x ; when $x = 0$, if we are at end of string A then A is a proper variable, otherwise, it's not.

20

In .26 ff, if we use Polish or (Polish) we should be able to do same thing.

In my "One, Two tries" psm, I have these nodes and info about nodes. Could I use that info to detect common sub-trees? (or a particular sub-tree "Context"?)

57th

00: : **SN** On S-CFG deriv: Stolcke uses 2 operators Merge, Chunk to fill (L)RHS to prod.

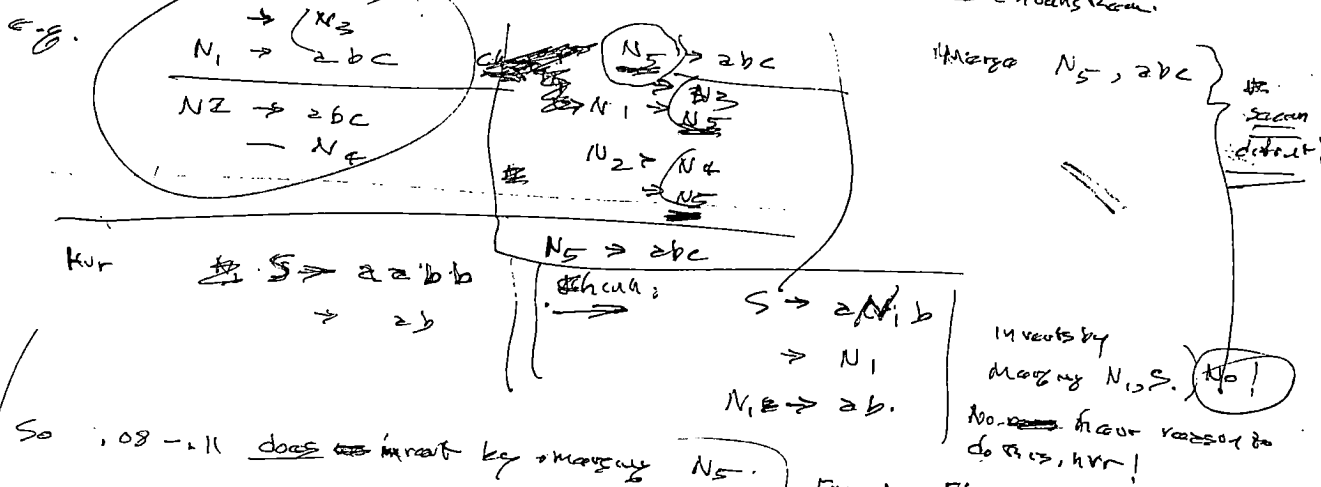
If should also be possible to do inverse operations: (1) inverse new NT's is what is chunking!

So Merge & chunk may be inverses! Maybe no need for new operators?

Any Merge can be undone by a suitable Chunking.

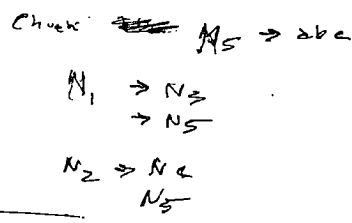
Any Chunking, however, need not be inverted by a merge!

Say 2 concrete chunks, find 2 RHS's of deriv NT's in the chunks deriv.



08
10
11
13
15

So 08-11 does not invert by merging N_5 .



Chunking chunking success replacing N_1 , sep, by its derivation in all cases: This means that N_1 is not part in recursive down.

4 ways of chunking is some form of a Non-rec NT by substitution in all cases.

In fact Merge, would be breaking up the NT into 2 NT's; (looking sub at rec all cases!)

What may be complicated is since we have have to double up every

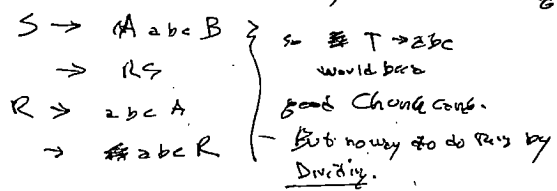
role in which the "derived NT" occurred. Have to see why this would cover PC.

But it must sometimes, because some merges & PC!
 \rightarrow If a NT has many RHS's, there are many ways to divide it.

So "Divide" is inverse of "Merge"
"Expand" is inverse of "Chunk" (= "Derive")

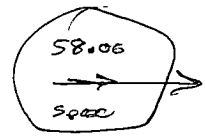
Chunking can sometimes be a kind of divide. — Both cross to new NT's.

But all chunking can't be done by dividing! e.g.



Merge is Expand & no. of NT's
Divide is Chunk & no. of NT's
Chunk is Divide

Chunk & sometimes = Divide in non-recursive
Merge " " Expand " " " "



30

Submitted to Σ -XIV: 17 Febos' had at Royal Holloway! (1)

SN: Gen Models: Paper D/PS/VCOM - least-complex model of lang: 4/2/05
Says VCOM is K-comp. for "concepts": Also interesting (prouder) results on Σ of finding ~~models~~ inductive models. He seems confused about Hofstadter's scheme
~~Robustness of~~ MLP conv. time - being related to P & B's results on C.C. as function of problem complexity. He may be too closely tied to poly vs. nonpoly categories.

URL: Patterns in Mathematics URL: documents.wolfram.com/mathematica/book/section-2.3.1
2.3.1 thru 2.3.14: It's notation for "Patterns" seems very useful: Probly I'd want to learn these (definitions/notations)

I'm not sure that the notation is "consistent" (or that it means what one would expect it to mean)
But that these kinds of "patterns" are important, & count!

Could we devise a test so TM would discover patterns because they were useful, rather than because they were a soln. to a simple QA induction problem on "input data"?)

Anyway, for TM to invent "structural" rules/equalities: would be a good way for it to learn how to do algebra, & much of math.

Can all of this stuff be done easily in Lisp?
Do try to Gen. Mathematica Patterns" Dennis Ray, to Prouder so that they could easily be derived
V.2. Ray's "Primitives" in 4. lang.
SN: D/PS/ Grammar discovery 4/3/05 & SCS may be improper: pg 300

$$x = \frac{1}{1+x}$$
$$x^2 + x - 1 = 0$$
$$\frac{-1 \pm \sqrt{1+4}}{2}$$
$$x = .618$$

Improper means Σ all terms < 1 , due to many infinite strings that never terminate.
They get eq. $S_{n+1} = q + pS_n^2$ which leads to Σ poorly bound < 1 if $p > \frac{1}{2}$
Note I might get \downarrow

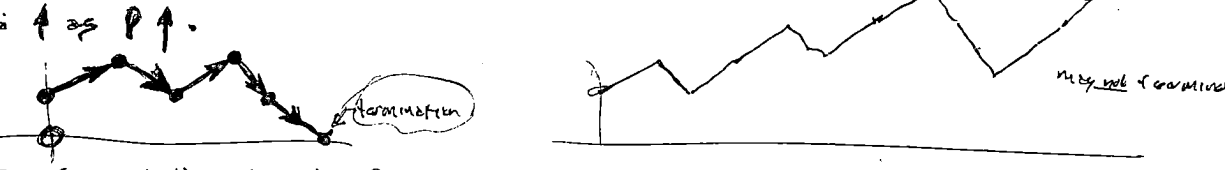
T Grammar $A \xrightarrow{p} AA$
 $\downarrow \frac{1-p}{2}$
If $p > \frac{1}{2}$ then we will tend to have intermediate strings of terms A having w. $\frac{1-p}{2}$ A 's in them.
(i.e. non-terminating - so we $p > \frac{1}{2}$ we will have a central

$$x = \frac{1-x^2}{1+x}$$
$$x = \sqrt{1-x^2}$$

Bob's drawing

fraction that will terminate is then fraction p as $p < \frac{1}{2}$.

An exact Random walk model: up steps = $A \rightarrow AA$; Down 1 step = $A \rightarrow \epsilon$.
we start at +1. If prob of $A \rightarrow AA > \text{prob of } A \rightarrow \epsilon$, then prob of never hitting 0 is $> \phi$.



Still, I'd like to know how Ray got the eq. (22) $S_{n+1} = q + pS_n^2$

It is clear that the basic thing is true: S.C.G.'s can be "inconsistent" -
But Ray says that Grammars that try to model real data will be consistent (non-terminating strings are not in corpus!)

Q: Say we have 2 corpus Generated by 2 distinct Grammars: It would be possible that

5-7M

SN, 30

γ Gamma | λ lambda.

20

the semi-martingale induced by f . (a Consi grammar (for Real corpus) would \textcircled{a} have a shorter
down length than the Consi grammar assigned to δ . corpus? ~~the~~ R_{δ} idea is that

the inconsi grammar should give v.g. relative priors of δ . class m to the corpus.

So what we want to make is $\left\langle \left(\text{pc of } \overset{\text{inconsi}}{\text{grammar}} \right) \times \left(\text{pc of corpus w/ } \overset{\text{inconsi}}{\text{grammar}} \right) \times \left(\text{normal constant of grammar} \right) \right\rangle$
wrt Real corpus.

would \textcircled{b} give a result different than $\text{Pr}(\delta) \left\langle \left(\text{pc of } \overset{\text{inconsi}}{\text{grammar}} \right) \times \left(\text{pc of corpus} \right) \right\rangle$
wrt Real Grammar \rightarrow i.e. would these 2 options exaggerate differently?

Maybe small deltas - only do fact to very large corpus!??

103
00

21

0

SN Admin Note \textcircled{a} :

When in College I had to ^{More rapidly} Make Real en size of Science - even when I
didn't understand part of it. ^(I would remember what I didn't understand, hvr) Outside of college I would often stop at a pt. & work
for it until I understood it: This could take a wk. (or ~~more~~!).

What I need to do now is: ^{Thru} Make the TM construction. Parts I don't understand;
write & clearly what the problem is, & move on $\left\langle \text{Unless it's a } \text{C} \text{ artificial problem } \rightarrow \text{so} \right.$
soo can't "Move on".

So here's ^{updated} list of Unsolved or "Slightly" ^{not completely satisfactory} solid problems: & notes to work on them.
Summarize State of present development

Tom
Murray
Doe?
Judy?
Alex
Ray
Grace

5

SN SM

In predicting a stock by its correlation w. ~~the~~ ^{an} pool of 1000 other stocks, a correlation
has to ~~have~~ give an up in predn by addn of 1000 drvr one can justify s predicting it as a
predictor! - Hvr long all cross predns in || has been at zero (it's all a pi).

This last isn't so bad. Most ~~drvr~~ ^{drvr} will have corrln of 0 any way

using a wt. for all 1000 drivers would give predn. of zero most of time? (wts would be by far
drivers that gave good (low) predns. - so maybe not so close to zero predictor pool?

This is
high - low
weight w.
= wts for all 12
predictors

Also, my idea of certain stocks getting "News" earlier than others. The slow stocks stocks
slow to respond are held by Big funds $\left\{ \begin{array}{l} \text{that can buy/sell rapidly w. moving mkt. (slowly)} \\ \text{they respond late and over a longish period} \end{array} \right.$
of time

STM

00

HA! By chance I ran across ~~98TM~~ 98TM 110; --- 132
 (Feb 110.00 is on the Mind Corpus Firm which I seem to have just discovered!
 " 132.00 is on SGA (Super GA). This idea was to model a GA core by a Stoch Grammar
 so the pc of the grammar was a function of core (or a function of pc of grammar was Fitness Func of the core).

In present work on QA industry, the core is a PE!

I once tried to fit a SGA to a core, but I really didn't know what

I was doing!

One of the (apparently) next tricks: If the core was not a pc, I would make the case count
 assoc with each candidate to be a fitness func. Of course, fitness func. had to be
 modified so that they looked like pc's — but in a common GA, making reproduction
 pc or "fitness" is common — so I suspect that the people that use it don't know
 how flexible it is! (i.e. that any function of a fitness func. is an "equivalent" fitness func.)

for the Max-M criterion, it is well to have a "linearized" core (i.e. fitness func.).

In ibid (129.33) ^{98TM} an interesting note: Given a large Corpus: One doesn't want
 to consider all of it! We have a certain cc available & we want to use it best to
 be on problems of most value.

I did get the "Encyc problem soln" (I don't remember just what it was,
 with problem & soln.). I think Proctor was that TM has this Encyclopaedia which is "free".

To answer Q's it has to point so particular: These Proctors have pc's (& bc's)

TM indexes & energy (a Meta Task) so that less cc is needed so that Proctor

21

Another Remarkable Old note: 3TM 425.00 - .06 Idea of using Maple/act, for d-induction
 problems for Ph.I. — objection was that Proctor inductions would not readily extrapolate to s-induction
 well (maybe Ray would)! T. way it works for s-induction we get functions that output a most
 likely "A". This is a "cluster center". The Random R inputs to 3 Induction modify
 the pm to give more distant members of the cluster. It would seem that the cores used for d-induction,
 would be good for getting "cluster centers". Testing final "R" values should not take very long
 because we are mainly interested in the pc outputs!

0

3TM 423.32 ^{424.29} on "Source (Swimming Machine)": Using BUS such! We want best 100 coins Proctor.
 Say we use Proctor (this augmented PTM) to construct a pd on possible trials. When a new problem
 comes in: we try all 100 "Best" on this problem. If solns aren't adequate, we use Proctor to
 generate trials until we get a satisfactory soln. We then pick 100 best coins for reference
 to be basis for the next problem. I think Proctor tends to keep by diversity".

I've ~~been~~ forgotten what the "kinds of PBS" were: Maybe one was a prop that threw seeds into.
 The other was a prop modified by experience, that is used to guide Levick (or even BUS),

As I see it now, there are several (>2) pd's of interest:

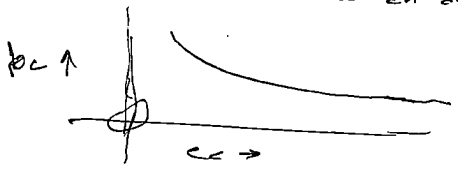
STM

- 10 : 5440 : 1) T. true expr w. zero corpus, obtained ^{from} reference vme, using CBDD. Two ALP.
- 2) T. initial assignment of c's to t. symbols of b. lang. used for (1)
- 3) After TM has solved n problems, we have a corpus, $[Q_i A_i]_{i=1}^n, [O_j]_{j=1}^n$.

This corpus via t. "two ALP" of (1) induces a new pd, which is t. "true & proper" for ~~the~~ next problem.

4) Since we don't have time for good approx of 3, we use PPM, augmented PPM, S-GFA, augmented S-CFG.

Some things I did work on: If we have a "not so good" method of induction w. a certain cc , we can ~~often~~ ^{often} ~~reproduce it~~ ~~again~~ (w. Boosky) to get a better PC values, but for hyper cc . [Assume for t. whole that this is true; It means that for each proby evaluator, we have an assoc set of expected pc's for corpus c & cc w. its own cc .



For a different proby evaluator we'd have a different curve: So, unless our curve was covered in a flat completely, there would be a crossover pt, where for $cc > cc_0$ we want to ~~change~~ ^{switch} proby evaluators.

For a simple L searcher for proby codes... certainly hyper cc / ^{almost} always gets better codes.

But consider PPM: w. hyper cc we could use lower contexts, but this doesn't help ~~that~~ much. (Using larger corpus helps, hvr.). Using more complex contexts, (subtrees) - we can use larger units of time to search over (larger sub trees).

Similarly, in S-CFG's we can ~~use~~ do fewer searches for better grammars, we can retain more ill parses. - Use context conservatively. Contexts can be of various kinds - perhaps contexts?

SN At many pts. in past I felt I had (as I do now) a v.g. picture of t. path of FM development.

One such pt. is 99TM 5.03 "PLAM"

How to use (Avg) 2^{-141} for prach/compression.

To compute the pc of the next symbol of a seq:

Computer the pc of seq. containing cont. S_i [$i=1 \dots$ Next symbol of alphabet.] using simple ^{as bit} ^(Dir's w/le) ^(Lars m/le) ^{Beun} ^(?)
Call this $P_i^{k_i}$ (i means sym. defn. of length i)

Computer $P_i^{k_i}$: This is pc of seq. ~~containing~~ S_i for next symbol, dividing a k symbol suffix.

$P_i^{k_i}$ sums pc's of all possl parses of i corpus, using that definition.

~~It will have a max value~~ \rightarrow This is when the suffix assoc w/ (k_i-1) has not occurred before in the corpus

The pc of S_i at that point is $\sum_{k_i=1}^{k_i^{max}} P_i^{k_i}$

This pc should be already normalized. If not then I have missed some code, $\leftarrow ?$

The actual pc of S_i at that pt. will be $\sum_{k_i=1}^{k_i^{max}} P_i^{k_i} / \sum_{k_i=1}^{k_i^{max}} (\sum_{k_j=1}^{k_j^{max}} P_i^{k_j})$

Actually the normalizing denominator in log is not 1.0 but if we have included all codes,

we can use, instead of that denominator, the pc assigned to the corpus w. S_i deleted $\leftarrow (2)$

It may be false. The coding for S_i containing the aux symbol S_i may be quite different from the codes for the seq. not containing that S_i .

At any rate, log is probably correct. The need to compute $P_i^{k_i^{max}}$ for all values of i seems like a ~~big~~ drag! There may be a way to avoid it! Would

using all be a good enough approx? We could try to compress, using it.

We could also study relations of log, to all: i.e. how does the denominator of log compare w. the pc of the corpus w. the lost S_i omitted

Anyway the denom of log increases the cost of system by a factor of "t" (\approx size of alphabet) ^(i.e. radix of the seq.)

One way to deal w. this is to use binary & (plate) - i.e. more symbols to compute? It's apparent \downarrow in cc is $\frac{2r}{2 \log_2 r}$ i.e. ~~we~~ cc by factor $\frac{2}{\log_2 r}$ for each symbol, but ^{is multiply} ^{prediction}

for $\log_2 r \gg$ many symbols. For an alphabet of ~~128~~ symbols (ascii), it's still a factor of $7 \times 2 = 14$.

There may be tricks to reduce this. Perhaps get rid of factor of 2?

In TM, I may be dealing w. smaller alphabets, but still $\frac{2r}{\log_2 r} \gg 1$ so I should use binary.

Use other log basis $b \log_b t$ min w.r.t. b.

$\log_b t = \frac{\ln t}{\ln b}$; $\frac{b}{\ln b} \cdot \ln t$ is min when b is min. $\frac{\ln b}{b}$ is max

So $b=2$ is best. I can't see how to do $b < 2$.

A closely ~~near~~ b is best for $b = e$, but

only by small factor $\frac{2.718}{2.88}$; using radix 3 gives $\frac{2.7307}{2.88}$

So using radix ≈ 3 is slightly better than 2. Using radix 4 is same as radix 2

$2/\ln 2 = 2.88$
$1.5/\ln 1.5 = 3.699$
$3/\ln 3 = 2.7307$
$4/\ln 4 = 2.88$
So the best is
$2/\ln e = 2.71828$

Hur. There is more: using binary may give cheaper operations in general.

NP: In considering various constant sizes: Only consider multiples of 8 bits (i.e. byte boundaries)

It's arbitrary that we're breaking up bytes would be useful here... The one could try it out!

STM

10: (spec) (5.40): The ongoing discussion ~~is~~ ^{is} directly applicable to sub-tree contexts.
→ We do have to find ways to desc sub-tree for decomp. ~~problem~~

Th. Method (Expressa for getting p's is for summing over all parses and identical to those for ordinary string contexts.

[5N] The horizon problem for RTM: My arg't v.s. leaving $Horizon(\in H) = \infty$ was that if we had a problem to work on in some (S.I.) ~~form~~ ^{form} to do, TM should always do S.I. first. If our horiz time T is do a bunch of probs in S.I., one would do all S.I. first, then ∞ problems. If $T \rightarrow \infty$, one never works on problems! So, ∞ is not good! A better approach, hvr, TM has many problems, but (Probs) has to be finished by T_2 . If there are ∞ problems, TM would spend almost all time before probs, in S.I. This would not solve probs very well, but in "long run" it would give good solns for probs \langle for $\infty \Rightarrow D \rightarrow \infty$. This would continue indefinitely, TM would always spend as little time on problems as poss. (so as to not get ~~stuck~~ for say 10^6 problems). We assume that mean ~~goal~~ is to find goal.

[5V] Mark Johnson: S CFG models of Trees: D: PSL = 4.12.04

21 This seems close to what I want. ~~But~~ ^{he} does use "Context" for probn, but I think his "Contexts" are too small. — BUT I must read his paper (20 pp!).
Also, I take it that Pinsky ~~only~~ ^{only} assume only one parse: (The corpus consists of pre-parsed \Rightarrow which is my corpus, also.

T. remarks on Context (2.1-2.2) may be the main point. They look at the corpus of trees. The freqs assoc. w/ each choice in tree are used as the pc values for a s-cfg. Hvr, I don't know if they use pos labeled trees, is pos frequent. POS transitions, — or just what they do!

30 I really don't remember just what I had in mind when I saw s-cfg theory as a great breakthrough! The basic drugs of the functions was the corpus. Just (how) CFGs would be good at deriving regys in that corpus, is unclear! The ass in that corpus are, indeed, ass of a CFG... but I know what the CFG is! — its Lisp or AZ or whatever I'm using to desc b. functions. So CFGD would seem to not ~~be~~ ^{be} relevant at all!

HA! A possl way to retrieve ~~from~~ ^{from} subtrees n to present one! The problem is similar to finding a certain pattern m in a dim picture. What we do is take a pattern of interest and convolve it w. the corpus. My peaks occur at pts n to the pattern. This need not

→ 59.00
Spec
→ 59.00

4.12.05

STM

0:57.40

Kieffer (ITC) Trans July 2000 (paper)
May 2000 (2 papers)

On heuristics for ~~rules~~ rules for compression & context.

NB: Stolcke has 2 rules Merge & Chunk. Some number with 100% of time, ~~rules~~ ^{Make} Previous link in STM. The inverses of these rules must occasionally compress. Try to discover conds under which rules & invs work. Similarly Kieffer (?) et al do "chunking only" rules. How many get 4 rules is not clear - it would seem that only 2 kinds would be possible.

They also have methods of extending 4. system to 2 dim. & context. Look into this.

They compress ~~down~~ their rules until context is "irreducible" - I don't know what they mean ^{by} that... (see?)... Actually, what they ~~say~~ should mean is

~~irreducible~~ Irreducible by any of the listed ~~rules~~ heuristics.

So Go over Kieffer (?) & Stolcke papers (2 papers in STM on ~~STM~~ CFB discovery) & see what heuristics they use for compression rules. This is perhaps one of most imp. ideas in "Hill Climbing Compression" (\equiv Coding & Recording).

On organization of "The Lecture Notes" - As of present LQ has:

- 1) Introd. What is being AI?
 - 2) Listing of Needed tools. (Lack, PPM, ^SPSG discy, GP)
- } this discussion can certainly be improved!

STM

5800

20: 59.40: In CAM: How do we arrange to get one context at a time? Actually, we are only interested

in the population of one of the subaddresses of the address locations.

It may be possible to get CAM, using a register Memory in an "inverted way"! - Hvr, only 32 (or perhaps 64) address lines. 64 may be enough for Zebra today (?).

5

0: **SNF** On context: Contexts are used previous corpus (+ external context).

Essentially, we want to predict next symbol of D^j : It can use ~~data~~

Process I to do this: Use area in $T_{D^j} = T_{PM}^j$ (= Process II)!

This is not, hvr, "optimal predic": TM is just trying to do trials "similar" to previously successful trials. He doesn't yet understand what "Optimal" means to any ~~particular~~ great extent.

Still, its potentially Much better than simply using ^{sub-tree} documented PPM. I do have to design

a TSG, so it could ~~be~~ usefully work on this problem. Presumably, there would be

"Index symbols" telling TM which is working on a different part of the sub-corpus from 1.

usual (Mainline) TSG. : i.e. Meta-predic. problem is an "inference" problem.

5: summarizing book

ON "Lincos" (1957, done a few yrs ago; born in 1905. (Harry Freudenthal) This is a book on Language's Logic; it attempts

to devise a language (w/o pictures) that could be sent to distant extraterrestrial life - - - that

they could learn to understand. He expects to create forms would get models of humans

in our culture from the TSG he has devised. I would like to try his TSG

on my Lang. Alg. He starts out by teaching Arithmetic, algebra, Logic!

He uses unary notation for ops., then binary - expects e. e.t.'s could punch a 1 or 0 to transfer between the two!

0: 59.39: While it's logically true, I don't know it's possible to arrange the addresses of various Nodes, so

that one can pick a node & its children & parents are always spaced to some distance from it.

i.e. if node address is N, then children are always at N+a, N+b; parents at N+c

a, b, c are same for all nodes.

} source impractical.

General Conclusion That TSG docy. doesn't seem very relevant to induction on functions of known structure. Extending PPM by substituting contexts is impractical, hvr, it is ~~done~~ by CC - unless I can find better ways. Try to find my original work on "How to best do PPM" & how to extend ~~Process~~ ~~on~~ ~~two~~ methods to substituting contexts.

(10) Is impractical: It means that I can start improving the meta-induction w.o. all of Process II.

But I do Need a Good TSG!

STM

Matt Mahoney: PAQ6 V2 - File Compressor: ^{the .paq file is readable by clicking on it or by "fontview"}

Some several files in D:\PS1 4.16.05: Also some Bookmarks PAQ6

This includes the .htm file which lists the contents:

There is a large text file in which I've written comments; but the format is

changed so the diagrams are off: The diagrams in HTML are OK, I think.

Anyway, he uses many pretty programs in it & uses utf-8. The utf-8 was adapted & adjusted to mime bc at output. He was a "successor" of 250 people says "success" note us. response time

One very effective method for using: Uses dictionary of words: Dictionary gives pc to each word (I guess) This method gets ^{about 1.66} ~~1.66~~ bits/symbol for ^{English} corpus.

Using mix of Dictionary & several other ~~other~~ models, Top gets 1.522 for Calc Corpus.

The mixing of prodn. models is certainly good for the Calc Corpus: Also if TM has no idea what external context is (unusual situation, I think)

He has several, apparently "A.H." "tweaking" things that he does in computing pc's. Motivated by his paper on ANN for Compression. His first thesis proposal (which he didn't ~~finish~~ finish) was in 1998. I saw a D\PS1 (4.16.05) for a proposal. ("Cost of Natural Lang. Modeling"): by "cost" he meant both price & SS.

He expected to do it in 4 stages:

I think he bit off more than he could chew for a PhD thesis! The part of which he proposed would probably be fine for PhD!

- Character (letters)
- Lexicon (words)
- Semantic (word assocns)
- Syntactic (sentences etc)

See Ref. in (1998) for detailed description.

How, indep. of Mahoney, progress in each of these 4 directions, continues & to some extent, certain of them are being integrated.

q.v. "Compression - Thesis Proposal - Mahoney" D\B\4.16.05: 44 pp

He mentions "Trigram topic" 11MB corpus → 1.211 bits/symbol

I will want to familiarize myself w. the ideas he mentions

They are unrelated to word General Compressor/predc. He mentions various areas of real interest. He didn't

See: Hells: 1996 v. 13 (4/16/05) Sequenced Natural (not compressed)

Some areas of work on PPM augmentation by the Girl (PhD thesis) - use of English words & also Grammar Markers.

SN Mahoney's method of ^{adaptively} ~~adapting~~ ~~adjusting~~ wts. to N grams, seems in various places, to be what I had in mind, in response to the King-Koon result on "wts" being better than "total pc" corpus for assessed by PEM.

Actually, there has been much ^{research} ~~work~~ on this problem David, (1986?), and Russian expert RHUC.

Matt does simple ~~recursive~~ recursive ~~updates~~ updates - with a certain "decay" = 250 cases. The exact details are probably very impl. - results a (most indep. of context parameters) he recenters. (I may want dist. ~~to~~ λ 's for each PEM) ^{in retrospect,} ^{At not it is an overestimate of how long to recent sub-corpus is}

As another point: "SSE" he modified proby steps, like curves used for

... I think I'm not sure ... it's using bits as symbols ... I think Mahoney ... 250 cases ... 250 bytes ... 250 cases ... 250 bytes ...

20: : "CODEC's in telephony" \leftarrow ^{non-linear} Mod of (d.c.) ~~Volts~~ Volts \rightarrow Volts \rightarrow Volts \rightarrow Volts \rightarrow Volts
 \uparrow (lossy) $\frac{S}{N}$

05: Anyway: To what extent could I use Mahoney's ppm almost directly for Phases I
Mata long? He does give to ppm itself Tho I'm not sure he covers all of the
parallel production pieces like Dictionary PPM.

05: .03-.05 is fairly by cc (depending, hms, on how big corpus is & how many of which ~~all~~ producers/writers)
Anyway, we may be able to get ~~some~~ results (for a corpus of 25 Lisp functions)
using ~~some~~ PPM w. sub-text contexts

10: It may be nice to make some of the PPM's in Mah's system "Sequential" - but
probably this would probably be easy for these pieces.

0: [SN] while I am interested in "General" methods as a ~~goal~~ goal for TM;
The QA prodn (of Phases I) is imp. & does do many really imp. tasks. Hvr. this is intolony
since "imp." is def. to
be compression = hyperfc.
In General, Prod'n is a special kind of Qz problem. There are several ~~tech~~ techniques
Prod'n particularly good for it (like CM). But in general, the "tech for reggs" may
Cover all Prod'n prodn. is! Hvr. a "regg" is a rather "general" type of "object": "finding subfunctions"
is a comp part of it, but I don't think, the entire set of works prod.
Probably the only way to get all reggs is exhaustive Lsrch. - Tho probably not and each corpus
~~the~~ best way for any specific CB. - e.g. for each CB there will be a specific method
Best is Better than Lsrch for Best CB. - e.g. for a long corpus, a short CB.
Lsrch is not best, BV such occasionally is better (unless the P.S. finding Lsrch is V.S.!)

0: One common way is "Coding & Recoding": One can get to a local Max, hvr., so Backtracking is
necy. Hvr, this is Diff from usual ^{with climbing} Backtracking, in that one goes back to the "Best code"
of an earlier corpus (\equiv earlier pt. in the TSCQ).

p 1606

001 : Re: PPM prdn! Would it A compression, if we were able to ~~produce~~ ^{define} "words"?
 There would be "Multiple parsing" since if words would not form a "prefix set".
 Lemple - Ziv uses ~~very~~ "very long" words!
 We could use words as predicates only w/o use them as contexts, (predictors).

207.09

- 0 **SAS** A poss. way to organize ~~the~~ Lecture Series: D could string it in terms of "subgoals"
- 1) Top subgoal: ~~the~~ Seq. problems
 - 2) Next: Sequences prodn. / QA prodn. Physical (Only ~~example~~ string \rightarrow PC prodn transfer for Phase II)
 - 3) Next: Sequences of ~~words~~ string prodn. for Phase I meta problem.

For each sub-goal, perhaps (for my own use as lect) give vars & soft conds for each sub-goal.

For 3) I need to give p.d. over $\sum PSM_i, Prob_i, CC_k$ ^{vars} $\{ \text{real } CC_k \}$ ^{vars} $\{ \text{real } Prob_i \}$ ^{vars} $\{ \text{real } PSM_i \}$
 (Rind?) 0 to 1 or ± 100 .
 "0 to 1" for PC prodn ... impl.

For 2) \langle string \rightarrow PC of next trials \rangle or next word or vars phrase or "rest of text/corpus".

0 ~~the~~ The strings will represent phrases: It in a list way will be function words.
 In other Computer langs, tags will be of different forms. I should study other langs like for PC is Mach Lang.

0 On the "2 kinds of PD's" Question: One P.D. is the original PD that is used to evaluate models - to weight phrases.
 The second is the PD used for search. It is the updated P.D. of the "Investigator" or the "Scientist".
 Hvr, it would Search Best for induction and should always use the updated approx!

0035 Update, If we start at time = 0 w. approx Po, & use given corpus C, then Po, C will induce the updated approx. Hvr, in inducing the Context of C, we will use a prop Po to evaluate models. If the Corpus is augmented by D, we can use Po, i to evaluate these models.

Hvr, for corpus C, D we still would use Po to evaluate Models. \rightarrow Spec! 500 65.00 65.10 imperfect!

STM

Soln. of "2 PD's" Question/problem: (.00-.21)

00 : Re: 64:30-40 (on 2 kinds of PD's); In numeric seq. produ. ("regressn") a good way to do it would be to regard time series as a TSO. So after getting 2 very

02 : rough fit for first 10 items, one uses $P_{0,c}$ as a prior for next 10 items and so on. This is the "Sumac"; The summing for continuous params is a continuous PD - we have 2 means & var. (we could say PD we like... No free Gauss, Gauss, Beta d.f.'s are easy to use. - Much known about them, easy to compute moments, the discussn. of .10-.23 clarifies how this works (.00-.22) works.

05 : Also, maybe we can use it to deal w. equivalent codes for different...? \Rightarrow code 66.10 for v.g. partial soln (which is small)

09 : Well this is a diff. problem! Say we use a MTCor approach: for corpus upto t_0 , we obtain (i remember) a set of continuous params that did well "upto t_0 ". \rightarrow for augmentation of corpus, C , we obtain (i remember) a set of continuous params that did well "upto t_0 ". \rightarrow for augmentation of corpus, C , Re: 64:30-40 while using $P_{0,c}$ (64:35) as an actually equivalent P.D. (as well as a search guidance PD) may seem "dishonest", if \Rightarrow yes, we are still using P_0 for i & p of

11 : of the entire corpus

12 : If we feel jittery about .10-.11, we can actually apply $P_{0,c}$ directly to C (as a augmentation) - so actually, when using $P_{0,c}$ to search for code for future corpus (b.w. 6. param c already coded) we use $P_{0,c}$ but our uses P_0 for a search w. w. to aug. corpus.

If we used a "exact" ($c=0$) $P_{0,c}$ then there would be no difference in results, ~~but~~ \Rightarrow code 66.10 for v.g. partial soln (which is small)

21 : \Rightarrow using $P_{0,c}$ to search for ~~code~~ \Rightarrow using P_0 as prior for $(C + augmentation)$

22 : Hur, since we are not using $c=0$ (\emptyset), we don't do that! I.e. we use our approximate $P_{0,c}$ to search for code for augmentation, ~~but~~ but we use P_0 as prior for entire $(C + aug)$.

23 : NOW! T. forget can be used to find solns. for prodns w. models having continuous params. Which is what .00-.05 is about.

24 : .09 \rightarrow Then for continuation of corpus from t_0 , we use same ^{set} trials as before but w. ~~same~~ p.d. obtained by a prior ($t=0$) times $(P_0$ of corpus to t) for that set of params. We fast fit next corpus separately to some set of params, ~~which~~ using $wts = P_0$ of corpus upto P times a prior of set of params. So we are ~~actually~~ using same initial set of params, but w. progressively changing wts. - i.e. these wts should influence MTCor selection of trials.

32 : favors smaller - "Practical" approach! We start w. a very small subcorpus, i.e. we obtain $\hat{\mu}$ and $\hat{\sigma}^2$ which have i param values of i params that do best to fit $\hat{\mu}$ & $\hat{\sigma}^2$ variance matrix. \rightarrow Then when a new data process is added (which may be over may be 10 new data process), we use $\hat{\mu}$ on that data & we do trials near $\hat{\mu}$ - and so on we can still know much $\hat{\mu}$ & $\hat{\sigma}^2$ change due to new data. When new chunk of data arrives, our $\hat{\mu}$ ($\hat{\sigma}^2$) T. mean bad part about this technique is by computer way, we would need $\sim \frac{(n+1)^2}{2}$ ($n = no. of params$) values to get i components of $\hat{\sigma}^2$. It may be possible in some cases to do ok. w. direct elements only

57M

Getting rid of (degeneracy) or Models: 10.
Great Break Plan!

SM 24
35

00: So only n evals. We may need 2 ~~pts~~ extra pts for each eval, to get second deriv
so $n \approx \frac{n^2}{2} \times 2 = n^2$.

We could use this method to start out on many different ^{discrete} ~~rand.~~ models
Each model has diff set of continuous params. ~~At first~~ we can filter out
bad models by ~~the~~ using very small corpus & using a bunch of 'random' values
Params for each. The promising discrete models then are evaluated more carefully
by using n.l. optzn methods. We then do GS, 32 to continuous evals of each
of the discrete models. We return to best ^(best) few discrete models, - unless some
get way behind i. "best" model.

0 Att! Normally, several of the models will be really identical. (i.e. $xy = yx$
& other x's may change one model into another. In general, to start out
If \geq models have same no. of params, we will suspect they may be identical - - -
If \geq say rate for each no. of dims we will have a set of models - we will use \geq no
set of param values for these models (to start out). If ~~several~~ ^{several} models have same
outputs, we assume they are identical and we drop all but one & multiply by
a copy of that model by the no. of same models.

17 This (cheaply) gets rid of much redundancy.

For more confidence on the identity of \geq Cands: either use 64 bit random nos or
 2×32 bit or 4×16 bits: 4×16 bits is better because of coverage & larger range of
inputs space. 8×8 bit ~~trials~~ ^{random} trials would be even better ... but would probably take longer.
(Pro in a "properly" designed system, maybe 8×8 bit trials should take about
same time as one 64 bit trial.

20 Trying to find real m, real out functions,
for any other combination, of strings & reals, we can sometimes get some discrimination
of identical functions by trying random inputs - but I suspect that would
never to do say $8, 8$ bit ~~random~~ ^{random} strings or $4, 16$ bit random strings, etc.
One trouble is that the range of outputs may however have only a few values
(In extreme case: 0/1) - In this case, random inputs to \geq functions
would ~~be~~ ^{be} same outputs, would not be so certain an indication of identity.

31 To what extent does GS, 24 ~~to~~ ^{to} 66.17 solve the SM problem? It is a real \rightarrow real function.
The problem of redundancy of stocks can also be done by "tracking seq" - so as f. data
string grows, we are able to access more & more and stocks to use for "drivers".

34 \rightarrow A possi. "soln" to the "cost of calling a driver" problem: Say we have 5K arbitrary securities
& other indicators ^{possibly} useable ~~for~~ for drivers. Given \geq stocks to be predicted: list all
~~the~~ 5K indicators in order of abs value of cross correla-w. So. Then process as crossing
the one w. highest correla will be "normalized" k . (We may not have to normalize, since

STM

SN.01

see

0: 66.40 : all possible use drivers — so we only need to know relative pr's of drivers: so it may be ok.
(Pro it is a "mypro for prior")

Another possible pc of reference will be α^k , where α is adjusted so that
t (normalized) expected value of t. choice is some "1/2 powerpoint" of the distribution
of drivers.

On a simplest level, the "correlation" of S_0 w. indicator Z_i would be max value of
(sign of ΔZ_i yesterday * values of S_0 today). Other "correlns" could be on basis of how well
 Z_i predicts S_0 , given 1, or 2 or 3 ... etc. ~~of~~ previous days of Z_i data.

One useful cross correln: say Z_t is Z_0 's value on day t. $S^{(t)}$ is S 's value on day t+1:

Then a useful cross correln is $\frac{\sum (Z_t \cdot S^{(t)})}{\sqrt{\sum Z_t^2 \cdot \sum S^{(t)2}}}$ ← averaged over T values.

If this value is large, it means that a large Z_t value is likely to be followed by a large $S^{(t)}$ value.

T. System could be spending enormous sums of cc looking for good driver/drivers
& associated predictors. T. Q is, what is the "hard problem" it would be working on? — maybe finding good driver/driver
pairs? → 1.6
Unclear as to how GS, 2A → 66.17 (which is on how to find good
predictive functions) would help much — I usually think of cross correln as ^{using} very simple
models. Perhaps we could use more complex models.

16
Hvr. finding obscure driver/driver pairs ~~is~~ seems like not such a good idea, because
their low pc means they are expensive to use. Their 10th best pair has a factor of 10
penalty in "earnings"!

0 **SN** 2 "Exp": One way to Explain STRONG A.I.! Phase I is a ~~simple~~ Universal "Booster"
for any ~~other~~ sequential problem for that is used for its TM₂. Hvr, this is not nearly Ph II.
It's a system that eventually gets a v.p. TM₂ (= TM₁) for Phase I. It still doesn't understand
opten very well — but it is a v.p. "predictor" for Phase II. (predictor means/proby (G | PSM₁, Problem₁, CC₁))

0 **SN** 4/23/05 In my early analysis of class I assumed a "Franner Grammar" ~~was used~~
for a group of grammars to be used. If I actually were to use a grammar-grammar
to generate examples, a TM ~~using~~ G₁ → G₂, could it be used to "control"
Franners? (Note: T. G₁ → G₂ is usually very simple & a "Promiscuous" Grammar.)

SN2 What would happen if I used Rudy Rucker's "Bopper" pgm but randomly changed
some bits in the original pgm? Not very interesting unless I could get the pgm to actually
solve problems!

SN3 Perhaps Prediction of R. Rucker's "Bopper" ~~was~~ would be a useful TM turning problem?
A TQR could be derived from progressively more complex "Boppers"

5+M

On PHI ~~TM's~~ TM's problem is to get good pd for cauds. ~~Its~~ ^{problem} Its comput cook by ~~TM's~~ TM's $\{S, \text{comp}, \text{sol}_2\}$ pairs; possibly along w. ϵ . best/synthesized for best soln. As is to solns are simply bc's $<$ threshold given by Trainer. SN \odot we can regard ~~to~~ TSO given by "Trainer" as \approx "Training Set of data" / used in "Older Problem" (?)

On searching for substrings: \approx Worst case analysis.

Problem that has been analyzed much; Given 2 graphs to find common substrings

In our ~~other~~ ^{known} problem we are given a tree w. known root, (Euler's tour partial order tree we can possibly construct). We want perfect match between n trees (known) tree, n all previous "successful" trials.

Say ~~the~~ ^{known} the previous "successful" term is h symbols in them. The root of tree tree to match has \approx symbol w. $\sqrt{pc} = p_0$. So we have $p_0 \cdot h$ poss. matches. It will take $\sim \ln n$ to find each item $k = \text{size of alphabet}$.

How many comparisons for each of those $p_0 \cdot h$ cases? Each of those searches is a small tree search. \therefore no. of comparisons $\approx \epsilon \cdot \text{largest matching subtree}$

The number of substrings of t , known tree (1.12) that have 1 or more matches in previous successful trials can be very large. Can this no. be $> n$? ≈ 2.8

Of the $p_0 \cdot h$ substrings of interest, most will be very small. Each of them will have a "size number" telling how many nodes they have. If all trees only have branches of size ≤ 2 then we will have $\leq 2 \cdot \text{size}$ comparisons for trees. (trees is comm to "known tree" is ≥ 2 least one "successful trial". So "size" will be \approx majority will give no. of comparisons.

So total comparisons will be $\ll n \approx > p_0 \cdot h$.

If the known tree has L leaves (terminal nodes), then we have at least 2^L poss. sub nodes that match known tree at depth of known tree or depth - 1.

Further what if t whole trees of depth D we will have $\leq D \cdot 2^L$ subtrees. Since each "depth level" has \geq least 1 less "terminal", I think we end up w. $\leq 2 \cdot 2^L = 2^{L+1}$ subtrees.

How, I'm not sure of 2^{L+1} result, it may be much larger. Even $D \cdot 2^L$ maybe ^(very low) completely wrong. 2^L is a lower bound on no. of subtrees.

Def RT If there are N nodes in the "KT" (\approx known tree) then certainly 2^N is an upper bound to no. of subtrees. But between 2^N & 2^N is a large gap!!

Instead of worrying on .12 ff. Consider a \approx mouse problem! How many trees are there w. just k nodes? I may have solved a more general problem in heath (1989-91) i.e. how many functions are there w. k nodes? I'm not sure just what I said. I got a functional (recursion) eq. and then obtained a closed form soln. by guessing. ^{That} ~~That~~ Soln. may be closer to what I want/used than .12 ff!

5TM

20

for 1-time being,
 Dropping $\frac{Q}{k}$ off no. of poss. contexts: The practical Q is to use P size subtraces for produ. A fairly General context, consists of short spans that define sets of subtraces. Associated w. each such context is its $\text{aprop} (\approx 2^{-k \text{ length}})$ mult by $\frac{1}{k}$ pd. then assigns to 1. next symbol. It codes not all of f . corpus because its contexts do not occur in every case. So consider to total extra compression it gives for f . s. e. m. b. o. l. e. s. for which its relevant!
 I ~~need~~ big discuss on P is when I figured out how to use arbitrary context subtraces for produ. — try to find it! **4TM 223.31 - 227.40** is v.f. discuss. of this & related Q 's,

10:

For regular PPM's w. sub-trace contexts, its likely that we w. d. not get v.f. produ. unless Rep Corp are very large. This is essentially $\frac{1}{2}$ ^{Guess Memorization} Table look up (v.f. For much faster" (v.f. much less ≤ 2) we ~~have~~ ^{have} to use "mpm st": These are cheap ^{descriptors} defined sets of sub-traces ^{Subtree} contexts. Hvr, even tho. it would be "expensive" it training & in needed P (i.e. ≤ 100) we may want to do it for PHI "tree (v.f. \leq PHI's TM) — This enables us to "pool" Data from many distinct long Domains. In my early work on sub-trace context (perhaps ~~from~~ near by procedure 4 TM 223.31) I prot. of pools of trees of same size \leftarrow or that size \leftarrow (i.e. cardinality) would be key to aprop ($\approx \frac{1}{k}$). Maybe some things like $\frac{1}{k}$ for size k \leftarrow Note this is an "improper prod" & commonly used for "relative p.c." — But pooling trees of same size would seem to remove near whole point of sub-traces for preserving specific functions \rightarrow .30

20:68.40

SN I think the idea was $\frac{1}{k}$ like PPM: Given k bases and a set of functions of various arities, how many trees of size at least size are there? **NB** Given k functions to identify function (of arity $= 1$).
 This $\frac{1}{k}$ notation will change when I find some 89-9 soln. \leftarrow $N_k = \text{number of functions}$
 $d=0$ $f(x)=b$ $d=1$ $f(x)=a$ $\sum_{i=1}^k x^i$ \rightarrow $\frac{1}{k}$ \rightarrow $\frac{1}{k}$ \rightarrow $\frac{1}{k}$
 For $\frac{1}{k}$, there are k possys. \leftarrow $\frac{1}{k}$ \rightarrow $\frac{1}{k}$ \rightarrow $\frac{1}{k}$
 that has k possys via identity function. This assumes no symmetry (like commutativity) in function inputs.

27

So maybe $F(d+1) = \sum_{i=1}^k (F(d))^i$ So just what problem did I solve of SARB?

28

I think the functions were not trees but just directed graphs: A "concept net".
 Actually, I probably want (in TM) something like this, rather than a "tree". **(NB!)**
 (Note I'm worried on contexts which are not functions.)

30: (19)

So it's likely that I will not "pool" in PPM way: but subtraces of same size may not be the same functions!
 Same aprop. — So, do ft deal w. cardinality of subtraces as a func of no. of nodes.
 $\frac{1}{k}$ may give $\frac{1}{k}$ value of $\frac{1}{k}$ (i.e. cost of coding). — Hvr, Note **4TM 227.33**

T. analysis of .27-.28 will be needed for PC comp. — also I have devised a "context" data str.

So I can get p.c.'s of them.

While f. w. of unrelated functions functions contexts is quite large, I can assign arbitrary order no. to each of them and store their case counts in a hash code.

STM

30

There are lots of courses/differs w. fuzzy stuff! ~~the~~ hvr, I think analysis of 4TM 223.32-227.40 (i.e. for $\alpha, \bar{\alpha}$ coding) is correct: but f. pc of derby f context is v. hard diff. T: contexts I'm using do not have to generalize of f. functions I analysed in SAARB. Cifs. 69.28 R.

04

05

Second, f. contexts run backwards: They are trees but different from f. forward full tree

See 59.03 ff [Note 59.39 on "Content Addressable Memory"]

36

A method of analysis that I don't remember course: Analysis of LSP expressions themselves. — "All strings of LSP symbols are logical functions" (is true).

06

u. be because of in country no. of functions.

0

If I do get into some differences in Estimating costs of definitions of context: ~~is~~ possibly use Monte Carlo runs to fit curve to data — as I did on SAARB.

Re: f. diffy of .04-.05; since f. contexts are trees: perhaps I rendered them, by ~~the~~ Polish notation!

In which case it may be easy to list/count them. Looking at 59.03 ff: It looks like f.

"backwards" desc of "contexts" is same as for word notation; Backward Nodes have some "arity" as forward nodes — so desc on desc. them in Polish.

14

An exact way to get down pc of any tree (forward tree or backward (CONTEXTS) tree) is to multiply f. pc's of each char used in generating tree. Tho v costs at each pte will be integers — so: product of integers. There are some errors in picking one single tree. — since many trees are identical, thru symmetries, commutativity, assoc... etc.

20

Re: Implementing 2 eq. of (4TM 225.14)! Various approxs are poss.

Usually only largest, or a few large terms are needed. Also note (4TM 225.18) || f starting (simple) approx can be used. Even for $n=1$ it's not too bad!

Not long approx to $N!$ Better than starting for small N works for $N=0, N=1, N=2$

1! $\frac{1!}{e^1} \sqrt{2\pi \cdot 1} = \frac{2.5066}{2.71828} = .922$ v. close than! $\approx 7.8\%$ low hvr. g. was much closer to 1.0

2! $\frac{2^2}{e^2} \sqrt{4\pi} = .91$ 9% low.

Usually, n will be $\gg 1$ & in general, errors in ests. of w 8% would usually be acceptable.

3! $\frac{3^3}{e^3} \sqrt{6\pi} = 5.836 : \frac{5.836}{6} = .972 : 2.7\%$ low.

30

My guess is that f. exact analysis of 4TM 223.37-227.40: would be better if it was used to define WMS v. other Plan contexts.

I think better profa/zunderside. If you don't define terms, you don't get many multiple purposes — much simpler... but is it as accurate?

37

T. unmethod of f. analysis is eq. v. to a more exact formulation of the "escape costs" i.e. f. pc for change in going to a diff. "context length" or "function tree size".
It suggests that to be complex problem of weighting "func trees" in it to be done empirically by an adaptive weighting method of "escape cost function".

Maybe wrong!
See 80.00

5 TM



The analysis of f. dom. cost of any tree (i.e. context Huffman) as in 20.14, is very neat!

While it would seem best ^{M. sp. of} ATM 225.14 would be a "exact" way to do prodn., it superficially seems at odds w. the usual PPM "escape" pc's. My impression: They use the longest matches "reference" from the previous production pc, they shorten down context ~~to~~ to first-occurrence valid prodn. Each bit of context length is wt. by exp factor of α . $\alpha \approx \frac{1}{n}$ where n is corpus length, up to that time.

AZ 141 \downarrow pc for longer contexts, \uparrow pc by freq. of occurrence of match. It does not seem to be getting a monotonic \downarrow in wt as context length \downarrow . (!!) - This seems like a big difference betw. PPM systems & AZ (41)

SN How to get hyper-precision pc's from compressed systems w. 1 Byte resolu:

To get pc of next symbol; consider a (k-symbol) continuation of each pass. "next symbol". I don't know how much precision one would get for $k=1$. As is to get pc for 1 symbol, one has to do 256 methods out symbol trees: - to my best first most likely 10 would be an of - so maybe 100 or (much less) 2 symbol continu. Look up trees, grams & grams. Also, it probably would be a lot easier to look at the context & get pc's directly.

could it hash code trees?

Re: Subtrees: All I really need is good way to store/index sub-trees. I can use them directly (as in PPM) or as "definitions". Definitions have advantage of speed, compactness... but disadvantage of definite "commitment" - that can only be "uncommitted" by "reversing" - which is usually done for only a few recent "levels" of definition.

30 **SIV** 4-28.05 Hash coding of trees? Each node in tree has a unique path to it. When a new tree occurs, we want all path trees w. same paths to any of its nodes. The paths can be hash coded. Hvr, this means that every node has its own hash code. Well, put these paths in Lex. order? ^{No need for Hashing} We can be satisfied in a "content" parts of each path. Each node will have an address of its entire tree. I'll have to work out details of this: It seems to be interesting, but it may not be even marginally practical from an ec. pt. of view.

20

30

33

STM

o: ... T. trials, there are an enormous no. of "paths". Consider a corpus of O_j functions!
 This is a set of prodn algms (cond. pds) that have been successful int. past ...
 int. sense of "Accepted by Trainer".

Each symbol in t./corpus has an associated tree, w. leaves are "set of primitives" (see
 discussn of "AZ" pd in one of appendices of I D S I A report.)

Some of t. trees are short (of depth 1 or 2) others are full (depth of t. function depth).
 If M is t. max no. of nodes in a function, and there are N functions to NM nodes each.
 Each node will have $w \log_2 M$ possible nodes on its path.

Q.K. - Say we have A^d many trees: depth d $2^{d+1} - 1$ nodes: $0 \rightarrow 1$
 $1 \rightarrow 3$
 $2 \rightarrow 7$
 $3 \rightarrow 15$
 $4 \rightarrow 31$
 $5 \rightarrow 63$
 $6 \rightarrow 127$
 $7 \rightarrow 255$
 $8 \rightarrow 511$
 $9 \rightarrow 1023$
 $10 \rightarrow 2047$
 $11 \rightarrow 4095$
 $12 \rightarrow 8191$
 $13 \rightarrow 16383$
 $14 \rightarrow 32767$
 $15 \rightarrow 65535$
 $16 \rightarrow 131071$
 $17 \rightarrow 262143$
 $18 \rightarrow 524287$
 $19 \rightarrow 1048575$
 $20 \rightarrow 2097151$
 $21 \rightarrow 4194303$
 $22 \rightarrow 8388607$
 $23 \rightarrow 16777215$
 $24 \rightarrow 33554431$
 $25 \rightarrow 67108863$
 $26 \rightarrow 134217727$
 $27 \rightarrow 268435455$
 $28 \rightarrow 536870911$
 $29 \rightarrow 1073741823$
 $30 \rightarrow 2147483647$

So sum $\sum_{d=0}^n (2^{d+1} - 1)$ is \approx No of nodes in net mult by depth of root.
 $d \cdot 2^{d+1} - (2^{d+2} - 1) = (d-1)2^{d+1} + 1$ nodes. in all of t paths.

So $\approx (d-1)$ no. of nodes in net + 1 \rightarrow So \rightarrow t. storage needed only $d \times$ no. of nodes in corpus

Say it is true: We could do prediction using "all path" contexts for each path symbol to be predicted.
 For z / error of depth d - we would have a 2^d associated path "paths". Each path would have an
 associated d from the alphabet as an associated SSZ . From this info, using ϵ ideas of

18 ATM 223.3 (1-227.40) we can get a p.d. for f. next symbol. ~~that~~ - all off. "path contexts"
 working in ll to give t. final p.p. How does this compare w. my using "context substraces"
 for predn? Are these methods about f. same? If not, just how do they differ & is
 one much better than the other?

It would seem that "Context path" method would be easy to implement ~~by~~ by storing
 path contexts in Lex order.

HVR: Using these "paths" for predn is not f. same as using substraces for predn. To get substraces,
 I need to AND of a bunch of ~~paths~~ paths. What I've been considering is a "OR" of t paths.
 Say I had a list of all off. maximal paths ~~for a node~~ for a node, not had cases in path. Certain subsets of
 these paths will have a non-empty AND i.e. ^{common} substraces as cases. What I'd like is ^{all} common substraces
 that have ^{max} SSZ as a narrow predn & much wt. (\approx short decs).

1.29.05 SN For ME-Colo Lsrch! Generate trials randomly, then ~~then~~ use hashcode of complete
 d.crn. of trial - to see if there is previous work on that trial. If yes, continue next work w.
 \approx construct $d \cdot c$, do a st, make new hash sites to do $d \cdot c$ on that trial. Hvr. one will have many duplicate

SN T. trials in AZ are not "your standard tree" / Functions can be generated w. ~~any~~ of
 any arity, & used in subsequent parts of t. causal. d.crn.

Can Rice "Common path" trick of 71.33ff be used to ~~do~~ do a PPK on AZ?

[AZ ~~is~~ as in t. an appendix of I D S I A report ~~is~~ ~~to~~ ~~not~~ ~~do~~ ~~reconstruction~~ - but it's
 See do Ling (2 form), so it can do "If ~~for~~ $x=1, f(x)=3$ (or $f(x)=G(f(x-1))$ "]

functions w. default
 d.crn! One can
 suspect this is t
 Goes over idea how
 to solve problems?

5 TM

20: [5N] Differential Eqns are a way to define ~~some~~ recursive functions (Prim. rec. functs).

Since they are for continuous functs, ~~the~~ they would seem to give them a different domain of "defined" than usual discrete defns of recursive functs.

E.g. is there a differential eq. that generates $f(x)$ funct? - look in Bur. Steady - Book. (Hypergeometric D.E. ?); ^{idea related to Shannon/Moore} IST's continuous func. work on defining Universality in Continuous Mechanics?

[5N] In ~~my~~ my fault of PPM to function trees ($\Sigma \geq AZ$), I've always been using "direct" nodes ~~and~~ and "nodes accessible by continuous PPM" as "context".

In general, if we allow Neurons as contexts, this will not be ("effectively") so.

E.g. In vector representation of Σ (columns) dim. images. - ~~the~~

One "heavy context" is on nodes that are ^{directly} just k nodes behind - where k is

the ~~long~~ length of a vector line.

[5N] To prove that any CPM (Computable Problem) can be realized by Σ (3 output bits)

~~w. random input!~~ For each output bit, use ~~the~~ Method of AVG Coding.

Peter Elias \rightarrow David Williams \rightarrow J. Reissner.

20: Main Big Problems for TM:

1) Σ writes a sort of popular explanation ~~the~~ present ~~the~~ paper: Partly as introdu. for discussion w. Marvin. Also mention: ^{search/TSA} 2 kinds of ~~input~~: Lsuch v.s. BU.

2) Work out ~~exact~~ / approximate PPM / AZ methods for sequential (or pos): Compare w. standard / methods. See if big diff. distance.

3) After 2) work out how to do (w. subtree contexts; ~~possibly~~ subtree sets). Both exact math ~~is~~ / ^{deriv. w/} Approx. are imp.

4) Designing LZ-type (Σ LZP) for long. in PPI.

5) Investigate ~~for~~ Machine lang. w.r.t. learnability of regys.

6) Design PSM's ^{for phase I} (Σ maps, Metrics ... etc. to start) - Eventually try to get to Learn. Hierarchy of various Early AT. plans.

7) From good set of PSM's: Factor them down into Σ Loopless GFG ... later, perhaps (copied) CFS.

8) Design of TSA for Phase II.

9) Lsuch ~~is~~ ^{on} contin. param. functions: Every have this under control - but can not see.

10) Lsuch for ~~stun~~ ^{on} string output. (Idea of "counter plus insts in 'danzhus'" good, but needs work. Idea also need good functional forms.)

STM

20

11) For BU such: Use of S-CFG: Then \Rightarrow S-CFG is other enhancement
 BU such is \approx improved GP. $M_x - M$ may be good - for "measured base" - is G. Goves I normally use (for my (local) induction Prob solving) are \approx linearized ... e.g. Macey, PC, time, CC.

12) Long "English".

13) T. Problem of when to Make Definitions: How to reparse using Rom; How to "underline" in Backtracking: When & how much to Backtrack (\equiv recoding).

Ritehand ■ 73.21 & 23 are must-read: 21 is also part of "Lecture Series"

So outline of 73.21: "Road Map to Strong A.I."

Early work on Heuristics and Expert systems: Inspired by M.S. Logan's Recursion & returns of the GPS, etc.

In 1956 ~~my~~ work had ideas for AI, but no actual program (M.C. (why some) chosen by working on LISP (a Meta Project for A.I.),
 However, my main interest has been ML (I wrote first paper in 1956, then 1960, 1968, 1975, 1983, 2002, 1975: table).

20

SN In case of sym. args of a function - for args, say (commutativity), consider function to have 1 arg: t- unordered pair x,y. Can this help in detecting, utilizing Symmetry?

SN In applying CFGs to ~~any~~ functions in any unknown lang (like APL, LISP, A2) we already know f. lang, & it's a unique parsing lang! "Contexts" are normally trivial & involve a depth of 1. However, we can consider deeper stacks for contexts, which brings us to our present exam. of PAM. So it looks like CFGs would add ~~to~~ nothing to f. PAM approach. (or vice versa @).

SN is "Classic PAM" approach to lang. We find 6. context adj. cont. sub trees to present prediction that have occurred in past at least once. We use (in II) all of them that give > 0 pc to next (known) symbol (for compression). If none of them work by all sub trees that have one fewer symbol, & use "escape proby" to approximate pc wts. If this yields no correctly predicting sub-tree contexts, form smaller contexts by deleting one symbol (in all ways) from previous "deleted" subtree giving another escape penalty ... continue until a context is found that produces x. next (known) symbol.

The lang. can be used to get pc's for all of y. poss. "next symbols". - Maybe not perfect" but it is "inspired" classical PAM. ~~It is~~ it probably has a much lower CC than the correct ~~APL~~ way, but still much less CC than full APL (\equiv AZ14).

5 4 05

5 TM

.09-.23 & .29-.31 Somewhat of a "Breakdown" in Prodn w. string accept
Applying to both Ph I & Ph II. It really "Boostraps" Ph I!

SN Rémi Béraud: Complete Algebra!

PSX LISP-Mathematica 4.8.04:) I made word copy that I put in Wolfram's "Mathematica" book.

This Guy is called a "Matrix" as a fusion of Function, Symbol, & Algebra languages —

is able to express approxs. in a nice, useful way. 2 page paper.

It discusses 2 very General Approx methods that I don't know about:

"Perturbation Expressions" & "Variational Formulations"

It is clear that there are several parts of the Math that I should learn — so ~~not~~ to be able to do better approxs!

SN An early development in Ph II: Int. Grammar of PSM's: To decide on next symbol
found in a derivation, TM considers the "recent state" of its systems & Q is the next symbol

in A is used, this as a corpus for a Ph I production.

Unfortunately, the QA problem of .09-.11 is string \Rightarrow string prod. — I know that I'm least able

to deal w. (Prodn. has been recent progress). A suggested soln. to .09-.11 Part's "Not Bad"

Try various ~~Q's~~ \rightarrow O_j (deterministic): Each will have a production score for post

how many fractional fitnesses beavered we then take some kind of wild guess of x. prodns

of these O_j 's. One needs a lot of O_j 's to cover all possible A's (say x.A

alphabet is ~~256~~ symbols... (or probly we could do w. many fewer!))

Another way to do prodn is by Njists. ~~function~~ F^i would recognize

a certain set of contexts. Each set has a production vector for next symbol ("A")

we do wild guess of (pc of F^i mult by pc of its corpus) \rightarrow we can use the π , π method used in AZ 141 to get ~~relative~~ relative pc of corpus wts for various F^i 's.

These F^i 's are like "definitions".

Of course .09-.23 doesn't solve Ph II's problem — even at a policy level —

It does provide for V.G. ~~some~~ sub-langs for PSM's... an early part of Ph II's work.

For the rest of Ph II's work involves doing (Probs, PSM, $CC \rightarrow G$) productions;

a string \rightarrow Num prodn, which we know how to do, to a large extent.

HVRI

.09-.23 can be used in Ph I's meta problem! We normally use a

dogmatic PPM to get results for Ph I. We can use .09-.23 to get much better, much more general kinds of prodns for Ph I.

ABIG Problem w. F^i : We accumulate big large set of F^i 's. Trying them all out.

first corpus would be (eventually when we have Very Many of them), very CC expensive.

We can either arrange a "hyper order problem" to decide which ones to try, or find some Standard way of ordering F^i 's

One good thing about F^i 's: They can be frozen II.

\rightarrow See 27.04 for one possible way to use F^i 's.

00 : IN PH II, I'd like +M to devise a D.F. over PSM's for a particular problem,
So Rora would be an emphasis on PSD's that would "do work" on that problem.
As is, I got a D.F. on PSM's, w/ a s-funct that relates ~~to~~ by (PSM's to, prob fns) to a
~~df of conc v. s. cc.~~ — which seems very time consuming

04 **EN** For BU such we may want to do strict Monte Carlo trials: In Rig case, in 76.32-40,
we can choose the F^i functions in pc order until we get one that sees "yes" — then generate the trial
according to the pc's assoc. w. that F^i funct that was chosen.

07 Actually 76.29-40 isn't so easy from another pt. of view. Say F^i is a possibly useful
generator of Ngrams: To validate F^i for predn, we have to try to find cases
where it
in past where it was relevant ... Not so easy (cf. big cc!). Trying all possys in
past is too expensive. We may try Sampling - Monte Carlo-wise.

13 One way to solve cc.: say α following symbols α . We examine all
cases α which α occurred & we look at Set of (previous contacts). Can we find some characteristic
within that set? ~~See~~ we can try various F^i 's on that set
Using Lex Ordering (as in PPM) we can easily find sets like. We can also possibly find
probable contacts of 2 or 3 or more symbols using Lex ordering trick.

00 Hm. Re-13ff: It would also be nice to have negative cases (In addition to the corpus of positive cases).
Say we simply try to compress the corpus of positive instances — i.e. find ways in it — This is induction
on (ordered) set of strings: One way to do this is to consider it as a "language" needing a "grammar".
We can try CFG, CSG etc. This seems to be drifting towards the rather flat corpus ~~is~~
is a language that we know the d-grammar of. ← But what about the idea that
new approach P13 would give us the power of "hearts"?

30 → Actually, we can apply PH I level to PH I level directly! Say Level zero is t .
So $(QA)_1, (QA)_2, \dots, (QA)_n$: Level 1 is $O_1^j, O_2^j, O_3^j, \dots, O_n^j$.
We give the system various kinds of QA problems — including those of ~~ex~~ trapology

34 Sequences of strings or repression of strings. (i.e. either $Q = \text{string}_i \{ \dots \} \text{string}_k$; $A = \text{string}_{k+1}$
or $Q = \text{string} \& \text{symbols}$, z_1, z_2, \dots, z_k , $A = z_k$ from such capabilities, we can
then apply P13 P14 O_1^j, \dots, O_n^j seq has a Q as or we can represent 78.00 Spec

00

77.34 ~~████████~~ $z_1 \dots z_n$ as f. beginning of O_{n+1}^j ; if we want z_{n+1} , the D.f. on
 the next symbol.
 Of course, TM must know ~~what~~ an appropriate TSD (leading to such a problem —
 or it will do poorly!
 Also, even if 77.30 ff were done perfectly, we would still have only a u.g. PHI, but
 didn't "understand" what "optimization" means. It would simply find ^{very} good "2 word"
 for O_{n+1}^j , based on past O^j 's: T. D.f. wouldn't understand/know what the goal
 of it ~~is~~ all was! — i.e. good product of QATSD.

0

(EN) Re: ↑ I'd like TM to be able to divide up its corpus into sub-corpi (by indexing the
 sub corpi, say ~~by~~ — which frames could teach it). Say it did this (or had ~~it~~ done for it).
 A q. is just how to ^{and how much} have info for sub corpi? This seems to be an Inpt. problem.
 One way: for each sub corpus, we have a vector, telling how much wt. was given to data
 from ~~one~~ other sub corpi: so we have this $n \times n$ matrix for the sub-corpi. (I don't know
 if matrix is sym. or has other simplifns.)
 Anyway, T. system is able to look at ^{new} problems, assign ~~sub~~ them to
 sub corpi to modify wts in sub corpi.
 Also its able to Merge old sub corpi & create new sub corpi by ~~the~~

20

~~select~~ & merging together a subcorpus based on former members of other subcorpi.
 Re: 77.30 ff: Its not clear just when in development to learn this kind
 of self improvement. One could just start out by periodically testing to see
 if TM can compress ~~of~~ f. O^j sequences. If it can compress by
 some Δ then we can figure out what fraction of time to spend on self improvement.
 See ~~47M124.33~~ ^{stuff leading to} for how to do this. Go over this now: it's imp. idea

0

! **SR** One "not bad" way to do S.I. for PHI: Do search for predictive & substructure sets.
 T. search for single context sub-frames is perhaps simple. T. search for "ngrams" may involve
 "recog. functions" of 76.20. is more diff (>>cc),
 This doesn't seem to be "Booster" (Regenerative) hvr. — its just a way to improve
 the zero level PHI
 Certain kinds of subframes & perhaps ~~subframes~~ subframe sets can be represented as vectors
 like Zero Codes. — And we can use a kind of dot product w. present situation, to get
 wts of relevance.

5TM

in Row 100, including 120
This is on 5 kinds of string predictions

10 Summary of various kinds of predn of strings (usually individual characters),

1) Simple contexts as in PPM.

2) " " but PATS thru a tree

3) Can ^{take} ~~be~~ that are single substrings

4) " " " sets of substrings

5) Fⁱ ("recognition functions" type contexts) they over a variety of τ : A way to decrb a set of contexts.

for 1) ... 4) we use Bern predn: Each of τ contexts decrb a sub corpus: I guess we could use τ of α

Method of coding: But what about the multiplicity of models of D.F. and wts of contexts?

The Fⁱ are ^{also} ~~also~~ τ (sets of characters) — but τ characters need not be ~~set~~ ^{subset}

τ d τ cont ... they can be placed anywhere in τ tree. Also they can include set of substrings ~~which~~ ^{to check} that have never occurred before, see (11.13)

How, for just about all of 1) ... 5) ^{each} poses a problem.

Also, we can view 1) ... 5) in 2 ways: τ as pure context w, and P.P. on τ symbols to follow

AS you've defined objects, Part τ to include τ char to be predicted. As defined objects, the

remark of (3R) must hold, since we have all these parsing methods. So does (3R) hold for contexts that are not "defined objects to be used in coding a corpus"??

6) An older string τ type of predn, looks at τ context ~~from~~ previous part of τ function being generated τ gives (as a function of that "previous part") a d.f. on τ next character.

7) see 92.10-93.10 also 93.24-25 for a combn of 79.20-40

One way to ~~get~~ ^{put} a good(?) universe(?) d.f. on strings is a function of τ part τ . We have a bunch of τ -functions, each with its own wt.

These functions are on τ plus one or more other symbols of τ .

So τ is generated like ~~string~~ τ corpus τ PPM, possibly one symbol at a time.

Each of τ s funct has a d.f. on one or more of τ chars in τ alphabet. Like in PPM, the τ -funct that have a d.f. on very many chars, tend to have less wt.

If I could work τ off out, it would seem to be more or not so exp for me to find

pc's of ~~any~~ conceivable "A" — just τ PPM can assign a pc to any particular

contn of τ corpus.

In thinking about "string output predn": I tend to think of a function of τ outputting a τ string w. by pc.

~~latter~~ aspect, not so much τ former.

A τ -func Related to τ off: A function looks at τ tail end of a corpus (represented by a

locative symbol) τ assigns a pc to "next string" (also maybe a wt). It could assign wt. to partial strings completed w. "dot care"

Any funct that assigns pc's to sub strings τ a "s-language". For strings w. condition .32R, we can conceive of a new kind of "grammar"

for $L \rightarrow R$ on line parsing. If we had a ^{stochastic} grammar for complete ss only, we could define

a grammar to deal w. .32R τ essentially $L \rightarrow R$ parsing. \rightarrow see 92.10-93.10 + 93.24-25 for contn of this

STM

How to do Sub-tree Contexts in the PPM manner: .10-.19 w. trim of .19-.21

10: 70.33: [SN] I have been trying to look at PPM as an approx of AZ(1): But it looks like it isn't!

It's an independent method of coding. ATM 223-227 Discusses AZ(1).

PPM certainly seems much easier (eccc) than AZ(1) - Too complicated. ATM 223 I think

I worked out a way to do \approx AZ(1) - but I'm not sure it was right - it may have been just an approx of PPM!

06: Using escape codes, not 2's place coding costs \rightarrow 10. I have been considering my simple function trees; but the kind of Z generates 13!

NO!! ATM 223-227



I think they are easy to count.

[SN] I'm not sure I included this type in my analysis (at SAARBS)

0: 06: Anyway, though I haven't been able to find it: The way it works, i.e., we find

- 1. largest contexts in past consi. w. present context. Say there are all of size k. If
- 2. case counts of all these contexts include the correct pred, we use total case counts for
- 3. pred. of (next symbol) (A poss. problem w. multiple counting. Maybe not. With contexts of size k, we must be exclusive.

14: [alpha] If the size k contexts don't provide zero pc for (next symbol) do escape to contexts of size k-1 & make pred from their case counts; if pc = 0 for correct pred, loop to [alpha]. (14)

A more exact (perhaps) method starts by looking at all contexts of (max) size (k) that fit.

It then uses pc to specify which of them give correct pred w. pc > 0. If there are r contexts of size k, then there are 2^r distinct poss. subsets of context to specify (rather expensive!!) - one of the 2^r is None of them -> "Escape" to size k-1.

9: Another approach is to use pc = 1/r for each off. If codes that do give pc > 0 for rite pred. so that r_i of the r contexts do predict ~~correctly~~ acceptably; we have r_i codes: each with pc = 1/r.

22: So: we just add to case counts of the r_i contexts (they are mutually exclusive) and give a sim, with 1/r. (The r_i contexts that predict a "not a") also each give legit codes of pc = 1/r. Which might suggest that we have to use all k sized contexts in 11 - & simply add them from case counts!

Is .22 Lay.1? Do we have to use all 11 codes? - My impression is "No": Just the n-n_i codes would

"Not code to carry" - That .19-.20 is correct.

26: For .10-.14 is ok. ~~with~~ if ~~the~~ of .19-.21 is added.

The major problem w. using ~~the~~ sub trees as contexts is the (apparent) difficulty of keeping track of them, storing them, getting access to them when needed.

One possible way to do this is to have say 16 (i.e. 2^4) parallel stores of context: one for each 16 paths backward (for binary nodes). Each one is in Lex order, like PPM

The store is index paths backward to distance 4

Another way is to try to find short tree contexts. How they could be stored & accessed when needed, remains to be discovered!

Another way is to use definitions of F^2 trees (i.e. sets of sub trees). They have advantage of needing less space since they include sub trees that have never occurred. Ideally, we could devise a function that would look at the suffix of a context & set a pred over F^2 that is most likely to be relevant!

STM

00 : A way to get ~~good~~ predictive 'P' function. Collect a ~~large~~ ^{sub} set of context (trees) that precede each poss. symbol, & look for ways within each of these sets - try to compress each of them.
 Maybe that kind of compression is not so good, I'd like ways to use a sets for compression.
 - is a given context (z tree) so find the def. of sets it is likely to be in? A categorized problem.
 Perhaps SVM could be used - for elements that are strings, I don't know if SVM's are usable. If I could find good "Features" of the trees, I might use "decision trees"

06
 Actually in the same spirit of the above is regarding a "pd on the next character" to be a QA problem for PHI: "Q" is the ^{context} (postfix), "A" is the next character; "T" corpus is a set of all (postfix, following symbol) pairs that we have in corpus (if the ^{context} corpus has n symbols, we have just ~~just~~ n-1 pairs. .06 off is a v.g. way to regard PHI's self measurement! (≡ SI) —
 That we want α (the next symbol) to be a pd on the postfix ~~of~~ of the corpus (R is R) — so we try various d functions to get R's, a ~~tree~~ d function will be like a function of f, times p (→ pc).
 Or, we may be looking for ^{position} functions on the post (→ T/F) (≡ Obs) ~~in~~ which the T's (or F's) have a narrow def. on the next symbol, α.
 In general, I want to try to interpret the QA problem in the form of ~~the~~ ~~Q.20~~, a f. problem of generalizing what GPM does — prediction of next symbol/phon.

10 : A really ~~not~~ good way of prod is using Data & periodically Reparsing all of your corpus. — Think more about that!

12 T. example of Clark & Sauer not ~~using~~ ^{Working map} way have seen a State Tree (MDL) followed by Pruning (which I don't understand ^{or?} but ~~page~~ ^{buffer} ~~proceeds~~ ^{Decision Tree}).

My Impression of Pruning Dec Trees: (pp 662-3 in H.I. & Modern approach ~~2003~~ (2nd Edn))
 is that normally (in pruned) Dec Trees, uses the "Best" split available — it does not, however, consider ~~SSZ~~ & pc of f. Model form. In the pruning discussion, ~~the~~ SSZ is considered, but Model ~~pc~~ is given a default value — ^{'is same value as "Null hypothesis"'} — they actually use Hypothesis testing formalism to decide if a split is legit.

The "95% confidence level" (that's "standard") is ² accurate. (I think) to a certain ^{approx} of Model. — Curiously, 20% is better explains. NB. 2.39
 → T. way Hypothesis testing works: it decides the relative ~~prob~~ ^{prob} of the Hypothesis vs. Null Hypothesis. I would (in ALP) express this as relative wts. of H vs. Ho. The 95% conf level means 5 to 95 or 1 to 99 ratio of wts of Ho vs H for acceptance of H. This is actually not v.g. for prodn. (Pro maybe o.k. for "Hypothesis testing" which is ~~bad~~ ~~idea~~ ~~anyway!~~)
 I guess they assume all hypothesis H have same prior: ~~the~~ which simplifies things!

34 I would probably do ~~all~~ splits down to "50% confidence" i.e. when H & Ho have ~~the~~ wt.
 36 Another term for normal Decision trees would be "loss Greed" i.e. Comparison of orderings of splits.

On 'splits' say before split we have n_0^0 0's & n_1^0 1's.
After a split we get n_0^1, n_1^1 & n_0^2, n_1^2

The original pc of data was $\frac{n_0^0! n_1^0!}{(n_0^0 + n_1^0 + 1)!} = P_0$

After a split pc of data is $\frac{n_0^1! n_1^1!}{(n_0^1 + n_1^1 + 1)!} \cdot \frac{n_0^2! n_1^2!}{(n_0^2 + n_1^2 + 1)!} = P_1$

If $P_1 > P_0$, the split does some thing useful.

My impression: if $\frac{n_0^1}{n_1^1} = \frac{n_0^2}{n_1^2} = \frac{n_0^0}{n_1^0}$

Then $P_1 < P_0$

$\frac{n_0^1}{n_1^1} = \frac{n_0^2}{n_1^2}$ always implies $\frac{n_0^0}{n_1^0} = \frac{n_0^1}{n_1^1} = \frac{n_0^2}{n_1^2}$: reflects equality at any 2 of 3 ratios implies 1. & 2. are equal to other 2.

Other way round, if there is any inequality in the ratios, $\exists \epsilon > 0$ for all $n_0^0 > N$,

$P_1 > P_0$

10-11 is probably provable but I'd like to be sure! also know how large N must be

various deviations.

On second thought, I'm not sure 10-11 is true! - it maybe may that deviations in children differ in a "usual" way from that of parents.

Try some experiments w. now before trying a proof!

start w. $n_0^0 = 5, n_1^0 = 5$: show do various n_0^1, n_1^1 pairs (when equality n_0^2, n_1^2 values)

4. Using only size of the precursors = 0.5. which is split into 1 for 0, 1. Using split into 0.5 would give slightly different results. (probably smaller size needed to get equal wts of 0, 1. = Dirichlet DF)

On 10-11: If the corpus really was found from a deterministic Decision tree, then

This gives a way to try to find that model. However, no such model may exist ... i.e. it's a stochastic system. In this case, we can take a note from PPM:

We will say, end up w. a specific symbol, result, classifier for each (state) of "features". If we are doing predn - trying to compress a corpus, this is not a good. We have to be able to get pc's of other poss. symbols. (Whenever we make a split, we are not 100% sure of that split.)

Consider only binary classifier. Each final class will be 0 or 1 (w. pc = 1 only).

Say it is 1 we have to code a zero! We can do this by backing up only one level. The

(If we have > 2 classifier symbols we may have to go up > 1 level), previous level will give a pd over 0 & 1. Prob is not deterministic. Alternatively we can stay at the final level

& make prob ≠ 0 or 1, because of finite size.

However, when we have > 2 cat. symbols, we do not have to back up at least one level to get all parts > 0. To backup one level we may use "escape" pc: But

I'm really not 100% sure this is the best way.

But anyway, the merger of Decision trees & PPM seems (like a V.g. idea). [It may be somewhat also related to BBNs.]

The merger enables us to get a "pc" for every label symbol from dec. trees. So its own kind of useful Model for string \rightarrow string prodns because very poss. corpus is assigned \rightarrow pc $>$ 0.

Normally, the way Dec. trees work: We have a set of (strings/objects) that we want to classify. We have a set of binary functions that do (partial) classification. Each R funct. will classify an object into a (or more) categories. Each category is then acted upon by another R funct that further categorizes the Obj, etc. until we come to last R funct that categorizes Obj into a set of categories. T. under "next" type of R to use depends on it.

Output = "category" of previous R.

We visit the Training Corpus! In each step we apply a binary set of R functs to each member of the set of examples. At each point, a R funct will have created a sub set of examples ~~and~~ divided up among its own internal categories: each internal cat. will have a population distribution of final classed cats. These distributions can be used whenever an "escape" character = "unknown" = problem into that R funct.

[50] Remember, when estimating pc of "escape", the value of α vs γ size of context escape from: Often α will be quite small: < 10 . Even if corpus is very large, α will still be very small, because γ context will be long string. If there is only 1 symbol type following α largest context its pc is $1 - \alpha^{\gamma} = 1 - (\text{escape pc})$ (not 1.0).

Each F_i function has 2 or more output categories: each cat. is a terminal F_j . Each cat. category will feed into another F_j . If we use Baum for compression of several texts, the cardinalities of the sets of terminal cats. in each F_i will increment in each trial, as soon as the final (set of) is known. So we can use it for text compression & get compression comparisons w. other methods.

Useful kinds of R functions: 1) For "ob" functions (ob) of a (space 2) Is it about Math/Chem/Phys...?

2) If it's a math: subcategories Algebra, Geom: Is it eq. solving?

(These) cats. can be (used) (by name) as a part of the TSO. - Then put into situations in which these cats would be useful.

One kind of cat. method has been studied: "continuous" Param X is $> k$. We have a continuum of Rets. - one for each value of k . We can do a nu way split by putting in $n-1$ k values; each category will be $\approx k$'s. This could be regarded as a "quantization" of data - to put a set of continuous variables into useful "bins" - what are good/bin boundaries?

One can do another continuous of "bins" using continuous functions. An example: Age;

All s's (subjects) are ordered \rightarrow x_1, x_2, \dots, x_n \rightarrow $x_1 < x_2 < \dots < x_n$. After such a "split" \rightarrow (See 84.12)

SM: General of Convex SVM scheme: 34

SM & Big Criticism of PPH (Some Argued POM) of Dec Trees! That it concs. used depend on ~~the~~ contexts that have explicitly occurred in the past. Unseen elements that have never occurred, are never deduced. This means (I think) much larger ~~size~~ needed for learning (i.e. slow (very)).

For PPH, slow (very) (large ~~size~~) is not so important! — what is important is that \rightarrow universe of concs available very large, relevant to poses to GP solved. — i.e. PPH can need excessive ~~size~~, but it fits objects ~~and~~ concs for problem PPH — that's O.K.! \leftarrow But! Check on this!

Decision Trees: Any decision (classification) loses info. The more of a good decision (tree) is that it loses as little as possible.

1.2: (B3, 40) space: How do we choose continuous or discrete splits? Well age is height (is wts, etc) could map to a pt. in a dimensional space, ~~rather than~~ a pt. on a line. Then later, ~~the~~ partitioning of that space would give "splits" ⁽¹⁹⁵⁾ SVM's is one way of partitioning an n -space (by \rightarrow by dimension). A particular cut can be unconnected ("irreducible" is \leftarrow is ≥ 40 yrs.)

(195) brings us into two ideas of linear (SVM) non-linear partition functions. — A very large space of "Ksch". But, unlike SVM, we have a continuous space for partitioning (is \leftarrow maybe not continuous, but certainly not a function of only a few "support vectors").

ANN, radial basis functs, etc. can (perhaps) be usefully ~~regarded~~ regarded as decision trees on a continuous space: The idea of (10-11) (to use as little info as poss. by many partitions) is a universal conc.

Lots of ways to do partitions in n dim space: One is (maybe Radial Basis funct?) to cut a pt. ~~the~~ Gauss is pts in space of value +1 or -1 Do what \geq distance from all pts (Distance can be by \downarrow func: Gauss, $\frac{1}{r}$, etc! It need not be ~~the~~ be. normalized, since ~~only~~ finite sums \rightarrow vs ∞ . If wts. $\begin{pmatrix} < 0 \\ > 0 \end{pmatrix}$ do $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ One can ~~use~~ \rightarrow a n pt space of \rightarrow by dim. before decup. 20-23 (at SVM)

One can do all discrete partitions first, then all continuous space partitions.

R-DANN: This may give a better (no overfitting) way to do ANN — But, one will have to minimize no. of nodes \rightarrow bits

② Way to get good results out of ~~large~~ ANN: Divide convex into $\begin{pmatrix} C_0 \\ C_1 \end{pmatrix}$ (e.g. $\begin{pmatrix} C_0 \\ C_1 \end{pmatrix}$):

-29 Using random initial cuts, ~~to~~ use ~~forward~~ ~~pass~~ ^{feedback} ~~pass~~ not to hill climb a stepwise way to max conc on C_1 . ~~do~~ ^{many} random ~~initial~~ ~~cuts~~ like .29: But if \downarrow of \rightarrow conc on C_1 will be good estimate for feature. I'm not so sure of my reasoning in .30! I. idea is that each initial cut gives a code for C_0 . But, it does not. — it gives a sequence of codes. Picking \downarrow best one, still gives us SOP. If \downarrow random ~~initial~~ ~~cuts~~ each corresponded to a code for C_1 , then ~~it~~ it would be possible to average them w.o. fear of SOP.

34. **SM** Say we have a SM strat. w. a continuous param: ~~we~~ ^{small} ~~we~~ for each hypercube in n space we invest ~~cost~~ ~~amount~~. We can then compare the amount of yield \rightarrow \rightarrow optimum of \rightarrow end of time. It is equivalent to frequent "rebalancing" among \downarrow ~~we~~ n dim param. components

STM

0: $x \in \mathbb{R}^n \Rightarrow$ [I found a very neat way to express all functions of degree d of n variables including all to new degree forms.]

$$f(x) = \sum_{i_1, \dots, i_n} a_{i_1, \dots, i_n} x_{i_1} \dots x_{i_n}$$

$$x_0 = 1$$
 Expresses all quad, linear, const forms.

05 $n \in \mathbb{R} \Rightarrow$

$$f(x) = \sum_{i_1, \dots, i_n} a_{i_1, \dots, i_n} x_{i_1} \dots x_{i_n}$$

$$x_0 = 1$$
 Expresses all cubic, quadratic, linear, const

Because of commutativity of mult., the "a" arrays are very symmetric or we just sum over $i \in j \in k$

~~XXXXXXXXXX~~

for each value of d a known degree of \mathbb{R}^n , we can compute no. of indep. terms needed.

We mult by \mathbb{R} to get no. of coeffs needed of a d dim space.

Say we chose d "a" a random unit hypercube. (we may want a hypercube). \rightarrow 88.01

N.B. If there are d datapts in a d dim space, then ~~represent~~ it ~~with~~ each data vector

is a basis vector, we can probably use a hyperplane to separate it. Set into any

derivable 2^d distinct dichotomies. ~~XXXXXXXXXX~~ If we have d dim space \Rightarrow 2^d ~~we can only~~ (2^d)

(50) in \mathbb{R}^d , ~~XXXXXXXXXX~~ we get 2^d times as many a coeffs as ~~we can~~ ~~cost~~. used for n

But, we could just use one set of 2^d coeffs & choose ~~sub-set~~ sub-set of them for each dimension of \mathbb{R}^d . An easy way would be to have $a_{i_1, \dots, i_n} \equiv 1$ for all indices,

then choose subset of $1, 2, \dots$ at random to get one \mathbb{R}^d component. Easy to do because

we don't have to mult by 2^d . We could also make it less costly but two off: \mathbb{R}^d component

would be identical (No negative weights if they were!) by allowing for coeff "-"

n of d times \rightarrow so even if two components ~~choose~~ choose same set of coeffs, they would be

identical to all have the same signs.

24.14 \rightarrow (mostly separate them. But, what's the chance of a hyperplane $\mathbb{R}^d \Rightarrow \mathbb{R}^d$?

It would seem that ~~plane~~ ^{plane} would be a lock for d data pts in \mathbb{R}^d space.

It would seem that ~~plane~~ ^{plane} would give better separation. — The Vapnik would say

that the "separation margin" is to bring that margin ~~good~~ ^{good} separation. — This way, but,

we ~~sub-optimize~~ ^{functionally related to} ~~so~~ for fixed SSZ choose d

to maximize separation distance.

35.31R \rightarrow This is actually double! ; Maybe easier than SVM?!

We have S , k dim vector many samples. We want a param function that will separate them w. max

distance betw. 2 sets (Int. ~~sample~~ ^{sample} we know ~~where~~ ^{where} we have each case n)

To start we want function that ~~minimizes~~ ^{idm.} minimizes/overlap of 2 types of cases:



we have 6 overlaps here. We want to \downarrow to zero, then

try to get further distance betw. 2 groups of cases.

So the work consists of 2 non-linear optimizations.

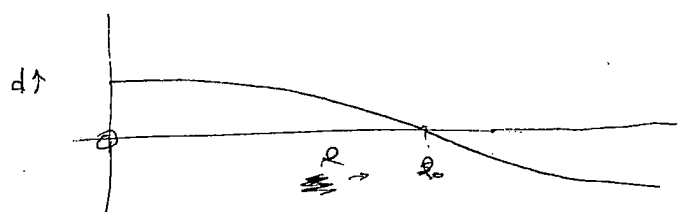
The first is discrete, the second continuous.

39 We can actually do one ^{continuous} optimization — we want to distance "d" to be as small as possible. — ~~eventually~~ eventually as negative as possible.

SM, 32

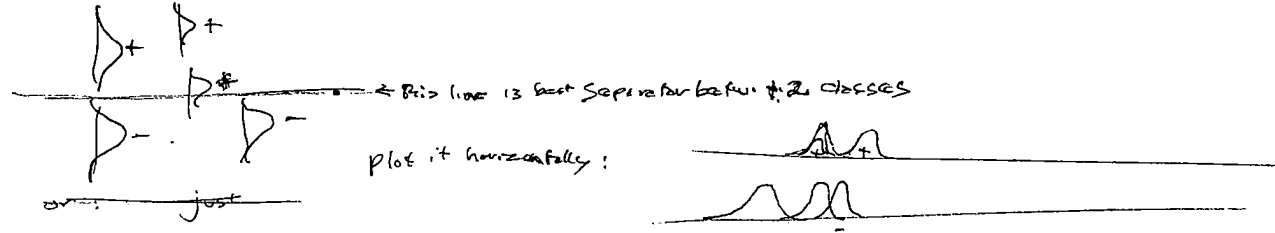
00 while optz of 86.39-40 is continuous, it probly has discontinuous derivatives every
 01 once in a while. Because of disc. deriv., its not so easy to optz. Process functions. $\rightarrow .13 \rightarrow$ 88.20 for better way to deal w. discontinuous Df/Nd/Nd
 I guess we want small d & small no. of ~~parameters~~ parameters in function d

A possl plot of d v.s. ~~parameters~~ d

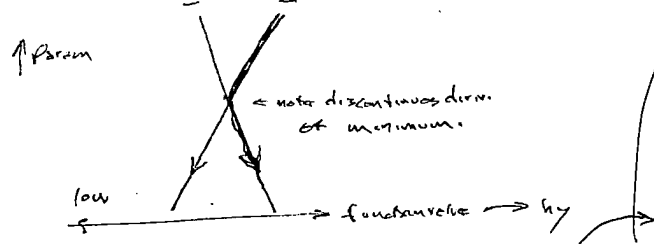


for common data sets, separation is impossible, since we have same vector labels $+1, -1$.
 We may want to remove such cases from the corpus.

3: CC \rightarrow Cheery to deal w. this is to use relatively large x_i when computing derivatives.
 We'd like to be able to get second derivs but probly ~~not~~ discontinuity makes it impractical.
 A way to smooth out the derivatives: Instead of having points for each case, we have Gaussian (or whatever) "means". Unclear as to how one "overlaps".



say no "smoothing". Vertical coord's variation of ~~param~~ a param of f. classifn. function.
 Horizontal is motion of $+d -$ cases.



How can we smooth out discontinuous derivatives (which I haven't yet shown how to do).
 Another way is probly what Vapnik does:
Move along edges toward optimization.

Hor. ~~is~~ in his problem, one always moves in the lines - but not in the present problem.
 But I'm not really sure of the either! He does solve a Quadratic Optzn. problem, hor.

-32 For SM $\textcircled{1}$ Daven has 2 (up/down) drivers (including Driven) have 3 (at least)!
 Try to get 2 drivers (previous day only) but if not available, use, or use - using "escapes"
quality

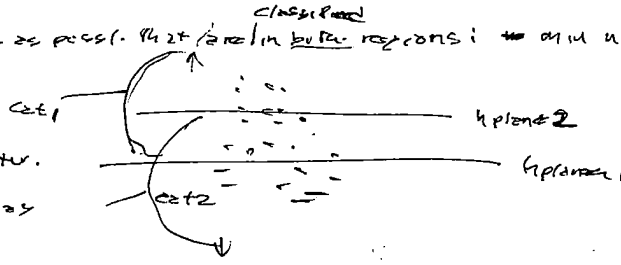
$\textcircled{2}$ TM is constantly looking for "compressible" "spacio-temporal" entities of securities & say other available time series. Main thing sought is ~~temporal~~ diffusion of info in temporal direction. When one has a "home computer" over short time TM always in background, looking for used CPU cycles that it could put to work in f. SM such!

00 : Q: Is it true that problems in which \exists a classifier better deterministic - are usually "Top Problems"?
0005 Σ all of MATH hw, today (w. "Top Problems")? @.

01 | 86.11 : To do summation of 86.05: The condition $z^2 \in j \in m \dots$ is easy to get in nested "D" loops.
The values of $z^2 \in j \in m$ can be stored serially (nearly) in a 1dim array: Together next value of z , ~~the~~ rather incrementing $z^2 \in j$ or m) we just increment z 's 1 dim index.
If z is 0 or 1, we can fill the 1dim array by generating 64 bit random nos - $z \in \text{array}$ & the bits into the array (or other means).

07 For SVM's w. linear Separation: The problem is to find 2 hyperplanes

- ① All cases in cat_1 are on same side of h plane 1 and classified as ①
- ② " " " " cat_2 " " " " " " ② " " " " " "
- ③ P_2 2 H planes are \parallel ;
- ④ There are no new cases as poss. that z^2 or m better separates: ~~no~~ min no. of cases better for h planes:



- ⑤ If there are no cases better for h planes, we want P_{max} as far apart as poss.

0 For non-linear separation; perhaps include $z^2, j, m \dots = -1$ i.e. one needs **exponent**.

Note: If there are no neg exponents, sum of $z^2 \in j, m$ have to be < 0 , if we want ϕ to separate the 2 classes.

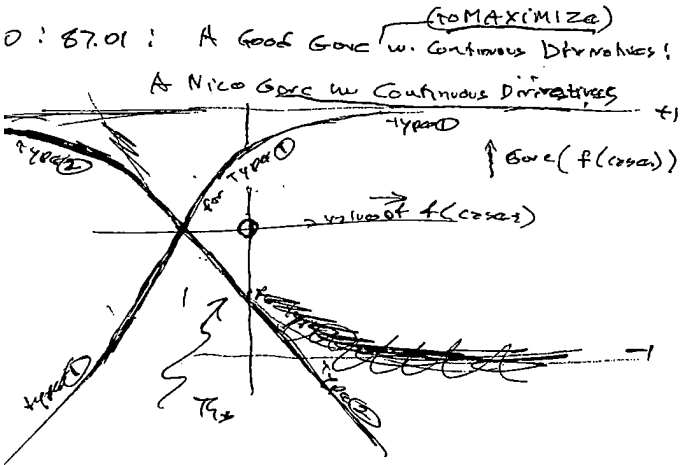
In general, the separation criteria in SVM are not to Best Cost, not probabilistic.

Using a more realistic / ^{practical} Gove may give nice, continuous Gove!

Perhaps if one sticks to quadratic functions, the no. of local maxima in the optzn,

will be small. ^{one?} ~~262~~ \Rightarrow one probably a random initial pt. \rightarrow local peak.

So one repeats ~~263~~ many times and to get a good one



0: 87.01: A Good Gove (to MAXIMIZE) w. Continuous Derivatives! Instead of b. Gove of SVM of .07 = .17!

The idea is best for type ① cases, much motivation to modify to be $f'(x)$, if $f'(x) < 0$: Hvr it separates as

$$G_1(f(x)) = +1$$

$$\text{for type ② cases } G_2(f(x)) = -G_1(f(x)) \text{ so}$$

$$\text{we never see } G_2 \text{ at } G = -1.$$

$G_3(x) = \dots$ There should be a cheap way to compute $G_3(x)$!

$$G(x) \approx \frac{x e^{-x}}{1 + e^{-x}} + 1 \text{ for large } x \text{ it is } 1$$

$$\text{for large neg } x \text{ it is } 0$$

$$= \frac{(x+1)e^{-x} + 1}{e^{-x} + 1} \approx \frac{(x+1) + e^x}{1 + e^x} = \frac{1 + e^x}{1 + e^x} = 1$$

ETA

504 3K/mm

2hrs.

88

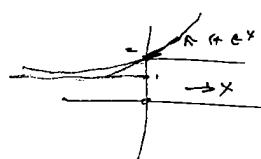
$$G(x) \approx \frac{1+e^x+x}{1+e^x} = 1 + \frac{x}{1+e^x}$$

look at $\frac{x}{1+e^x}$

$$\frac{x}{2+x}$$

$$\frac{1}{1+2.7} \approx .3$$

$$\frac{1.5}{1+e^{1.5}} \rightarrow .100$$



$$\frac{x}{1+e^x}$$

has a peak at $x > 0$

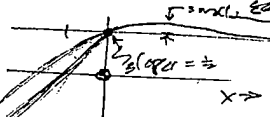
$$\ln x - \ln(1+e^x)$$

$$\frac{1}{x} - \frac{e^x}{1+e^x} = 0$$

$$x e^x = 1+e^x$$

X	0	.1	.2	...
$\frac{x}{1+e^x}$	0	.1	.2	...
		$\frac{.1}{1+1.1}$	$\frac{.2}{2+1.2}$	0

so it looks like



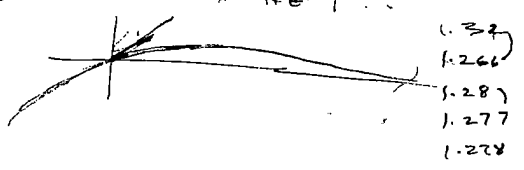
$$G_2(x) = -G_1(x)$$

$$G(x) = 1 + \frac{x}{1+e^x}$$

is a good but it has a slight bump. Probably not so much

It might be cheaper to use $\frac{x}{1+2^x}$ because power is

very rapid approach for 2^x (Power is for \log_2 ... but is slower for 2^x)



89

90

Q's about "degenerate solns" use $G_{1/2}$: The function F concavely expands, contact is most near.



$$(x+1)$$



91

$$G(x) = 1 + \frac{1+x}{1+e^x}$$

looks better!

No!

$$\frac{1+x}{1+e^x}$$

has no max values

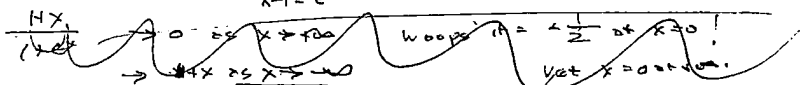
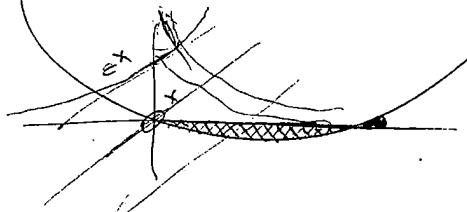
$$\frac{1}{1+x} = \frac{1}{1+e^x}$$

$$x = e^{-x}$$

$x \approx .567$

$$\ln(1+x) - \ln(1+e^x) \rightarrow \frac{1}{1+x} - \frac{e^x}{1+e^x} = 0$$

which has no real soln. = peak!



so 2 peaks at $x \approx .567$

(is this Euler's const? ... unlikely)

$$\frac{1+x}{1+e^x} = x \text{ at } 1 \text{ peak which is } .567 > f(\cos) \text{ so its worse than } .05$$

Hvr, I should check to proof main eqn.

$$1.278 \times 2 = .556$$

with got peak just about .05 just about 1/2 of .56714

$$L \text{ test try } G_1(x) = 1 - e^{-fx}; G_2(x) = 1 - e^{+fx}$$

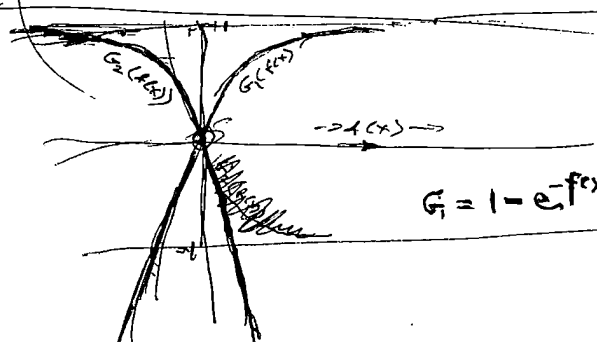
As is this binary descent scheme would seem to be feasible, programmable;

I'd have to ① find a good problem to compare w. SVM.

② Get a good non linear optimizer; The one in "Num Rec." may be adequate

but not sure.

93



There is extreme punishment for very bad cases!

This seems to be main complaint about

this Gorc.

$$G_1 = 1 - e^{-fx}; G_2 = 1 - e^{+fx}$$

10: For cases in which it's same vector in both class 1 & class 2, t. base piece
 at $f=0$ $1 - e^{-f} + 1 - e^{-f}$ is min at $f=0$ $e^x + e^{-x} = 2 \cosh x = \min$ at $x=0$

$1 - \ln(1 + e^{-x})$

desired
 is such a function - But on 92.10 - 93.10 I end up not using it!
 TM could put them in

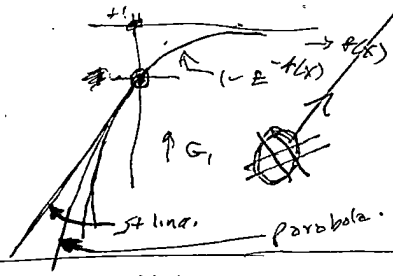
I still haven't related to fuzzy classification schema to PC!

Say we had a perfect classifier: How much info work would be needed?

Att! A way to fix it. G curves.

$G_1 = 1 - e^{-f(x)}$: patcher a st. line. tangent to G.

Corro, so b. first deriv is continuous (no 1st second deriv would not be!)



Att! Match curve to second deriv.

Curv will give a parabola - which is much slower than $-e^{-f(x)}$ (!).

They may be using for certain N.L. or (Optim. Alg.).

Will this make it much more time consuming to calculate e. a 2nd, first, & second derivatives?

For N.L. opten, use full matrix method because all second derivatives are available

→ Since we have a full Hessian is a great thing because we have invert it, matrices - so getting their eigenvalues is not much extra work

So we can get Hessians for "flat min such" - but maybe be sure I have a published

paper for using the Hessian. → See 92.10 ff for a better way to understand this stuff in a v.g. probabilistic way

Partial derivatives: $G = G(f(x))$. $\frac{\partial G}{\partial x_i} = \frac{dG}{df} \cdot \frac{\partial f}{\partial x_i}$

$\frac{\partial^2 G}{\partial x_j \partial x_i} = \frac{\partial}{\partial x_j} \left(\frac{dG}{df} \cdot \frac{\partial f}{\partial x_i} \right) \cdot \frac{\partial f}{\partial x_i} + \frac{dG}{df} \cdot \frac{\partial^2 f}{\partial x_j \partial x_i}$
 $= \frac{\partial f}{\partial x_j} \cdot \frac{d^2 G}{df^2} \cdot \frac{\partial f}{\partial x_i} + \frac{\partial f}{\partial x_i} \cdot \frac{\partial^2 f}{\partial x_j \partial x_i}$

$\frac{\partial^2 G}{\partial x_i \partial x_j} = \frac{\partial f}{\partial x_j} \cdot \frac{\partial f}{\partial x_i} \cdot \frac{d^2 G}{df^2} + \frac{dG}{df} \cdot \frac{\partial^2 f}{\partial x_i \partial x_j}$

I'm not sure if sum of this formula!

If I use finite diff know no derivatives > 2 I could check this numerically, exactly.

$G(f) = f + f^2$, $f(x,y) = 1 + x + y + x^2 + y^2 + xy$
 $\frac{dG}{df} = 1 + 2f$; $\frac{d^2 G}{df^2} = 2$

at $x=y=0$; $f=1$. $G=1+1=2$

$x \uparrow 1$; $\frac{\Delta G}{\Delta f} = 4 - 3 = 1$ $\Delta f = 1 + 1 = 2$

$y \uparrow 0$ $\Delta G = 3 - 2 = 1$; Also $\Delta G = \Delta x + \Delta f^2 = 2 + 1 = 3$ so it's correct.

$G(x,y) = 2 + 2xy$

Say $x \rightarrow 1, y \rightarrow 1$; $\frac{\partial G}{\partial x} = 1 + 2 = 3 = \frac{\Delta G}{\Delta x}$ $\frac{\partial G}{\partial y} = 1$

Back to $\Delta x = 1 \rightarrow x=1, y=0$; $f \rightarrow 1+1+1=3$; $G = 3+3^2=12$

$\frac{\Delta G}{\Delta x} = \frac{12-2}{1} = 10$

$\frac{dG}{df} \Big|_{f=3} = 1 + 2f = 1 + 2 \cdot 3 = 7$

$\frac{\partial G}{\partial x} = 1 + 2x = 3$

so $\frac{\partial G}{\partial x_i} = \frac{dG}{df} \Big|_{f(x)} \cdot \frac{\partial f}{\partial x_i}$
 $3 = 7 \cdot 1$

$\Delta f = 3 - 1 = 2$

Try making x and y vary small & ...

so $x=0, y=0$: $f = 1 + 0 + 0 = 1$

$G_{f=0} = 1 + 1 = 2$
 $G = 1 + 0 + 0 + (1 + 0 + 0)^2 = 1 + 0 + 0 + 1 = 2$

$\frac{\partial G}{\partial x} = 3$

so $\Delta G = 3 \cdot 2 = 6$ $(1 + 0 + 0)^2 = 1 + 0 + 0 = 1$

5.16.05
BTM

GOOD!

A New Genl. soln. of S. Predn for strings. 33ff Also note: 30 may be "developable"

for $x=y=e$: $f = 1 + e + e + e^2 + e^2 + e^2 = 1 + 2e + 3e^2$.

$G = f + f^2 = 1 + 2e + 3e^2 + (1 + 2e + 3e^2)^2 \rightarrow (1 + 4e^2 + 3e^2 + 4e + 2 + 6e + 13e^2)$
 $1 + 10e^2 + 4e$

I hope ~~to~~ distinguish betw x & y , — just call them x & y but remember properties!

See post analysis
92.10

Expo: Start Lecture 1 by explaining how just about all Ling pems P vs P have been walk: P vs P . Ling. used has not been universal!

Process various kinds of ~~the~~ claims to 'universality'.

- 1) LZ compression (four states bugs)
- 2) Feed forward Neural nets: Can approximate any ~~function~~ almost continuous function.
 But large of data can be very big.

2) Fourier ^{expansion} ~~series~~, Taylor expansion: needs many data pts: each method recognizes certain forms easily.
 Fourier recognizes sine waves; Taylor recognizes integral powers.
 ortho ^{sets} funcs recognize funcs \sim to over set: ~~we~~

4) Any formalism for predictive functions that has no "halting problem" ~~is not~~ cannot have a complete set of be complete — it cannot include partial recursive.
 One way to detect absence of P r. funcs. P r. funcs ~~in~~ ω will occasionally run out of computation memory. If your system does not have this problem — it is not using P r. funcs & it's incomplete.

On the other hand, if your P r. funcs occasionally run out of time or memory, this is not really ~~a~~ proof that it's using P r. funcs. . . .

SN On string subset funcs (P vs P Problem): One way is to have O_j be a general set of P r. funcs. That for each Q we want to find shorter functions that give to A def. O_j operates on Q to "simplify" or "condense" it. — P r. R gives A .

* SUMS (Gazd) could do job of string prediction in folg. way(s)

1) Use Binary string corpus, t . p . of next bit is given by funcs. Next look at past's classify past's predicting O or 1 . Each function gets score, depending on how many bits are correctly it got, wtd by its apriori (\rightarrow psu. length). Since each bit has $p > 0$, all products of sequences between p 's > 0 .

2) we can use radix > 1 (say 40 radix: 26 alphabet to nos space, period, comma, hyphen.)

33

5PM

00:

We can do PCs by functions on post that have 40 character output & score from same words & binary is scored. For short corpus is large alphabet alphabet: because many words are used & symbol of only symbols that have occurred up to now, & use an escape pc, for new chars that have never occurred before.

One nice thing about f. forgg. is that this s-prodn. is an easy general of d-production!

This f. forgg. is descr'd in terms of set-valued prodn, it would work as well for ordered & usual Q.A. prod.

79.40
90.17
91.03

Going back to the binary classification problem for real vector input: I really want a more probabilistic approach.

If categories are sharp edged, the params of predictive functions give discrete changes in classification errors.

for very large corp: apparently not so big a problem - but still, second derivatives can be very noisy!

perhaps a better approach! we have these corpus of vectors, we classify data \vec{x}_i (to ± 1) to new h vectors. $\{\vec{x}_i\}$ is + (corpus data set). If we have functions that map \vec{x}_i onto $+\infty \dots -\infty$ ($\in \mathbb{R}$)

we want to devise a (symmetric) function that assigns a pc to a $f(\vec{x})$ value, such

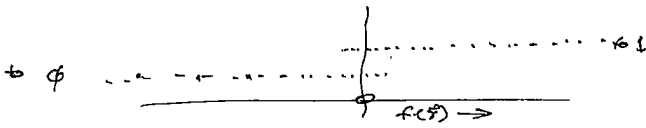
its being mapped to ϕ v.s. ± 1 .

Say we get an $f(\vec{x})$ like so: — that's almost separated.

2 classes — but some overlap.

I'd like a function of f that gives to pc of ± 1

\vec{x} having $f(\vec{x})$ value between ± 1 .



presumably the 2 pc's would sum to ± 1 .

we could use a Gaussian d.f. — decide on its μ & σ .

we could ~~use~~ $\mu=0$, $\sigma=1$ then try to find an $f(\vec{x})$ that gives best

a psip. Well, if it did give a function that separated the 2 classes, then

± 1 . we would want to modify f . function so the Gaussian in the "middle" was

very small w.r.t. distance between the 2 classes' "support points", so if complete

separation occurs, we get no pc's — or pc's = ± 1 only — but function itself maybe costly.

If there is no complete separation (i.e. overlap exists) then we have to "scale" $f(\vec{x})$ or change μ & σ to optimize f 's predict.

But, often, we can reduce the no. of params of f , until separation is no longer "possible" (or we have not been able to separate w.r.t. given CB).

On, there may be cases in which a given vector has a classification of ± 1 sometimes

& at other times (for the same vector) — so separation must be impossible.

If complete separation occurs, we can measure the pc. Other functions with fewer parameters may give a greater pc for the data, but.

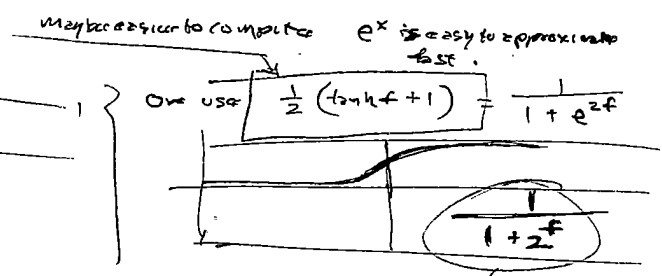
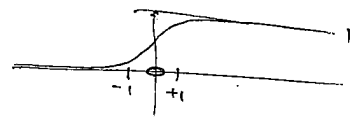
So we take a weighted mean of the predictions given by all possible $f(\vec{x})$ functions — so we always end up with pc's $\in [-1, 1]$ Good!

~~BBN's can be used as a kind of f(A, B) ...~~ 94.40

STM

20

Say $G_0(f) \equiv \int_{-\infty}^f e^{-x^2} dx$
 $G_1(f) \equiv 1 - G_0(f)$



23

We want $\sum_{k_i \text{ class } 0} G_0(f(x_i))$ $\prod_{k_i \text{ class } 1} G_1(f(x_i)) = \text{MAX}$

So, w.r. discussion of 92.10 ff, this would seem to work fine! — also note that G is not in present case.

its derivatives are continuous.

In some cases, we have large SSZ & complete separation; we use n params of f . If we used $n-1$ params we wouldn't get separation. Hrs, $n-1$ params gives a slightly lower pc for f corpus, so while its pcs for ~~cases~~ cases are not all 1, they get much less wt. (The not ϕ out.) That f func. ~~is~~ actually separated the corpus.

21 92.22

About the fixed algo, $S=1$ model. NOT BAD! What we are trying to do is get an $f(x)$ that combines n Gaussians, & uses good pcs for the corpus. Also, in f cases in which complete sep. does not occur, we get a Gaussian in n params that define f , & f 's determinant gives us a measure of the wt. of f in ϕ . — So we know how much wt. to use in combining it w. other solns. For n params, we can try for max source solns. & use their params in Π using associated wts.

good!

20

Actually, getting a P.D. for 0 or 1 as f next bit in a binary seq. looks like a standard T.S. problem for ALP. We can use many of procedures in Π (EXLP) or use fewer. spread across Π . The I'm still very uncertain about this last! One way is simply to try to get a numerical output betw 0 & 1 (as a prod of f's of \pm) & we could use a numeric output from $-\infty$ to ∞ , then squash it down betw 0 & 1 — T. idea being, that any function from becomes legit. I guess LSP can express an arbitrary string to read.

22K

25

Is there a close relation betw .24-.25 and 92.10-93.10? T. "squashing function" that brings $-\infty$ to ∞ down to $(0,1)$ is .00-.03 : $G_0(f) \equiv G_1(f) = 1 - G_0(f)$ Perhaps .24-.25 is a small (if at all) Gamma of 92.10-93.10! Is T. squashing funct. any or somewhat arbitrary.

2 A Big difference betw .24-.25 & PPM

(They both do binary prod), but PPM requires a very large corpus & enormous storage facilities for f regys found.

32

HMM: We could use SVM to do string \rightarrow string prod! Use string to read ~~the~~ functions in Π \rightarrow LSP. Use support vectors to decide if 0 or 1 is next symbol. — or use 92.10-93.10 for classify. of next symbol — or b. gamma of 92.24-.25

In using SVM, as we add data to it. Time series, say, we keep a same hybrid space & assoc. functions, but we slightly move the separability by per plane. To get probit we, we will need a way to deal w. partial "complete separation" — It's likely that Plane is much SVM LSP for ways to do this. (92.10-93.10) could work also — maybe better, maybe worse

0: IN SVM, large distance betw Support vectors means we need less info to specify f.
Separating Hplanes.

[SN] CIAP \rightarrow UCP \equiv Unintended Consequences problem.

03 Expo ID "Why I am not now working on Renft. Medicines". Short paper.

04 \rightarrow How to express "Cure Cancer" as a "well defined problem" (related to 1)

05 ③ The science into "science in ^{startree} ~~startree~~" (activity is epistemic "begin ^{startree} ~~startree~~")

Experiments that have assumed our present (narrow) theories of physics ...
Unlikely to ever even "succeed".

09 Actually ② C.04 is not yet in to shape for "Expo"! I'm not so sure my "solution" is practical.

0 - adequate \rightarrow How do you explain to TM what an "adequate approach" is?

Much of the reasoning (rational/scientific) probabilities.

One way would be to ask TM to devise a decision rule (for what to do at outcome of each experiment) so that the probab of get a whetiz deduced to be "cure" (w. 10 yrs, say) is Maximized.

I.e. TM can view this as a "ordinary" (But very diff) OR problem.

A way to "killen" to various sub-problems: Have a human do the planning of f.

uscl \rightarrow see just what is being done & how it could be optimized.

A big problem is that the results of experiments can be many many different possys, so we have to find good ways to categorize them & use the categories to decide on "what to do next".

Back to string \rightarrow ^{down} ~~that~~ \rightarrow Being functions; The recent tricks (like SVM's)

have been like using ANN for ~~logic~~ digital logic. Like walking dog - surprising that it works at all - never mind that it doesn't ~~work~~ work well!

\rightarrow Try using methods humans used (Bernoulli's seq. counting pc's "log's role")
 \rightarrow try to organize them.

But do write Review of recent stunts in string set get a nice relation \rightarrow to ANN (SVM is also close to ANN).

20 [SN] Reading stuff in kernels, Rad basis funts, SVMs | ^{maybe read forward ANN, & fuzzy Logic} in "Pattern classification" (Peds, Hrb, Stark) book.

Try to unify this all w. my own version of SVM's: Can I write a EV program for All of these? - Showing just how they differ & suggesting new additions to "T. language".

Do BBN's fit into this?

Ordinary ff. ANNs ^{feedforward} \rightarrow a prob w. 1 output cu (perhaps) be re-used as a realization of 93.0 - 93.10! Also, as a "linear hyperplane separation" system, but it ~~is~~ is

also \rightarrow tank x as "synchronizing func.", we can translate it to 0,1 (by $(\tan x + 1)$) and get pc's! If we keep ending ourselves to feed forward configs, we can

get loops & oscillation (\equiv "non-convergence") may be approach Turing Universality!

BBN's can perhaps be regarded as a type of ff ANN!

STM

Firstly, I want a clear understanding of "kernel" as used in "Patt. class" "

It is used in a general way to describe many "Patt. desc" methods

SN CMOS Logic is a great squashing function: (Also TTL logic!) - But in both cases the "Analog" (i.e. ±0.001) regions are too small.

T. Foggy (94.28 Feb) suggests that I might use RANN or Chris Moore's (Humans)

"Universal Analog Computer" (these last can be simulated in Digital to experiment w. them.) For ~~RANN~~ RANN: I have a fixed set of "Nodes". Each node can have

Several ~~many~~ ^{simple} inputs, & its ^{simple} outputs can go to many other nodes. So say we decide on ≥ 100 wts. machines & we decide on which half the nodes is v. outputs & which nodes are inputs. (We could have special input, output nodes that do no summation or squashing.) Anyway, we $\textcircled{1}$ have to connect the 100 wts between nodes made pts (as input & as output) & we have to assign wts to the nodes.

It sounds like too big a problem to do a brute force search for each new ~~QA~~ QA to try.

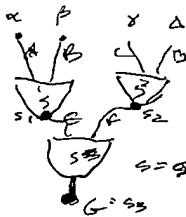
Or for each update (improvement of ~~TM~~ "TM" - or just for "Meta-induction" in

Phase I.

In usual RANN, the connections are kept fixed & wts are "updated" w/ performance on (at least 1 to pair. Simple first order partial derivs are used: $\frac{\partial G}{\partial W_{ij}} \propto w_{ij}$.

Second partials not used. (Too many wts: No. 100 x 100 matrix is solvable.)

How do compute second order derivatives?



$$\frac{\partial G}{\partial E} = \frac{dS}{dx} \Big|_{E+F} = S' \Big|_{E+F}$$
$$\frac{\partial^2 G}{\partial E^2} = \frac{d^2 S}{dx^2} \Big|_{E+F}$$

$$\frac{\partial^2 S(E+F)}{\partial E^2} = \frac{d^2 S}{dx^2} \Big|_{E+F}$$

$$\begin{pmatrix} D-C-(B-A) \\ B-E-S+A \\ D-E-(C-A) \\ D-B-C+A \\ S_2 \cdot D \cdot S_1 \\ = D \cdot S_2 \cdot S_1 \end{pmatrix}$$

N.B. I know of all certain of these formulae!

$$\frac{\partial^2 G}{\partial E \partial F} = \frac{d^2 S}{dx^2} \Big|_{E+F} = S_1 \cdot S_2 \cdot S_3 \Big|_{E+F} + S_1 \cdot \phi$$
$$= S_1 \cdot S_2 \cdot S_3$$

Actually I'm not entirely certain
into 1. T. input to third nodes $S_1 \cdot E + S_2 \cdot F$.

$$\frac{\partial G}{\partial E} = S_1, \quad \frac{\partial G}{\partial F} = S_2.$$

$$\frac{\partial^2 G}{\partial E^2} = S_1^2 \cdot S_3; \quad \frac{\partial^2 G}{\partial E \partial F} = S_1 \cdot S_2 \cdot S_3; \quad \frac{\partial G}{\partial A} = \alpha \cdot S_1 \cdot E \cdot S_3$$

at any rate, one

00 (59.03) ? This 2 dim convolution would not be exactly what I want, but maybe close enough to be giving fairly good contexts. A problem w. it is that it doesn't give

what are necessarily "properly connected", subgraphs. It may be that a 1 dim

63 convolution would do as well. The kernel would not really have adjacent nodes on it ... it's could have large gaps.



Note that a tree can be generated in (at least) 2 ways: Bottom-up or top-down. Bottom-up

was used in 8. Derhof AZ lang. in Appendix of IDSA report.

R girl ... Ferraris described a method of doing GA in which the notation was

10 deriving the graph of (03 L) was 1, 2, 3, 4, 5, 6 ... 15 !! She used normal ("GA" rather than "GP") crossover. She said she got unusually good results !!

The string 23 would have the "context" "9567" on one side & "1" on the other. ← used!

1, 2 would have context 3 ← useful.

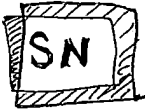
2 would have context 1 on one side (useful), 3 on the other somewhat useful

3, 4, 5 ... on that other side ... not so useful.

5 would have 4 as (useful) context 6, 7 as not so useful context.

T. notation Ferraris uses is not clear in her paper, but it is in

20 (later papers but are (dense (parameters, extensions) or Non Work.



A Maximum Entropy Approach to Sampling is EDA - The Scala (connected) case

0 chaos trials Solo, Mühlenstein

Ps old / Ciarp 2003, ps why this title?

9/12/03

This seems to discuss problem of finding proper context trees, act. — But I'm not sure!

8 pp. — try to read it.

Key discuss (binary) Bayes nets, & the

relevant trees.

Ching et al to

Tree context decay 2003

on HP 467

Physical Realization of string A.I. Look for various "technologies" that achieve

30 T universality. Find one in which Macros are easy to detect — i.e. Easy (educated) commonly

occurring contexts. Matchers need not be exactly correct. It should be easy to verify

that a particular "technology" is universal w. this property. Lots of technologies

are universal, but it's hard to prove they are universal. Ones w. this desired property

→ should be a lot easier to show they are universal.

One approx method for "context matching" is discrete convolution —

using analog of FFT. (Context can be derived). Finding FFT of strings plus new

(context) strings should be possible in an incremental way, since FFT is linear!

— Not majority of "convolution" is defined.

20: (SPAC) (97.29) : I can refer to the direct QA induction as "Primary induction": it uses Turing models.

The "secondary induction" ($\equiv TM_2$) need not be Turing universal (but can be). It can use a variety of induction techniques mentioned in 97.26-29. Eventually we can use the primary induction to do secy induction — Any improvement in this will have to be via Phase 2.

For Phase 2: The "Grammar of Inductive Systems" can be obtained by studying the systems in "Patt. Classifi." = "A.I. & Mod. Approach". In "Pattern Classification" systems are (I think) for induction only (ok. to boost secy induction of Phase 1), "A.I. & Mod. Approach"

Discusses other kinds of problems as well. — O2 probs: "Patt. Classifi." also discusses O2 probs soln methods of various kinds.

What I'd like to do is get a larger corpus for QA.I. — possibly use some "pre-proposed" sets of problems. Then see which induction systems can

compress most, "per unit time" (i.e. exact time may not be found related to eq. of $\sim 4TM^{1.24}$...)

21 : A major problem will, of course, be the initial TSO. It will be easier out. Trained if the machine is very "fast" (by Fpc), since the examples don't have to be ^{as} close together (perhaps) don't have to be put out so carefully. For the early TSO, TM_2 is it going to be very "smart", because it has such a small corpus.

22 → At this point, I really need to write a good Review of just how the ~~present~~ devices in "Patt. Classifi." = "A.I. & Mod. Approach" are to be applied to secy. Induction. → See 97.26-29 for a list.

Understanding Optimization +30

I Duda, H, Stork: p 319 - 324: § 6.9.2 & 6.9.3: 6.9.4
 6.9.2 gives Newton's method, but it involves inverting Hessian matrix & is $\sim N^3$ problem - too long.
 6.9.4 is "Conj Grad descent" which gets some gains Newton's method, but only requires N jumps.

[presumably, each jump takes time $\ll N^2$!] & gets some results Newton's method.
 As I suspect, this must be related to Inverting Hessian in $\ll N^3$ time - so this is a way to invert sym. matrices (perhaps positive definite) in $\ll N^3$ time!

Also, If it results are the same; ① why does one value of " β " give different results than another? — "One deals w. hyperorder non-linearity's better"!

If Conj Grad is an approx soln to inverting Hessian; It can be used iteratively to get a fast, successive approximation to the Inv. of H. Matrix.

When Hessian has no inverse (singular system, H is nonsingular at least one $\partial f/\partial x = 0$)

Does "Conj Grad Desc" give pseudo inverse?

This is the pseudo inverse of a sym matrix w. one or more zero ~~eigenvalues~~ eigenvals. Well; the inverse involves all evns but 1. zeros - which leaves invariant. — easy to see this in the coord system in which the original (symm) matrix is diagonalized (i.e. always poss. to get to this coord system by Rotating a symm matrix).

SN on BOOSTING: This is an induction method that uses several (or many) PEMS for prediction. When given corpus of many cases, it tries to find a way to combine the PEMS to work better on harder or b. cases in which the system had worked most poorly in the past.

This sort of technique would seem to be very useful for a QA corpus! So "check it out" more carefully... see that it can be adapted to more universal induction systems than are normally used in "Boosting".
 "AS2MA" makes very strong claim for "Boosting" over 6.66 || Also see Pitt. Chaitin "PLATON on Boosting".

Re: **Understanding Optimization**: Knowing about local minima & how to deal w. them is an essential direction of "Understanding". Various optiz techniques in Sim anneal & GA-GP have ways to deal w. this. In ANN research "Momentum" & "Accelerators" Newton's Method (in Ndim) is also a way, & it's that Conj. Grad method (0.09) Also, fittly good Hessian (good statistical fit), then using Newton's Method.
 → It might be possible to get v.g. induction in Ph.D w.o. going to fit it!

5TM

{ On Ordering A's in pc order - .30.
 On ~~some~~ ways to generate B's - .20 } Classic Problem:
 Expt: Binary Solns

: 99.20 & 99.30 2nd both very encouraging! Suggesting that I can get a very good
 phase 1 w.o. Phase 2, using Grishin tried techniques!

SN on ^{emergence} ~~conv~~ of Newton's Method! Use stabilizer (A.M.), but have size of "Box" of
 def. accordance" be k . Large k is by temp. in. Sim. Annals. —
 Start w. k
 Start w. k is slowly \rightarrow if due to smaller k values.
 start w. k

: On s-ways: we can have a fixed or growing no. of symbols (including delimit symbols
 = "tokens"). Flow termination One of them is to terminate symbol, (Δ , say).
 T. pc of Δ will be constant in simple ways, but more generally it can be any
 function of the previous symbols. This gives a new way to deal w. Exp comp.
 If the prod $\Delta = \alpha$ for all contexts, we have to "continue" (univ. perhaps) d.f.
 A lang can, using this convention, have some finite strings in its context & some
 infinite strings.

On putting outputs of a stack generator of symbols (including Δ (.20)), in pc order.
 In generating strings (sequentially), keep current bins of w , Δ \in — \in can
 be marked. This \in is \uparrow in n th so \in (in length),
 we first generate all strings of length n symbols; Δ put current bins according to
 length. Δ No expand all strings into shortest bin, by (symbol Δ ~~current~~)
~~from~~ transfer from to new bins (usually not same as previous lowest),
 loop to Δ (.34). As soon as one string terminates, we have an upper ~~bound~~
 on length of shortest string. We may want to no longer keep bins longer
 than that. We will periodically find the shortest terminating string known so far.
 when the shortest bin = length of shortest terminating string was stop. \rightarrow (spec 101.00)

T. Bins. can be "linked lists"
 we later go thru whole list
 when we "develop" its
 strings.

JTM

practical vac. v.s. prim recursive funct.

Normally a for loop could never give an "undefined" output.

However, here is ~~the~~ "(3 ÷ 0)" is one of the functions, ~~which~~ have a "undefined" value.

A "for" loop w. large index (say approximately n "until", in this way. Say the until condition is $f(x) > 0$. Each loop sets new value for $f_1(x)$ and $f(x)$: when $f(x) > 0$, we jump out of loop w. $f_1(x)$. To simulate this jumping out of loop

The recursion eqs. have modified $f_1(x)$ & $f(x)$. for each recursion

~~the~~ Here, ~~if~~ if $f(x) > 0$ sets a flag to 1

If flag = 0, recursion proceeds normally for each ~~cycle~~ cycle of loop.

If flag = 1, ~~each~~ ~~cycle~~ cycle of loop does not change

$f_1(x)$ or $f(x)$ — so when $f(x) = 0$ we keep $f_1(x)$ constant.

Just before end of loop, we look at loop index n : if $x = n$,

~~flag~~ and flag is still zero, then output is "undefined"

Or ~~output~~ $f(x)$ is initially set to "undefined"

only ~~if~~ when flag is set to 1 does it change to a real value.

Alternatively: If in the loop simulation, we end up at end of loop w. no acceptable value for $f_1(x)$, we set $F(x) = \frac{1}{0} \rightarrow$ undefined.

1) Defn. of Strong A.I. (ie. what my goal is).

3 → 00124

Indicators of Current A.I. systems.

Induction is impt. part of capabilities of A.I. (It is not All that is needed - it's major problem)

1) Inadequate (arg. for induction).

How many are inadequate! FFANN, forward pass, poly expansion...; GA, GP (usually).
They are "Universal" in the sense that given exact parameters, they can approximate any function with any desired precision. However they need not be particularly efficient in this - One may need an enormous no. of passes at very big precision to get desired approximation.

There is another kind of universality that's much better suited: it's the universality of an ordinary computing machine, or any Turing computable instructions.

If you try to approximate a sine wave with a FFANN you will need a vast / huge no. of nodes and wts to get a good approx. At no point into the process will the net realize

How, using T. universal approx: ~~it can~~ after a certain no. of pts from to be part of sine curve, the system is able to recognize this fact and make a very accurate approximation with relatively few parameters.

FFANNs are unable to do this. - Any net is not truly universal, is not able to do this in all desired cases.

The inadequacy of most "universal systems" is in their need for an enormous no. of data points before they can approximate a function.

The advantage of Turing universality is the "small size"

I want a kind of outline of to talk:

First: Define what my goals in Strong A.I. are! That if they can be met, a system could probably be trained to do ~~things like passing the Turing test~~ understand spoken speech and gestures and even do tricks like pass the Turing test. - However these are aspects of A.I. ~~that I'm not particularly~~ in which I have little interest.

I will list several important deficiencies in current A.I. research. While many systems are free of many of these criticisms, there is ~~not~~ no system that I know of that is free of all of the deficiencies. ^{In the final lecture} I will ~~try to~~ describe a system that I believe deals with all of the difficulties in an adequate manner

One of the major problems in A.I. is Machine Learning

I will identify the strong A.I. problem with that of Machine Learning: Given a corpus of ~~relevant~~ relevant historical info. - How best to solve a particular well-defined problem. How can the machine best bring to bear the info in the corpus to help in the solution of the problem. We will begin with pattern recognition problems, then later show how an understanding of induction can be used to solve more general problems ~~more~~ very effectively

STM

First, the induction problem. In a simple form, we are given a sequence of symbols and we want to find a pd. of E. not necessarily

Discuss Universal limit of universality: L.E., R. SMITH, polynomial approach
List of pts. I want to make ... in proper order (or note this secondary!)

Some poss. attacks:

The bad news is that you really have to ~~use a universal~~ consider a universal D.F. - which is strictly incomputable?

The good news is that you don't have to ~~use~~ actually compute a Univ. D.F. - you can get good predictions and good estimates on how accurate they are w.o. computing the Universal D.F.

More bad news: Most A.I. researchers don't know how to do this.

More Good news: I will tell you how to do it.

It is possible to get very nice, exact results by ^{narrowly} ~~completely~~ restricting the class of concepts considered. In a benign Academic environment, the instructor will tell the student that he need only consider linear models for ANNs with 1 hidden layer or whatever. The bad news is that "our world" is not given any such information. Its pretty much an "Any Day Goes" situation. He doesn't know how to deal w. this, so he imposes the kind of restrictions his instructor imposed. Not good! The Good news is that I will tell you how to deal with this problem.

"Well defined class of Academic Concepts."

In predicting the stock market: NYSE does have rules about trading stock - but you will not find any rule that says only linear models for price variation are allowed. Its really "Any Day Goes".

To us, that it is not only the stock market that has this (unruly) quality... just about anything else in the world has the same features: There is no real limit on the kinds of laws, regularities, and spontaneities that might occur.

[SN] A poss. understanding of Kolmogorov's Program about representing N dim. continuous functions by ~~the~~ N or $(N+1)$ continuous functions of 1 variable.

If a function is continuous in a region, then ~~it~~ it has a ∂ -derivative over the region. It is given a simple Lipschitz condition satisfied by f. function. It profoundly

limits the no. of functions poss. on a 1^n cube w. precision Δ . Given the point (derivative).

We can say how many functions there are of this Δ \rightarrow sub Δ compare w. point deriv. = ∞

Say $f(\Delta, d)$ is no. of functions on unit cube of grid size $\Delta \geq$ peak derivative = d

$f \rightarrow \infty$ as $d \rightarrow \infty$. Hvr. n. dim. functions can have lots of info in them if they are allowed very

large d. values!

10: On ANN for digital ^{Binary (the N-ary conversion: No. of bits)} time series predn. If p_i is outputs, we want

01
$$\prod_{i=1}^n \frac{x_i}{p_i (1-p_i)} = \max \quad ; \quad x_i = 0 \text{ or } 1 ; \quad \bar{x}_i = 1 - x_i \quad ; \quad x_0, x_1, \dots, x_n = \text{e. f. no. of bits}$$

03 If output is $-\ln p_i$ or $-\ln(1-p_i)$ we want to minimize \sum_i of outputs which is \sum_i can be viewed as error function. So we may be able to use standard ANN fcn, w. non-standard "output squashing function". So we use Newton's method or Conjugate Gradient to find wts.

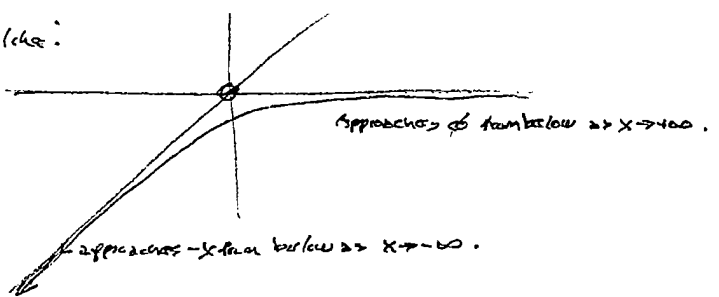
06 Advantage of this ANN over PPM is that it allows more complicated functions.
07 Disadvantage is its small memory cap (200 no. of wts in net). long time to converge w. no. of wts. or (wts). PPM can remember very long, complex contexts.
04 We can use data over & over again - (not really for accuracy.)

Unclear in my mind as to details of this!

0 Also, we can use recurrent ANN. For small no. of wts, + config of wts, we have large search space. As Net begins to converge, we can get "This is probably a loop because it's been running for \approx seconds". How do we handle "undecided" decisions in these trials? - its default from a UTO machine, in which all subsequent outputs become "undecided".

10: 90:10 : $\left[\frac{90 \cdot 10}{100} \right]$ a way back, I wanted a function like:

- $\ln(1+e^x)$ is such a function:
It is the \ln of the squashing function, (.03)
$$\frac{1}{1+e^x}$$



10: .06-.09 seems imp! ANN can use only a few parameters. (That's usual way is to use very many wts & cross validation.) But in general, contrast, say linear/n.l. regression (low param) w. PPM (which stores whole corpus). (n.l. regression can store correct matrix & update it via Kunita (fitted Memory Needed) or Runita (Much more memory needed). PPM w. Kunita or Runita would be better but has its drawbacks ("small") corpus, even though not so good compression.

6.6.05
STH

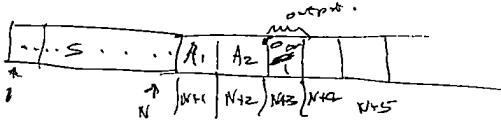
T. HALTING Problem

106
6-6-05
STH

Proof of unsolvability of halting problem:

S is a program of N bits. When S is executed, at address 1 , and addresses

$N+1$ & $N+2$ contain start and end of prog X , then S if it started at address 1 , will decide if X will eventually halt or never halt, a print 1 or 0 at $N+3$, then ~~it will~~



It will go to $N+4$ for its next inst. (which can be "halt")
But need not be.)

Consider: prog X goes from 1 to $N+5$

Set $A_1 = 1, A_2 = N+5$

Start prog at 1 . Eventually it must print

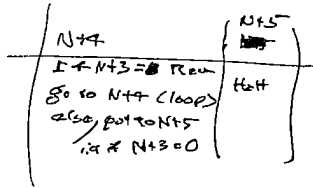
0 or 1 at $N+3$, since we asked it about prog 1 thru $N+5$

If it prints 0 it has decided: prog will never halt. ... but in this case, prog continues to $N+4$ then $N+5$ & $N+6$

" " 1 " " " will halt

$N+4$ and loops to $N+4$ infinitely

So it looks like if prog didn't work it.



STM

Rev:

00

106: Helloing problem: Quick, easy, complete first

105:20 Easy way to get ~~func~~ func, for es. 10! But surcharges to what is unhit!

105.01 may be better. I will have to decide on just what shape of ~~func~~ func for ~~func~~ func.

00 On use of ~~ANN~~ ANN (or possibly RANN) for discrete data — i.e. id. on strings : s. d. on strings

.30 D. func. betw. ANN & PPM in kinds of mto. in models. Beats more security!

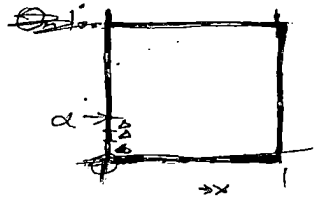
104.30: A poss. way to understand ~~lots~~ lots param approx of N dim. func's by ¹N, 1 dim. func's, in terms of MB^{1-c}.

0

0

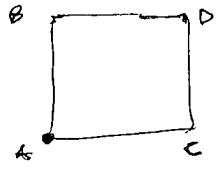
0

10:104.40: On Hell's theorem approx. of functions. If there is a real value of derivatives in a given region, then whatever Lipschitz condition: Consider 2 dim region $0 \leq x < 1$: Grid size $\Delta = \Delta$.



consider α_j , α ; $x=0$, $0 < y < 1$: In distance Δ , func. can move Δ ("3/4 from Lipschitz"). So, in vector of Δ , α can have $3^{\frac{1}{\Delta}}$ values.
In each Δ dist α is function f goes up, down, or stays same. (3 way prob)
However now to α has loss since its change is limited w.r. α , but also w.r. near by Δ 's. The main restriction is if α is ± 1 .

2.05: In 2D case



$B \neq C$ are to same! Then P number $C \neq 0$, or ± 0.2 .
If $B \neq C$ differ by 2Δ , then D can only be \pm max of $B \neq C$.
Figuring out how many contiguous $n \times n$ regions seems difficult!
Say $S \approx \Delta$.

To simplify, consider lattice w. interval spacing and $B \neq 1$.

$\neq 1$	$1, 1 \rightarrow (1, 1) \geq 3$	$0, 1 \rightarrow \pm \frac{1}{2} \pm \frac{1}{2} \geq 2$
$\neq 0$	$0, 0 \rightarrow (0, 0) \geq 3$	$-1, 1 \rightarrow 0 \geq 1$
$\neq -1$	$-1, -1 \rightarrow (-1, -1) \geq 3$	

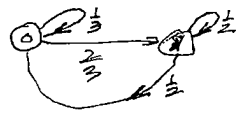
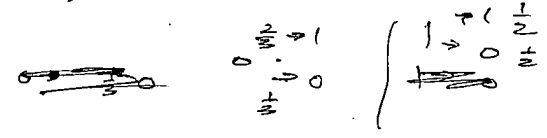
To state the problem: we have an $N \times N$ integer grid ≥ 2 in integer function at $N \times N$ pts: f assignment on f -function: No pts can be > 1 above or below adjacent pts. How many functions

20 Δ possib. ? May depend on adj. grids: Say all α_j are ± 0 .

$0 \rightarrow 0, \pm 2^{\frac{1}{\Delta}}$; $1 \rightarrow 0, 1$

$0 \rightarrow X \times X + 1 \times X$

So a Markov chain! $\frac{2}{\Delta}$ states: $0, \pm 1$:



So with computer can take percent resistance of $0 \leftrightarrow 1$.

$P_0 = \frac{1}{2} P_1 + \frac{1}{2} P_0$

$P_1 + P_0 = 1$

$P_1 = \frac{1}{2} P_1 + \frac{2}{3} P_0$

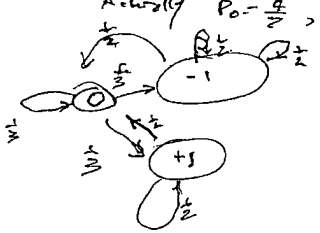
$\frac{2}{3} P_0 = \frac{1}{2} P_1$
 $\frac{2}{3} P_0 = \frac{1}{2} P_1$

$4 P_0 = 3 P_1$

So we can get "entropy" of row next to all 0's.

Let's say $P_0 = \frac{1}{4}$, $P_1 = \frac{3}{4}$, $P_{-1} = \frac{1}{4}$.

Thus ± 1 result follow ± -1 .



$P_0 = \frac{1}{2} P_0 + \frac{1}{2} P_1 + \frac{1}{2} P_{-1}$

$P_1 + P_{-1} + P_0 = 1$

$P_1 = \frac{1}{2} P_0 + \frac{1}{2} P_{-1}$

$\frac{1}{2} P_{-1} = \frac{1}{2} P_0$

$\frac{2}{3} P_0 + \frac{2}{3} P_0 + P_0 = 1$

$\frac{7}{3} P_0 = 1 \Rightarrow P_0 = \frac{3}{7}$

$P_{-1} = \frac{1}{2} P_0 + \frac{1}{2} P_1$

$\frac{1}{2} P_1 = \frac{1}{2} P_0$

$P_1 = \frac{2}{7} P_0 = \frac{2}{7}$

$P_{-1} = P_1 = \frac{2}{7}$

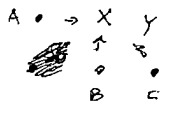
Is it entropy of "t. next row" given by transition $P_i = \frac{2}{7}$, $P_{-1} = \frac{2}{7}$? Reasonable see not indep.

Entropy is more like $S = - \sum_{i,j} P_i \ln(P_i | P_j)$

± Rank P13 will work!

0:10840
03
Total Entropy of next row. Oh, that's it!
t. value of f at a pt. They will have certain pt's. We have a mapping of a set of pt's of confcat into itself. For this confcat. to be "stable", gives a (solvable) set of linear eqns.

for f. can map of 2 ~~states~~ "by" other row!



determining X

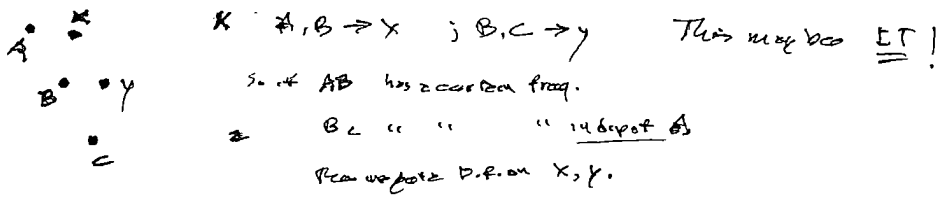
$$\begin{matrix} A=+1 \\ B=-1 \\ C=-2 \end{matrix} \text{ is possible.}$$

Since ignores ~~next~~ relative.
 $X=0, Y=-1$

So even $A=+1, B=-1, C$ is constrained by B only.
So knowing A, B , we can get D.F. for X .

I still don't have a clear idea as to how t. (discrete) maps to.

Maybe diagonal strips for states!



- pass. pairs
- 1 → -1
 - 0 → 0
 - ±1 → ±1
 - 1 → 0
 - 0 → +1
 - +1 → 0
 - 0 → -1

No: $A \rightarrow B$ can be $A=-(A+B)$ or $X=0+ B=±2$

so 5 jump points

$$\begin{matrix} \pm 2 \\ \pm 1 \\ 0 \end{matrix}$$

one way to get $\Delta = \pm 2$

Another way: $P_{12} \cdot P_{11} ; P_{12} \cdot P_0$

$$\begin{pmatrix} P_{12} \\ P_{11} \\ P_0 \end{pmatrix} \cdot \begin{pmatrix} P_{12} \\ P_{11} \\ P_0 \end{pmatrix} = (P_{12} + P_{11} + P_0)^2 = P_{12}$$

~~Another way~~

$$P_{12} = [(P_0 + P_1 + P_2) \cdot P_2]$$

$P_0 =$

$$\begin{matrix} 2 = 2+2 \\ 1 = 0+1 \text{ also } +1+2, +2-1 \\ 0 = 2-2 ; 1-1, 0+0 \end{matrix}$$

5-TM

20

On sets of functions that may be $\overset{?}{\text{partial recursive}}$: "Fractals" is a 1.1.10.10 (broad set).
 Suggests ways to get (perhaps ~~partial recursive~~) (large) sets of functions, that may be occasionally
 "undecidable".

10

The Mandelbrot set does $f(z) = z^2 + c : f^{(n)}(z)$. ~~we start~~ for each value of
 (complex) c , we start w $z=0$ and find $n \rightarrow |z|^2 > 2$ (circle 1).

~~for~~ $n \leq 255$; if $n > 255$ then assign $n = 255$. $n =$ a color assoc. w. c pt. used.
 So there is a way of assigning n to each c value. Of course it takes
 time n to get "breakout" value n .

Another problem: $\mathbb{Z} \subseteq \mathbb{R}$
 \downarrow complex \rightarrow real \rightarrow integer.
 So $\mathbb{C}, \mathbb{R} \rightarrow \mathbb{N}$.

w. a different \mathbb{Z} value we get a different output, so \mathbb{Z}_0 could be proven.
 $\mathbb{C}, \mathbb{Z}_0, \mathbb{R} \rightarrow \mathbb{N}$ integer.

15

To get ~~positive real~~ \mathbb{Z}_0 as a function of \mathbb{C} .

16

$\mathbb{C}, \mathbb{Z}_0, \mathbb{N}$ in use
 $\mathbb{C} \times \mathbb{C} \rightarrow \mathbb{N}$ complex.
 $(c + z^k)^n = \text{complex}^k \times \text{no.}$
 start \mathbb{Z}_0 (largest real (radius), real θ angle).

20

for $y < 2$ values maybe use $\mathbb{R} \rightarrow \frac{2}{1+e^R}$ by \mathbb{R} means near ϕ ,
~~near~~ \mathbb{R} means near 1.

We don't have ~~equal~~ (diffy approaching 1 or 0) — so not so hot.

A better way: take the real (or im) part of $15-16$. It goes ~~for~~ \rightarrow $y < 2$, so

spreads it to 0, 1 by $\frac{1}{2}$ (if tank \mathbb{R}).

T. escape radius R of z cannot be usefully increased. As soon
 as the iterates get outside $\mathbb{Z}(c)^2$, it ~~must~~ must $\rightarrow \infty$.

So essentially \mathbb{C}, \mathbb{Z}_0 (2 complex nos.) are used by $f(z)$ into positive integers, \mathbb{N}
 that are used to represent \mathbb{C} (arg). They could be ASCII.

30

So \mathbb{C}^2 into ~~pos~~ pos integers.
 \downarrow function.

So how can I get induction out of $\mathbb{P}(z)$? \rightarrow 111.02

This is of course kind of full of Wolfram's CATS. \rightarrow Rudy Rucker's expansion into
 Real outputs, \neq (PDE, approx. Solns).

3A

The interesting thing about fractal pictures, is that \mathbb{Z} real or complex region, we get
 more and more "continuous" info as \mathbb{Z} coords of \mathbb{Z} go to higher order bits!
 So ~~all~~ all of the bits (even if most ~~are~~) of \mathbb{Z} , ~~do~~ give good quality info
 in \mathbb{R} . — So it's like a sequential code!

5 pt
 \rightarrow 111.00

10: 110.40: so (10.34 - .40 is an unpt. remark)

Q: How many bits precision do Puz normally use in seq fractal?

2: 110.30: T. for prediction! Say t. entire 2 dim ~~space~~ complex plane as t. entire 4 dim. 2 complex Mand/Julia space is dominated developed into 256 ("colors" = ASCII symbols).

Consider a certain fixed resolution, so t. space is divided into fixed size pixels.
To do prediction of t. sequence x_1, x_2, \dots, x_n (x_i are ASCII or redy 256 integers)
we look for as many cases of t. horizontal (say) sequences x_1, \dots, x_n as we can find.
t. P.D. over t. next symbol is our predn d.f.

To find t. seq x_1, \dots, x_n : we just watch what t. fractal is being created, pixels by pixel. ~~the~~ Puz is actually practical for small n. \rightarrow spec 12.30 for ~~practicality of~~ ^{256 element v.s. 256 element v.s. practicality of} (02-10)

Say we have "coded" x_1, \dots, x_n Puz way: we have 1000 codes that we've

bound: we now know continuations of x_1, \dots, x_n if we want to ask for pixels

Prac: Well, we just look at Puz ~~for~~ points that code ~~is~~ thru x_n , if we

increase resolution, to see Puz in many poss. continuings.

Is this legit? 102-11 seems legit. But .12-.15 very questionable!
Actually, it's an interesting Q! not nearly redicious!
T. log. base 2 of .12-.15 is suggestive (14.10 off!)
Also, it's suggest .02-.11 is not legit!

Say we a priori fix ~~at~~ 2 corresponding to $2^6 = 64$ poss. sequences.

or: just 1 byte = 256 direct possgs. So we want Puz, remember all pixels of "color" x_1 . Then we \uparrow resolu. by 8 bits, 2 we continue all oth. color x_i cases to find ones w. ~~color~~ continu. x_2 ; we remember Puz, then \uparrow resolu further by ~~bits~~ ~~to~~ to get matches to x_3 , etc.

This may be a legit. method of coding t. corpus! — bob. fits 2 24 byte 8 byte sequences & Puz \rightarrow 64 bit precision, 24 byte \rightarrow 192 bit precision.

Since we only want colors that correspond to t. x_i 's, we can save some time if we code ~~the~~ x_i so that the low no. pixels are "most likely",

NOTE: Puz is muddled or a simple minded version of SUMAC!
i.e. t. trick of .12-.15 is what t. simple minded SUMAC normally parts - by decreasing diversity. fivr, if I use t. full (256 incompatible) SUMAC, it's just regular ALL. — which I can approximate by (.02-.11), i use it for small n, for "secondary predn" (Puz). It will not be as good as PPM, because it has so

little many of post: But use it in ll w. PPM for finding really crazy regions! (spec 12.00)

6:18:05
5:30 PM

J.C.

0: Th. Frick of 11.12-15 ^{mito} could give us much larger corpus possi., but suffers @ from diffy 11.32 (Simple minded SUMAC). Also Q of what ~~super~~ superhy resolu needed could be achieved w. reasonable cc. The 11.12-15 tree can be used every once in a while of diversity: let to keep no. of codes for corpus > 1000, say — normally are just continues w.o. ↑ in resolu. The Ptz would deal w. diversity problem (perhaps), then it's Q about what it would one good prodn!

T. problem of how to do > 1000 bit resolu. economically — This may or may not be a problem. The details would have to be worked out. Fracture does seem to be a resolu. arbly, no, dirty, but check Ptz! ← I really don't know how this would work! → f.22

SN In QBasic 4.5 there is a Manual on Creating & Maintaining Quick Basic Libraries. ^{Power Basic & Consul Compiler} ^{version 4.0} ^{Home 1.800.780.770} ^{ord.no. 372117} ^{Give Email addr. if they w} ^{Ship it.} ^{Customer 31126}
Read on for Ptz's & Stray Ptz are used in writing new QB Ptzs. or 4.5?
See if E have Ptz in QB 4.0 & see if its in PB.

In our copy version of ^{Visual} Basic 6.0 (2 CDs): Product ID is
82892 - 884 - ~~76~~ 76 991 - 92385 CD set: X04 - 12942 } 1991-1993
= CD key

PP85 In WinXP: to remove password: Control Panel → user Accts → pick task → change an Acct. → select an acct to change. — click Remove (The My Password) or click on one of the Acct. that appears with "Pick an acct. to change" section of the Users Accts window — enter your current password to confirm your identity. Then click "Remove ^{the} password" or "I remove the password".
"Removing a Pwd. Locks protection for the acct., so it is again freely accessible to anyone." ← ditto xp book.
"Removing password is not same as setting Pwd to blank"

2: .09 → It is conceivable that one could keep a Ptz resolu, but diversity would not ↑ — in which case, one would probably be stuck, & the prodn. could not continue!

Int. idea: I had Ptz of early bits in code defining the function table used (f.22) & it was getting stuck w. a particular f(z) housewife — but a truly universal funct should have subsequent bits available to modify f(z) in an arbitrary way!

At present: several apparent difficulties. 11.12-15: (a) Ptz f(z) is "universal end" (b) perhaps that subseq. bits can effectively modify f(z) (.24-26) (c) Ptz f. idea of ↑ resolu. of 11.12-15 will work at all! (d) Ptz Ptz is a way to do Ptz w.o. Ptzs cc (possibly Ptz is way used in previous fractms), (e) Will the system afford of stuck w. not enough diversity to be able to predict or continue coding of corpus. (f) How to adapt Ptz system to all QA corpus (maybe unaccy: Use Ptz system for TM2 only). (g) Any way to try if the corpus w. simple register — (h) Say we have Ptzs 2, bc adjacent: all same "color" = G If we code base 2, then ↑ in resolu. to any amount w/ (1) given source color, G. T. way I

11400 spec.

5-PM

30. (Spec) 113.99 expected to do (-w. this): That I would originally hexadecimal codes for 4. corpus ending in \mathbb{R} . Some would ~~estimate~~ predict a new color at same res. & some would predict \mathbb{R} . Another way to deal w. this diffy is by "BACKTRACKING".

Section
Covered by Pass 2
Miss (Recruit)
"Discover"
Summer June, 05
possibly May, July

0 : T. Sequence of Real (say) pts. would have coords like 1010⁰⁰, 101
11 101.011⁰⁰⁰, 101011^{83bits} (instead of 3 bit chunks, use 28 least 8 or 10012 ^{3bits} ^{space})
over 8 bits to get a larger number so that colors w. small pts could get reasonable estimates. So R. codes of .10-11 look like a sequence

14 of 3 bit "bytes" — so looks legit! — ~~in fact~~ in fact, it looks more legit

15 than using ~~utilizing~~ Backtracking (think!) [One way to expand diversity is to use different rows (cols) of pixels. (≡ Different values for complex syst.)] → 23

Another possibly relevant idea is Chris Moore's demo Prod 2 & simple (1992 or 1994 or 95) ^{Time} ^{or EBC?} System can represent an arbitrary Successive orbits representing (↔ correspondence) to Time Tape Motion.

3: (19) : T. way I'm thinking about it now, there could be no "Backtrack", because all poss. codes are investigated! We average k "Radix" (= 8 = 3 bits in .10-14) so that a reasonable no. of each/symbol ~~is~~ ^{output} ~~is~~ ^{is} produced is adequate to get good precision and small likelihood of "dead ends" (≡ no cases of k needed symbol occurring).

A way to introduce "Backtracking": say k is the no. of symbols in the alphabet of the corpus. R is the radix being used in the coding (R will be an integer power of 2 usually, but R can be k). k is the no. of symbols used to represent ~~the~~ ^{the} complex no., \mathbb{Z} .

~~the~~ R will be (about) k no. of codes per character of the radix R codebook. $\frac{R}{k}$ is the no. of times a ~~particular~~ ^{output} symbol occurs for each input symbol.

If we develop all input symbols, one of k output symbols will not occur a fraction of k times w. $\frac{R}{k}$. — Thus (looking at) — Since we "develop" all poss. input strings, if a ~~particular~~ ^{output} corpus symbol does not occur, we have a "dead end" — we cannot backtrack because all poss. codes have been tried. — A way to deal w. this:

say k is $\frac{R}{k}$ where $\frac{R}{k}$ is the max prob of a symbol.
R = input symbols
 $\frac{R}{k}$ is the output symbols.

STM

10: : Actually, if $\frac{R}{K} > 1$ we will have an exponentially # no of cases of most ^{corpus} out-symbols! So take $\frac{R}{K} \sim 10$ or 100; But develop only a ^{randomly chosen} subset of all codes for f. corpus. The fraction developed can be so that the No of ^{coded} cases of each corpus symbol is ~ 100 or maybe 1000.

So $\frac{R}{K}$ is comfortably large, i.e. we develop only a small fraction of legit codes. Then if we run into a of diversity (or "Dead end" which is an extreme case of low diversity (i.e. of diversity), we can back track by developing more of the codes starting 1 or more symbols back into the past.

So, as is, it would seem that I have an ideal SUMAC system, w. adjustable Backtracking when needed!

T. Frogg suggests that I miss an ordinary Sumac System w. occasional back track, but usually, the diversity is kept by a not so random as needed for back track!

I always plot of ~~the~~ simple Sumac as automatically of diversity — giv. not so!

18 \Rightarrow Each new corpus symbol can have several new codes extending many of the old codes. I really have to look into 18 more carefully: It ~~is~~ may well be that ordinary Sumac is normally quite good!

Fractal prediction contract w. normal Sumac in that fractal compression always has exactly one corpus char for each code char.

Re: Code Fractal codes: to compute $F(z)$ the progressive/increments of z , used derivatives! The increments are small but derivatives are Y.G. (we can get upper bound on error by knowledge of second derivatives).

20! Sumac ~~is~~ is only for the Series predn. Could I adapt it to do TM_c (on QA-TM_c)? Or could I somehow use a Sumac for TM_c(QA)?

Other realizations of Universality CORDIC \rightarrow (17.28) Chris Moore (Az, Lisp), FORT¹¹
 Traces, String matching systems — Eric Post (Wolfram's Cellular Automata) \rightarrow (FE, Dynamic Systems, Mechanisms (RISC))
 at RA predn, ~~the~~ BAC predn, time Series predn, quad 1, 2, Adm. Compression...
 Can we Go on. Cordic to do more Universal Prags? It does generate functions normally deduced by integration (see Chris Moore paper) — so can we use z.r. to include all (or partial recursive) functs?
 It may be able to do all prim. rec. functs.
 \rightarrow Consider just Sat of SONY PlayStation 3 (PS3). Low cost per inst ~~at~~ execution, by speed, ...

F3	320 x 200	256
F4	640 x 480	16
F10	320 x 400	256
<u>SF1</u>	360 x 480	256

SF2	no
3	"
SF5	Synch in morov.
6	"
7	"
8	"

SF9	320 x 200 x 16
SF10	no

0:15.40 : Away to think about ~~the~~ latest SUMAC idea: Say k is radix of corpus, R is radix of Z representation (of $f(z)$). ~~Then~~ $f(z)$ then maps Z space to Z space. If $R = k$ then for one $f(\cdot)$, Z spaces are identical. Optimization $f(z)$ maps Z space to Z space \equiv corpus space. This means that certain corpus strings will have no codes (download). If $R > k$ then we can have download occasionally, but as $\frac{R}{k} \gg 1$, it becomes unlikely... . Particularly affect corpus has many symbols in it. The redundancy of representation is $\left(\frac{R}{k}\right)^L$ where L is the length of f code plus for.

SN fractal for Dos V.9.2 allows zooming directly 10^{1600}

$$1600 \times \frac{\ln 10}{\ln 2} = 5315.08 \text{ Bits log.}$$

So $2^{5315.08} = 10^{1600}$: A question, how many do it, does it slow down much, can't be easily extended?

Actually, this could amount to adequate prodn. of radix & corpus. $\frac{10}{8}$ is only 1.25, but for long corpus, this can give adequate diversity. $1.25^{10} \approx 9.3$ So for corpus of say 100. 10^3 we have diversity factor of 100.

On Dos fractal 19.2 (I think): I kept zooming in on a st. line edge betw. 2 cols but it kept on being a st. line edge: So this can be an imp property of fractals.

CF	10
2	40
3	"
CF4	640 x 200 16
5	320 x 480 x 256 SF1
CF8	hour.
CF10	hour.
CF16	376 x 564 x 256 SF1
CF25	400 x 600 x 256 SF1
ATI	bluetooth 800x600x256 no
video	7 no

I downloaded fract 200 zip. This is a dos 2.0: In general, 4.57 k by/sec ≈ 26.56 k bits/sec. download by phone: Not 825!

NB On forpg. SUMAC: It uses not ~~ASCII~~ to "continuous" but to "discrete" unvli DF. —

i.e. each code tells Machine when to stop. So its noisy to use Cover's "Extension Complexity" which ~~you see~~ long to compute. I'm not so sure: I reference machine ~~not~~ necessarily have one way output to get ALL.

SN Criticism of Fractal prediction: The chars in corpus are not all created symmetrically at all! A priori, Permutator of corpus symbol names should leave f a priori d.f. invariant — but in fractal prodn, f chars assoc. w. small "Bifurcation numbers" are quite different from those w. large "Bifurcation nos.". Nvr, if one assigns f corpus symbols to Permutator nos. on a basis of observed past frequency — this objection is less strong. This objection has not been eliminated!

36 One way to get symmetric corpus symbol assignments, The first orbit to "breakback" will do so at a certain angle: We can assign corpus symbols to angle intervals of $\frac{2\pi}{S}$. Certain $f(z)$ functions could give some bias on Permutator angles, hvr. $\frac{S}{2\pi}$ (17.10)

STM

10. (116.36) My former objection to Source was (partly) just "errors" (omissions) early in corpus code, could never be corrected - except by a fortuitous ~~Architecture~~: Probably not so. - a good Source should be able to compensate for early errors - but a suitable tsp. is perhaps much search as is needed. This is (partly) the idea of being effectively able to modify early codes (or, the effects of early codes on present problems which is all that's really necessary).

2. (116.10) (SN) Some random Notes on fractal prodn: Using 2 Circular topology of output
 22 in 116.37 ff: We can easily divide up the output in complex plane by comparing its real and im. parts: $re > 0$ gives bit $im > 0$ gives 1 bit. ($re > 1$ gives 1 bit.)

- 1) C. Moore's set of Analog f(z)'s
- 2) D.F. in 0.

(SN2) A Continuous D.F. on finite circle is a 4th d.f. to add to $\pm \cos(0, \pi) 0, \pi$
 i.e. $0, 2\pi$. The main is a single "phase" of. - for any broad d.f.'s function
 $2(\cos(x+\theta) + 1)$ normalized will do it. Its peaks at 0 , (comp. at $\theta + \pi$).
 Now its variance is large. It can't have any small var.

Putting $\frac{x^2}{2}$ on the unit of world do it. - whether it would be actually
 manipulated Mathematically ... I don't know. Just expand $\frac{x^2}{2}$ in a f series.
 Very probably all the integrals are known: $\int_0^{\frac{x^2}{2}} \cos^2 kx dx = ?$ - is a function of $\frac{x}{2}$ probably
 maybe mult by n^2 or n (or none @)

23 (SN3) We can use C. Moore's (Shannon's) set of universal analog. functions for $f(z)$ in fractal
 Source. Perhaps D. code would start w. a short binary string that describes the discrete aspects
 of f. function, followed by 1 (or more) continuous parameters. Look at C. Moore's Paper in "P's" file.
 I think I have a first law P.P. in "Hard Copy" (B) If a "Integrator" function is expensive it
 done digitally.

119.20

28: 115.315 Perhaps try to get more General functions via Cordic.

Consider the umc-types of 115.32 & others: One reason I liked - Lexp was that 'regys' were easy to spot (Common Subtrees). So TM2's job was easier. However in the latest ideas about SUMAC, I didn't consider using a TM2 at all! So, say I used Machine (exp. is no TM2. How good could the system be? In General, I've been thinking about SUMAC's w.o. considering that. However, it would perhaps be a ^{kind of} ~~advantage~~ ^{for TM2} ~~compulsion~~ from the one I've been considering for QAM (?). Maybe not so distant! In both we have codes & lengths of codes, & compressing achieved for each of the codes.

If TM2 is not used/considered, then R-SC Machine code would probably be best, perhaps I could start out very slowly w. no TM2, then use TM1 for TM2 when the TSQ

5 TM

105

Was 2nd order. This may take long time! Also if the primary primitive lang. does not have easily recognizable regys in its codes, this $TM_2 = TM_1$ will still have 2 hard time!

Actually, the perhaps the best way to do: For each lang, try to see what kinds of regys occur in it. Existence of a detection of sub pms would seem to be major hurdle.

Perhaps just write a bunch of pms in each lang. & think of common truths in each.

See how strns are used/detectable

Here, for $SUMAC$ as if the brain thinking & brain it, we have a simple sequence corpus, ~~is~~ a set of trees that are codes for it. For a corpus that is many strings of symbols, the mechanics of $SUMAC$ has to be revised ... it applies to all!

On Parallel interaction TM.

Its easy to Parallelize TM_1 ^{secondary} ~~very little~~ cross communication needed. But for TM_2 (Secondary induction), all of ~~the~~ ^{secondary} corpus of TM_1 's are needed. These can be updated slowly, here. $T(\text{secondary corpus})$ is quite large - many parallel codes for primary corpus. Here, when TM_2 finds new regys they can be given to TM_1 rather slowly - as these regys are discovered slowly.

10

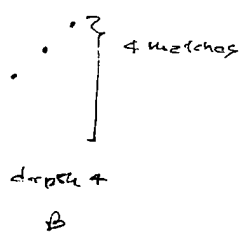
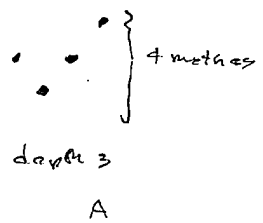
Context tree strct (for TM_2): Say we only have binary branches in the context tree. For depth 4 dis count, we have $16 \in 2^4$ separate 'PM corpi' (\equiv BTREE corpi). One for each path for distance 4. So for a particular context to be matched, we look at all 16 (depth 4) corpi. There will be much overlaps, of course. Could we simply pool all case counts of depth 4 matches? For each depth we will know how many matches occurred, i.e. d.f. of predicted next symbol. From this data, we can make a p.d. over predicted symbols. How good would it be?

27

Note that .20-.29 seems quite different from what I originally had in mind.

10

Formerly a "depth 4 match" was ~~an~~ ^{subtree} a match to the "query tree" in which ~~there~~ ^{matching} were n symbols. But were ~~each~~ ^{connected} path of matches. So: A ~~sub~~ ^{sub} B ~~was~~ ^{was} 4 old style matches; new style: A



- A has 1 depth 3 match
- 2 " 2 " 3
- 1 " 1 "
- B has 1 depth 4
- 1 " 3
- 1 " 2
- 1 " 1

37

10 : The "16 corp" could find all matches of depth $[z]_0^4$ individually, The optimum method of pooling to date, is unclear hwr.

My Test Bias is about 1/2 (on 118.37), A is "Better than B" - That more ~~helpful~~ compact context, small depth, but ^{many} contig. symbols is important.

11 : There were 2 poss. "objections" to Sumac: 1) That I'd not bother to women directly; the back tracking could do this 2) How to do TM2 was a major problem - what corps in 1. TM, code is how to recognize Recm.

To some extent, 1) is an empirical Q to be found by Expt.

2) Is something Empirical, but ~~was~~ ^{much} theoretical analysis to decide on what ^{2 particular} primatives Lang. is good / "adequate" / fast / cheap - \rightarrow 150.10 ft 151.00 for processor & Sumac

12 : 17.27: Dynamic Systems vs Times See Chris Moore Paper in 1996: "Actually, I Rely on Dynamic System" Any kind of continuous mapping from $t \rightarrow t+1$ to $t \rightarrow t+1$ would work for his analysis. This impt.

Q is just how he did for mapping from continuous funds to Time. States a tape storage excess.

Try similar maps using easily calculable Continuous func. Find out where "output tape" is.

Can I derive U10 mechan or equiv? The three ways to do induction on regular (non-U10) in Times. - But where is f. output" is an impt. Q.

36 To what extent is it practical pixels color & possibly useful output for a universal black? \rightarrow 12.20

13 : On Sumac (w.o. TM2) is "Phase I" w/ TM2: Why does TM2 seem help some much? Sumac, as is, is universal, but it may need backtrace (20 ways to effectively modify early code.)

32 In TM2, we put a rd on codes for TM1. This is effectively putting the output of a stack grammar (in general a Umc output or "random" input) as input to the "final" Umc.
33 So why is this TM1, TM2 combine better than the "simple" Sumac? \rightarrow TM1 (to start) is Sumac (20.09) for good UNDERstanding of bits !
good at finding a standard set of common vars in the corpus. It effectively removes them, so TM1 can concentrate on finding new kinds of variables. Once a new kind of variable is discovered, it should be incorporated into TM2's Stack Grammar!

It would be nice, if, indeed, we could incorporate all or most newly discovered
 rary ~~types~~ types into TM₂'s "Stitch Grammar". T. rarys in TM₂'s S. Gram
 would have to be detected by "TM₃" (\equiv TM₁, eventually/usually/always). "TM₃" does
 come into play until the corpus of codes for the primary corpus is large enough
 TM₃ has had a proper/adequate TSG ~~to require~~ to make ^{feasible} likely, the kind
 of rary recognition we need.

Still, I want to understand clearly, advantages of a TM_{1,2,3...} device v.s. a
 simpler "1 level" Sumac.

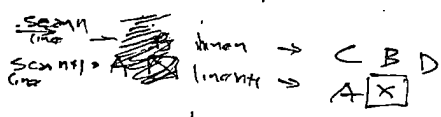
AHH!

T. TM_{1,2,3...} is well understood as a variety of "Coding & Recoding" —
 an ordinary, Lagit method of coding, equiv. to a "well coded" Sumac. In fact, this
 view, gives us a much better understanding of "TM₂'s corpus" in QATM!
 It enables us to more clearly define just what are "Lagit-Corpi" for TM₂
 in QATM!! @.

See if I can unify t. view of .09 w. ~~with~~ t. view of 119.32-33, as TM₂

providing a "preprocessing" of the "Random" inputs to TM₁'s "Reference Unit".
 6.16.05 → WOOPS! Not so fast! .09 maybe OK for BU such, but its not clearly relevant to
 Ls vch. In Ls vch we do ~~not~~ ^{use} ~~order~~ ^{order} coding in PC order, so "recoding" ^{couldn't}
 ordinarily ~~be~~ code length. T. only possy is that somehow, for recoding/ ^{by TM₂} USES \rightarrow 126.10

Note on Image Compression: One simple compression scheme: Code for distance betw.
 successive scan lines (non-lossy). Each pixel has a context of 1. Actually, this may not
 be any better than just coding ~~each~~ ^{each} pixel in each line as distance betw. that
 pixel & the previous pixel in the line. Each pixel has ~~at~~ 2 closest fixle contexts \rightarrow



X has ~~the~~ A & B as same distance \in C ^{to D} not much for Rec.
 X has ~~the~~ ~~same~~ closest context of \approx (A C B D). But how
 do we use this for coding? [X is closer to $\frac{A+B}{2}$ than to either A or B.]

Maybe try $\frac{A+B}{2} + \frac{C+B}{2}$ ~~is~~ \rightarrow a guess for X. It would be nice if we could store context
 in code order, i.e. in PPM. My impressn. is that the problem re. storing, retrieving

Subtract context. Look at more successful image coding schemes: Can they
 be modified to help code, variable sub-string context

In the case of images, we will be able to store the context ACBD's ^{corresponding} to p.d. off.

We could do a Lempel Ziv analog: look at previous cases, to forget
 common context in the past. Context is useful. My impression is that the simple
 dimension quality of text dit makes L-Z coding feasible. — A 2 dim analog...? How?

5TM

What about 1-dim Lossy coding? (I think lossy is easily modified to non-lossy, by adding code for amt. of error.)

One v.g. lossy code is MP3 for audio: v.g. select what info to discard, for 2-dim images, ~~FFT's~~ are used & (or Wavelets plus?). Info discarded is in higher frequency components: So loss at variance

I think there are v.g. lossy codes for movies 2/0 ≥ d. images. Also latest JPEG may be slightly better than normal JPEG.

Look at Dondari's 2000, 2003 papers comparing JPEG compression of ~~images~~.

Images w. LZ of images. Note this may be old JPEG - Not v.g. (?) - More recent

JPEG may be Much Better!

119.36: Dynamic Systems as UMS : New idea: I was concerned that "Real World" dynamic

Systems would always have noise "from External World" in by order bits of Mem - so t. Computer would have an effective upper limit on precision & RAM capacity per System param (≡ degree of freedom).

For prodn via a UMC w. random input, we can have 2 kinds of input for v.unc. Input by USR: Random input by Outside world. To do prodn, so t. USR wants to find many also short codes. (The noise from outside world is not purely random - it does have "repeats", but these repeats simply ↓ amt of info in that input, so we need more of it to qualify for "truly" random input.

How it works: USR finds (usr) inputs that give copies upto n bits. He then reports these inputs to get p.d. on next date. Each input may have to be repeated several times in order to get its bit for no x bit because of "outside world noise".

Certain ~~input~~ input codes will be less noisy than others. (I suspect short codes will be less noisy). Certain long codes will be so noisy that they will be effectively useless?

Also certain long codes will not be able to code. Known corpus w. much repeatability.

I imagine that 2 repeatable codes for corpus would tend to have repeatable codings.

In general, I think that prodn w. this system will not be difficult, but we

5.TM

may be limited in the lengths of codes that are usable — so corpus size will be limited. We probably can deal w. this by using many Dynamic Systems m.l. (2.08)

To do prodn, we can just run in many codes and if they code to corpus, continue w. those codes. For codes that code to known corpus, we will repeat them several times to find d.f.s on their continuings. We start w. search in (length) of codes. When we find a code for known corpus, we can use it for prodn, & repeat it for prodn to get p.d. for prodn.

Just how we use many ll systems to deal w. limited capacity of each Dynamic System — is unclear. How do we produce a very large corpus? How do we do larger T&Q's?

Note that 4. individual (micro) computers need not approach universality. They can be small parts of a larger machine that does approach universality. (w. limitations in addition to small memory)

One (sort of trivial) way to ~~not~~ combine ll machines: Each one is essentially

a highly variable sub-machine that uses only a few bits of each dynamic variable (degrees of freedom) — like a sub-memory or part of an Aritk unit. They are then connected in

"usual" computer ways. So the memories can add up to very large sizes.

How communication betw ll machines will be relatively expensive, so we want to reduce res need as much as possl.

Soln of QA probs; a seq. of induction probs by GP! → Note "objectives": 123.11

One of Dittys was need to verify each new cand. out. entire corpus of problems.

This was partly solved by ordering probs by difficulty & trying hardest first.

An imp. idea that I may have missed in ~~the~~ work: Dittys GP is ~~being~~

being used to solve BUSch problems in induction. As such, we always

start out w. a default code that the prog tries to improve by hillclimbing.

So every cand will have at least a pc & p code for entire corpus. CONTRAST w. LSch, in which PC = 0 until 1. first code is found — which is quite interesting.

A cand will not be tried for entire corpus, unless it does well for hardest probs thus far. Also, for ~~perhaps~~ almost all of the "very easy" part of the corpus, we don't have to test

to cand on all cases. — Maybe just some representative cases.

So, using a good T&Q is GP and CFGD for search, one might be able to do fairly well. This is close to what I had in mind for QATM,

(phases) but ~~there~~, I used LSch T&Q's, a TM₂ to discuss a grammar ~~to~~ get a pd for TM₁'s such. How, since LSch probs are already in a CFG — it's not clear as to what Grammar discovery would do. Certainly good pc's is our problem to be

word on! also Context ^{Sensitivity} Search is definitions of ~~...~~ sub-trees.

Just how sub-tree defn. would occur is unclear.

Also, perhaps the $\frac{\sigma}{Max-M}$ criterion for search ~~...~~ sub corpus.

As each case comes in, I get its G_{New} I insert it into a file w. G ordering. As I do so, I ~~...~~ keep track of G & G^2 : I can then get changes in M & Z for any subset of Cands w. $G_i \leftarrow G_{New}$ ~~...~~ so I can tell if $\frac{\sigma}{Max-M}$ has changed.

I would do this $\frac{\sigma}{Max-M}$ calcn only infrequently

I. $\frac{\sigma}{Max-M}$ criterion was originally designed for a PPA corpus. Using it on a CFG or CSG corpus may not work well.

I. discussion of 122.20 ff is a bit deceptive. The problems are not like each problem is an individual to be solved. It's not like finding a short code for each sub corpus.

The 2 defn problem units to be put in this form: we have an overall O^2 (for QATM)

The O^2 is a pd on symbols & macros (= tokens) a series of symbols a ~~...~~ sub trees.

From each O^2 we have to find a ~~...~~ short functions from $(Q_i \rightarrow A_j)_{j=1}^n$

So O^2 is a def over/function space, & the prob is to find function ^{of the PC} on this space that map from Q_i to A_j . To do a prob ~~...~~ from Q_{int} we just get pd. on functions,

(which is O^2) & go ~~...~~ f_{int} 's. - Actually prob is the "normal" QAT soln: O^2

is a pd ^{on} to set of functs from Q to A .

"derived numbers"?
"I forget just what the form is."
each string in derivation is an "acceptable function".

(NB) The pd of int is not a "developable lang" in which every ~~...~~ string of symbols is a function. It could well be the certain short strings of functs

have ~~...~~ max pe, but local shorter strings do not exist, & CFG can be of this form.

The p.d. of interest into be a finite set of functions: but such soln would not be Universal: i.e. it would not assign \Rightarrow $pe > \phi$ to all poss. Any - So we do want

a grammar of an infinite lang - so it's a CFG, it must have at least 1 loop -> see 124.06

An imp idea of 122.22 ff is 123.11 ff is that Boosting ^{for \forall such ϵ} may be relevant to these problems. - But normal Boosting Doesn't seem relevant to Lsrch.

A mod of the Boosting idea related to Lsrch in QATM: (perhaps).

In order to apply Boosting to an Lsrch: way we apply it is to modify the d.f. on f. "guiding P.D." that is used for Lsrch.

Perhaps certain functions are used in solns to certain QATs that get low pe's for their solns.

We have modify pe's of a primitive set of elts, so that the pe of f. ^{active} corpus is max. - How this is not closely related to "Boosting".

-> 32-36 sounds very much like Lopp's rule for a BIVT seq!

Lo temp G-luc today notes
on "bu Subgroup":
(for Lucy)



in 122.20ff One way to ~~complete~~ test a D.F. on functions, wrt a set of QA problems.
List the functions in PC order: try each functional all QAs in the corpus, & ~~write~~ list
which ones it works for. After envt functs have been tried so that all QAs have at least 1 found
that works it, we can then add up the p's of each relevant ~~envt~~ funct for each QA, to get ~~the~~
pc for that QA. Then get product of all pc's for all QAs to get ~~the~~
of corpus wrt. that D.F. on functions

124.29 I wrote about this particular model of QA induction before: Much criticism:

I don't remember if it worked ~~at all~~ correctly or not - i.e. was it universal?

If the p.d. on functs gave > 0 pc to every pos. function, then the Rank 1 system is

09 Universal. Actually, it corresponds to a symptom! A Big Diffy is that it's not normalized, i.e. perhaps very expensive to normalize!
? \rightarrow What was not universal was a "Search technique to find the D.F. on functions" \leftarrow (???) \downarrow .12
- I wanted a universal d.f. over mathematical p.d.'s on function spaces.

12 .09 halls may be harsher if not so diff! for 3 input ~~issue~~: T. "R" input need not be
interpreted as $pc = 2^{-|R|}$. I can use any interpretation I like, i get them to decr \leq to ∞ .
I can get them to be small measures, because certain R values will give no output for certain

15 Q inputs. I could just use the semi-measure product to get a "T" p.d., ~~but~~ i
Select functions to max R_i is value (for given C & B \in context). This would bias any
functional forms toward forms that usually summed to near 1! try 25

17 So: Anyway: to normalize, we just empirically find a normal constant for each
18 "3 input ~~case~~" — Or say R_i over the legal (integer) values of R_i . (R_i has an automatic
20 "end of string" marker) we use $f(R_i) \rightarrow \sum_{i=0}^{\infty} f(i) = A < \infty$ then A^{-1} is the
norm constant. This assumes all (R_i, Q_j) converge .151. ($f(i) = 2^{-i}$ for sum ∞) was our guess

Q: is there any particular case that would give very bad results (.151 \rightarrow .17)?
Q: In my earlier work with veld in .06: Did I really have any serious objection
24 to (.23.20) (i.e. symptom)?
25 .17 \rightarrow Also, remember the result on the "continuous" Univ. d.f. — that the ~~the~~ semi-measure
approached a measure rapidly as $\epsilon \rightarrow \infty$ to probability of "0" $\rightarrow 0$ rapidly
for large corpus: So maybe Normal is not such a big problem.

So I now have at least 22 ~~the~~ not very difficult solutions to the p.d. on
strings A in QATM: OT "vacant" solu of 123.20-29; 124.06 ff \approx

31 p.d. over $\{0,1\}$, end of string symbol as a function of $\{Q_j$ and the first k bits of $A_j\}$
On second Q: This d.f. over 3 symbols... I don't think I solved it! I guess I could not prove it
wrt. to d dim d.f. on $\{0,1\}$ to a 2 dim d.f. on $(0,1, stop)$ — but it's not so hard —
particularly since the stop symbol has different properties from 0,1.

Urr. in .18 (.24R) I discuss normalization of the "A" d.f. and it does give a way to
assign pc. to the "stop" symbol — but this is the prob. of "stop" in R , not in A ...
So the 2 "sols" (for 1 v.s. 2) (.30) (.31) are not the same!

Wrt Normal. see note on (.30R). V.G. Looks like it will not be a problem
when we get down to lay in the TSC (a q. to BU such)

Q: is there any particular case that would give very bad results (.151 \rightarrow .17)?
A: 2^{-i} "approx 2"
A: 2^{-i} E.
norm.
 $\approx 2^{-i}$ "correct value of 2^{-i} "
How I'm not sure that proof applies in present case! —
It does apply but I wrt. proof of convergence of QATM not proof of Sol 28 Norm 3