

special option He puts for induction: some portion of from Seq.

Doing Phase 2 first (perhaps never do Phase 1!)

0:324.40 : This list is a kind of "hill climbing" procedure to induction mobility (= comparison). Note that 324.38-40 is quite comparable w. 324.31-37: These Opten Methods can recognize induction problems & kinds of induction problems as part of a larger "context" ("correctness" context) of these Problems.

PSM'S

t: 323.19: Collin: (cont)! Phase 2 has (expandable) set of (prob Solving Methods).

35 An alternative Plan is (324.31-37) to do Phase 2 to start off! - i.e. no phase 1 at all!

029-04 } So: List various approaches w/ little detail; Run go back & forth more info: (possibly more "BRANCHES")

07 SN Consider "TRUE Lsrch" for (d/s)-induction in "Phase 1": This would essentially be same as starting in Phase 2 (324.31-37). This "Covers" all possl. ways that I could try to do Phase 1! E. Hvr, while starting off in "Phase 2" does look promising, it is a new direction, & I want to list the normal, "regular" ways first. IAP talk. Actually, considering many good ways to do Phase 1 is a good initial approach to "Initial Phase 2",

My impression: This one of major interesting problems in "Phase 1" is doing TSP's ... easier for d-induction; harder: more interesting for s-induction

Perhaps I've been thinking about s-induction in a too abstract way! Considering specific problems would probably suggest specific Solns. & some of these specific Solns.

0 SN d-induction & INV probs can be of any diffy: Can some/all [s-ind] probs be expressed as [d-induction] of INV probs? T. Some representative seems to be d-thus - which

1.2.05 } is one of fastest pol claims. (i.e. after is Universality) One claim is that the shortest R code for a Ai is usually not very different in pc from the A1 pc of Met Ai'. So s-induction is a lot like d-induction.

We might say that in d-induction, we need a super OJ w. min bc that works for all Qd's. In s-induction, we want a OJ & its bc plus sum of its "R" bc's for Qd's is min.

For this s-induction form, the search routine has to be worked out. We can get suggestions from the d-induction routine. We should consider both regular Lsrch TSP's & T.

13 "Problem Pool" TSP: Also consider BUSches w/ the Lsrch

31 So first outline d-induction for various TSP & such types: T. Corpus for OJ induction (APPMS) can be the set of successful OJ's off pack, weighted by the no. of probs solved by each (i.e. "total importance" of the problems solved).

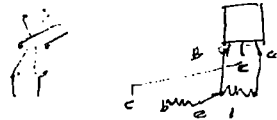
32 In simple Lsrch, we just hunt for a OJ that will solve all probs up to Qd's, using the corpus of 32.

As each problem comes in, we do a new search w. a slightly different OJ corpus. I think the s-ind Lsrch corresponding to 30 off: T. program is, how to weight the corpus of past OJ's.

The few that have been needed to solve the set up to now (= "n" probs) we could try various past OJ's on complete QA corpus up to n, to give them wts. I imagine that most of them would have very small wts, since they couldn't work many of the problems w.

any reasonable pc.!

ATM



"MCT"

Spec
 10: 223.40: The basis of f. pd in Normal OZ Lsrch, is not clear. T. great breakthrough of t. Mixed Corpus ^{from} "PBM"
 Was not we would do well by trying to find a d.f. for (one method, problem pairs): giving t.
 P.D. of solns (cc or "t" was included in t. problem dcm). Next breakthrough was t. WON soln

T. Simple form.

In Idsiz report: ~~the~~ idea of Phase 1 was to be good enough s-induction to
 Go to Phase 2. This ^{early was to} Goal ~~was~~ be obtained by any means possl. - i.e. not need to
 stick to Lsrch. ~~the~~ I think I considered (unusually?) ~~the~~ "Lsrch" over part of a s-induct
 dcm, & something "n" like "Edit dist" for t. rest of f.p.c evaln. "Augmented G.P." using "PPM" was
 another direction

Perhaps my "error" in ~ 30M (37, 138) ^{→ 319.00-323} was not to realize that Lsrch was not
 to be applied "exactly" to finding s-induction solns! - That, in general, searching for
 solns to s-induction problems was diff., & perhaps had few features to make it easier
 than general OZ problems. ⁽¹⁵⁵⁾ ONE myt feature: If solns of by pc exist for
 entire TSO, then listing find could O's in pc order is an acceptable technique.

It becomes an INV problem (But INV probs aren't Lsrch in that way either!).

Consider .155-16... is it true? What is expected cc. of soln? Consider O's to be a
 simple. Consider Phase 1 s-ind. tsqs. If we do a complete backtrck ^{from} ~~at~~
 t. nB QA, we search for O's in t. n R ^{level's} - which can be a very small pc!

Hvr, what is t. actual cc. of finding soln?

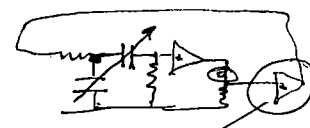
I also considered PPM for t. "Adaptive" part of adaptive Lsrch - for both INV &
 END problems. At this point t. main difference betw Lsrch & G.P. was that G.P. used Bsrch
 & Lsrch used Lsrch. G.P. was thorough; Lsrch was not (well, Lsrch could be Mf Carlo, but when
 it is, it is quite different from Mf Carlo in G.P.). Both are guided by a P.D. that uses P.M.
 → "Lsrch" lists O's in pc order (pc's & O's) & t.ys them in corpus.
 → G.P. such creates a subcorpus of O's w. $ES > X_0$ ($X_0 ES \geq \frac{O_k}{sum - M_k}$ is Max). Then Mf Carlo
 selection from that subcorpus. Corpus is continually (or periodically) updated w/ PPM & X_0 .
 Lsrch works well if PPM is a good O's of very hypc that solves t. TSO.

31 [SN] I had been thinking of doing Phase 2, which Phase 1 had ~~not done well~~ done well
 on problems that Phase 1's matrix controller ("TM") would find familiar. i.e. when TM was
 up & working on ~~the~~ Phase 2's Meta problem.

Would it be possl. to start TM out by having it inn params of a program that generated
 a bunch of OZ solvers that were known to be v. on? (I would give these solvers to run
 in "factored" form. [This is a REALLY SEPARATE APPROACH to TM] ^{One B&C! How much "Lateral Reasoning"} would be needed for "Phase 2" problems,

33 ~~the~~ [EN] Int. forging pp. I haven't really mentioned some ~~very~~ option methods that are very
~~special~~ special for induction. ¹ finding ~~the~~ contacts that have start dcm's is a good s-induction.
 40 ² looking for ways to ~~run~~ a corpus, then recursively (Lsrch replacing it) until the new corpus ^{is "Coding" & "re-coding"}

4 TM



or 20
compression
of (lower)

blac
322.40

Must ask what is pc. of ti corpus is, if w/out a 2 certain $m \geq$ a certain width of def. We want to
 make pc of "code" of
 corpus. — So the big Q is: Do we get more compression than w. straight PDM?

Remember we must be codingonly useful Lisp-like prms. This method may work for prms like Proc,
 but not for English Text.

There is another possy!: That it may not compress (is, well, but that it is a good way
 to generate code for by $\frac{G}{G_{max}} \rightarrow G \rightarrow \epsilon \mu_2$. So this should perhaps be criteria for

choosing m & distrib. width as a function of m . If the criterion is, indeed identical to Max corpus compression —
 Then fine!

So: There are 2 directions of improvement of t. prgm. Trick(s):

- 1) Elaboration of "Context" to be sub-terms rather than just RPN or PN strings of symbols.
 This is not easy (or apparently) not cheap.
- 2) Using pattern sub-terms & entire branches or just a few symbols of sub-terms. — ~~to be used~~
 for "production" (or more exactly Monte Carlo cond generation). — Also what assoc ~~and~~ Probab distrib
 "width" to use. The concept of "width" may have to be modified if we use contexts
 that are not simple strings.

OK: List various paths to TM: Give diffys mach, refs to attempts to fix them.

Plan: Phase 1 is d-induction then s-induction: This is a sub-problem to enable phase 2, which removes these probs for phase 2 (is phase 2 325.04)

1) Lisp on 2 d-ind TSQ in Main.
 The (Recn) Funct may not be needed. Checking on Post Corps comb: Harder probs first.
~~extensive external data synch. This is a major problem!~~
 Writing usable TSQ's may be major problem. — But maybe easy if large TS is allowed.
 Lisp can be used (w. AZ).
 Possy of TSQ based on Lisp to do part of what Mcler / Wilhousier do (But note)

325.04-05
 322.320.28
 It may be that Lisp
 on d-ind for Main TSQ
 will not be much good
 w/o ability to do Logic2
 Reasoning
 Also Note 320.29-30

2) "ATSQ for s-functs" is a major jump from 1: Some diffys:

- a) choice of s-funct. representations: (Probably use several).
- b) How to use d-function capabilities to help w. s-funct lang. (i.e. "translation" from d-to-s lang).

One way is "EDIT Distance". Its not "Univ" lrr.

Also note: d-induction is not INV Lrr. s-ind is a OZ Lrr, which is really diffnt!

On Lrr OZ problem, one searches over Optim. Methods, v.200 from trial solms. — Method my recent
 discussing of Lrr for OZ problems is wrong! (But not really irrelevant!)

Anyway: This is a true Lrr over Optim. routines / each search strictly Phase 2 problem!
 My/discussions of GP for TSQ's has a little to do w. Lrr. It is a search method related to Lrr.

In a list of trials of Optim. Methods: "Lrr" would be one simple method.
 Perhaps GP would be another, ANN another. — But there are really an enormous no. of very heterogeneous
 Optim. techniques.

OZ

Actually, Phase 2 doesn't use Lrr. Lrr simply lists OZ methods in their PC order — The arguments
 to read it consider.

4 TM

0: 321.40 Way to do 321.36 ff: From M, choose random pts, not nearly on defunct corpi. Keep list of contexts found at each point. These contexts can be of varied complexity; they can be simply ~~suffixes~~ suffixes (as in PPM) or can contain branches, subtree contexts of various sizes. Any way, when one has generated lots of contexts this way - find way to tell if 2 contexts are "the same" or "the same to depth d" (d=1,2,3,...∞). Then exchange branches at the 2 pts w. same contexts.

8 The PPM version of 0.00 ff: So we are using PPM to generate probability counts w. assoc. PCs at a certain point, we have a known set of contexts (suffixes). For each suffix, we have a pd over t, $\frac{\text{sum of}}{\text{sum of}}$ that will follow. Here, by looking at the various kinds that have been in context, one can also get a pd over t. ^{compiler} branches that will follow. Well "pd" is perhaps not so appropriate - since each branch will ^{usually} occur but once. However, each cand w. a particular context, will have a complete branch at that point - so one could just choose from at random w. a likely hood. Having done so, we may have finished our cand down, or not. If we are not finished we simply continue

15 as for the branch was selected. Selection of a branch amounts to selection of a Terminal Symbol. ~~to~~ to be determined in (08-15): How much context to use? One way: In PPM one can end up w. a specific context to use: (partly on basis of consistency of products). One can then use that ~~same~~ context as simply randomly pick a branch from ^{past} cand pts w. that context. This is by no means the only poss. way to choose, however.

2 How frequently to chose branches in this way? One extreme is to ~~choose~~ make all choices this way. Another poss. is to make the ~~frequency~~ decision random w. a pc to be adjusted so as to maximize $\left[\frac{G}{G_{max} - G} \right]$ of generated cands. How to do this is unclear!

0: 0.26.04 A trick related to (08-15): In PPM we have this set of strings in a (physical) ^{suffix} "context" ^{order} "postfix" order. For a "normal" McCarlo PPM "pass", we generate a cand, we go to the point in the tree where the largest prefix common to the cand & the "context" is. Then using a ^{prob} distribution of contexts, we chose a point (corresponding to a pt. into active can pt) up or down stream, & we use the ^{single} character following the context at the corpus member at that pt. Instead of selecting a single character, we could (as 08 ff), select the ^{context} branch of the cand, using a certain pd, width. Or, we could select a seq. of m characters, & use a distribution with the prob as a f. of m, so for larger m, we use a narrower distribution.

Remember that the goal of the McCarlo technique is nominally production - so we are ^{SPR} 323.00

4TM

31.29 - 40
31.30 - ...

?? - It seems (most) always get a context
by using a shorter context.

0: gets "off track" it becomes very likely that it will have to use the contexts of cands of lower PC than O_j .
 1: Since O_j will not have these contexts. So the usual cands will tend to follow O_j exactly, until it "misses" a symbol - ... & then begins to use PC's from the rest of the copy of cands
 2: (Note STM.3.00 to 1!)

→ STM3.00

0: On PPM ... or S-gram based GP, v.s. regular G.P.: The usual argument for S-gram based is supposed to be v.s. G.P., is that while GP crossovers transport useful subtrees to offspring, ~~they~~ just crossovers just as often, destroy useful trees. I don't see this: It might be true of one always destroy parents: It might be done in a tournament selection scheme, in which the 2 best cands become parents & are replaced by their 2 children: it would be better to replace 2 worst cands by children of 2 ^{Best} cands. (which may be done).

In general, if 2 parents are much good, they should not be discarded!
 Of course, my method of using PPM does generate a S-grammar. There is still much Q as to whether a good (multitrees) algorithm in G.P. may not be better than PPM in certain important ways. Hvr. while GP can generate some parents useful subtrees in cands, it cannot use them in context at all. Also, GP does appear to generate excessively long cands (??) ← Apparently not!

3: If M is the average length of a ^{branch} subtree (averaged over all entire trees of cands), then replacing a ^{for random pt. in entire cands} branch of mean length by a branch of mean length, leaves mean length expected branch length invariant.

5: Hvr, one normally restricts the choice of parents branch: The usual way is to use a subsample (say cands of $G \rightarrow X_0$) ... in which case the opt of 23-25 applies & mean branch length is invariant.

0: Actually in G.P. choosing a branch at random may be regarded as "sort of" choosing a smaller branch only basis of same context! i.e.:



Choosing branch ABCDEF w.o. context is ~~not~~ equivalent to choosing branch BCDEF with context "A".
 But still PPM doesn't work that way! A context occurs first then
 i.e. pt. on what will follow is obtained.

36: A way to get context in GP that suggests a better way to use context of PPM 1
 To do crossover with any subpopulation, M . Pick a cand from M (at random say). Pick a point on C_0 .
 Look at context of pt. on C_0 . Pick random pts on random members of P , until one has same context as chosen point on C_0 . When the same context is found, exchange branches betw.
 i.e. exchange cands w. pts of same context. This produces 2 offspring. → Among other loc

→ 322.00

4TM

(Spec 0:319.40) : In Lsach is m Searcher Solus (both d & S induction): 2 kinds of Diff'ties: (1) Finding ind functs that will satisfy all of Cor ^{in "Problem Pool"} as much as poss. (of) f. corpus as poss. (2) Not ^(wanting to generate) having to check a new cand on f. entire Corpus, ~~for each~~ for each Cand. 1984

[SN] I don't seem to have a clear picture of f. state of affairs! - Reminds me of my state at f. start of my "Month of ~~the~~ IDSA"! Recently, it seems quite clear i f had started ppy. & PPM. That's a nice Detail that i think I understand somewhat, but the rest seems to be a muddle!

How to Improve GP ~~seems~~ clear, but getting GP to do TSQ's is still unclear!

BU Search is done by GP. What i have is a u.f. GP, but getting it to do TSQ's is not clear.

FOR a problem ~~space~~ space for which a TSQ is available, it would seem ~~that~~ that

Lsach should be vastly superior to BUSach....

So, ^{The Argument} ^{f. argument is!}

Say one has solved upto QAn int. tsq: It should be not so difficult ^{to reform all f. info & use it to search for soln to f. TSQ augmented by QAn}

~~It would seem to be true~~ it would seem to be true indep of whether one uses Lsach or BUSach. Also, ^{the Arg't} indep of whether one uses d or S induction

I. main thing that it doesn't go along w. is ^(perhaps) f. "Problem Pool" approach - i even think it can be modified to apply. - T. Problem Pool is as close as we get to a TSQ, in v. "post

TSQ" part of TM's Education. Its close to f. Activity of f. "Free Genera Section".

[SN] Re: RPN vs. PN notation. It would seem that we would transfer by reading code in ~~reverse~~ order.... But's all. RPN has "advantage" for cond functs, in that one can start of w. Qi as one of f. args. This is not so easily facilitated in PN. The one can simply put Qi in as one of f. "primitives" - as I did in f. Analysis of ~~the~~ AZ: which used PN. Superficially, it would seem that ^{f. cons.} either ~~then~~ because of context considerations PC's using PN & RPN would be identical -

But if context is considered, they could be very different!

[SN] On Logical Reasoning: A real, reasonable motivation for my not considering other my models of induction (yet), is that others have done lots of work on this. - I mainly work in areas in which I feel that insufficient work has been done. → 323.20 R

316.03 125.04 [SN] A poss. Ditty in a "T. Problem Pool" approach: Int. Corpus of O's, O's will ^{be} much more frequent than those usable for hard problems. So comes to deal w. these problems is acquired if O's that solved many probs were given Much w. - (which would lead to be true, if f. 'default' solns (of essentially "unsolved" problems) were given default PC's of very low value. →

If we use PPM, this may not be so bad! say O_i is ^{an} O_i containing more "real" solns than any other previous trial. New trials will be add to be like O_i to start, but there will always be a reasonable chance that a new cand will not follow O_i exactly. As soon as it

monitor the rank
 we want to by harder to find cands that work f. probs that have, so far, been unsolv.
 5TM 2.00 → more complete discussion of BOOST/NO.

4TM

o: (317.40) We may have to modify the R's so ~~that~~ they accept/reject the new Q in a proper way.

Or we may simply use the compressed sets assoc w.t. R's as the effective R categorization function.

R. "Compressed sets" constitute a "soft" R function. (E. Probabilistic). I've written a bit (not much) on soft R functs. ---

I don't know where, byr.

Back to 3TM 319-323: Just how far did I get on the "S-function representation problem"?

Ibid 320.17 ff discusses 3I case. ~~Also ibid 319.00-05~~

323.11-18: Bibli. summary (ibid 327.00-30 said to be good but is not v. hard to use)

321.00-09 is outline how to use parts of previously known to be good Q's to make new trials - "PPM" probably does not better.

Ibid 138.23 ... about ~~some~~ Lsrech based on $|O_j| + |R_1| + |R_2| \dots$ order is excess, is appended.

It certainly seems nice, but looks like it would take enormous time! (Vn, for Lsrech, this is a necessary time needed. T. only Alternative I can think of is BU ("nonoptimum") such.

One should do Lsrech BU such in // - T. Lsrech soln is best (if one has a way to find it) ... if one does it then BU such is one alternative: is it the only alternative?

Re: f. Lsrech order of .08: If Q fits several Q's well, we would want to retain it for more

2/0 (over Rg) trials. Ibid 138.28-30 discusses this a bit. also ibid 319.00-05

ibid 318 may be of much interest (but is not in present "Page Packet")

T. "Lsrech" of .08 may not be as simple as it sounds!

[D.24.04]: So for Basic Problem 1: How far have I gotten in planning Phase I (or "Alpha")?

317.30-40; 319.00-19 are certainly on this problem.

In Phase I, I will (very probably) want to use a TSCQ: Whether I will use pure Lsrech;

Pure BU such (or a mix - Most likely) is uncertain. Also uncertain as to exactly how f. 3 would work.

Sub-problems are: ^{Fans} Representation of S-functs.

How to go from D-induction to S-induction in TSCQ. - What would be some good problem

examples?

Another view of TSCQ's: Learning Math "lang" is "problem solving" from examples in Maple, Mathematica etc. - In these cases, for Not-So-Used problems, Phase I may be adequate & d-induction may be adequate. Ideally, we next go on to S-induction.

T. mechanics of this Gauss. are unclear! Depends, in part, on what kind of induction we're

doing ... is numerical or string. (315.29; 316-19) is one way to go from D-induction to "3rd mc S-induction"

the details are unclear & it's not certain that it could be workable.

T. next BIG problem is search over an S-funct. space for solns to the TSCQ.

I've considered 2 kinds of TSCQ's for this: 1) Normal TSCQ for Lsrech 2) "Problem Pool"

- TM selects probs to solve from unordered "pool" of problems.

spec. 320.00

Topics of interest: Working Backwards:

1) Role of Details in Lisp, Fortran, AZ, Z144 ... That we may be able to synthesize not searching or re parsing fast. ≈ 16.04 ... ~~317.30~~

→ 2) $(316.31)^{.16-.19}$ The only way to MAX Efficiency: Combination of DD can ~~be~~ Lang for TM \approx Heads of finding pages in Corpus. So far we've considered Langs like Lisp, Fortran, AZ, Z144, Machine lang. Reyes searchers: PPM (anomaly): Subtree Searchers: "Context generators".

3) Search for Solms to TSQ's: T's "problem pool" approach ≈ 11.29 , $\approx 11.30ff$

← Simplex, Deva Triplets then flats ... 315.09 a new tack using EDIT DISTANCE. (310.21) (312.00) ← on Lisp v.s. BU search
→ (315.22) - T. head (again) of "R" ~~input~~ ~~of~~ EDS it.
→ Edit dist and "R" input of Biome. Not 315.24

315.29: 2 exist. fam from D-funct. induction to S-funct production - Prado more general

S-functs → 316.16-21 more (like truly universal); Biome.

313.25 - 314.24 another Go at time "to map $\approx \frac{n+m}{n-m}$ into $C + m \cdot \ln n$ of sd TS".

313.10 Proof that usual norm defn is best for normalizing sum measures.

312.00-01; -30ff: On Litch v.s. Bsch. (Maybe no real dichotomy!).

312.02: ANT Colony System (ACS): What is it? how it is similar to ALP, ANN, "Wisdom of Crowds"

→ This ties several prob solving methods together. ∴ useful in writing PSM grammar.

run
807 with cellid

Also, Discuss: ask Q's via email, if you don't ask them in class.
Also discuss my behavior "Measure God".
Mention it in Lecture Series.

316.31 is a very imp. point! Th. rapps (ie. subtrans) in the Lispish notation seem to be the most common types of rapps recognized by humans (other than Direct PPM type rapps). Since all universal/large zero rapps are probably perhaps every universal/large will be internally "coded as the Lisp.

See what other kinds of rapps (that are useful for compression) exist. Look at various kinds of problems, - various Break Parus in science as compressions. - try to design search schemas that would have made them more likely. Look at poly 2 & other books on prod. solvng, creativity.

SN Dream Has a micro phone feeding to you. Also, make sure there is no other Mike input & checkers. Also, try different Mike/earphone comb. to see if that's it. Is this from Mike Bustard SWITCH

316.07 Re: Definitions: Consider, say a Fath/Oops: When we make a data, say δ ; and δ occurs n times in corpus our P.C. for δ corpus is $P_{\delta} \cdot (P_{\delta})^n$
 $P_{\delta} \equiv$ P.C. of discovery of δ depends on length of δ ; P_{δ} is the probability of an occurrence of δ in corpus.
If, on the other hand, we are in ≥ 1 & δ occurs n times in corpus, we post-hoc, deduce δ ; its total P.C. in corpus is $P_{\delta} \cdot (P_{\delta})^n$: same as before, but w. additional cc of discovering δ & its compression advantages, if-true. T. main apparent difference is 316.20: repairing previous corpus. This would appear to make a Big Advantage:

Hur., this may correspond to Old, "wiser" Scientists reusing parts of old theories, V.S. Students, not having the older models, try more genuinely radical, "creative" solutions to problems (that older, "wiser" Heads would discard as "Unlikely").

Q: Is .20-.22 a correct analysis of Older, "wiser" v.s. "newer" "beginners"?
If maybe that the "Older Wiser" do not evaluate P_{δ} properly: That they consider P_{δ} very by PC hypothesis because they are too disturbed when they are wrong

Any reasonably good search routine used not make that mistake

(Rob could make other equally bad errors for other reasons?)
(Modern) Removes on 3TM 319-323: "On S-Function is how to seek over them"
Final opinion 22300-03: That it is a problem of extrapolating a function from corpus of (O_i, G_i) pairs.

At that time, I hadn't considered PPM to help w. this: MEJ ($\frac{O_i}{G_{max} - G_i}$) to help work on PPA!
T. present problem is (1) How to assign (comparable) G's to O's that have different ^{sub.} corpi (315.09 is one approach)
(2) T. problem of not wanting to test a new O's over the entire corpus. (Rob's Recognition function of IDSIA Report, was a strong step in that direction.

One approach to (2) After we have a set of O_i 's which R_i has recognized (so a set of sets) for each R_i we compress its O_i set. When a new O_i comes we see which compressed set it belongs to (w/ set & PD). We try to solve for new O_i using O_i 's of R_i most likely to be assoc. w. it. If no success, say δ success, most likely R_i . (319.00)

4PM

"Problem Pool"

10: 315.40 : If we use the "problem pool" is 317.36 : say we get 0's. PGM can work small
 subsets of probs. w.r. to TSA. These functions could correspond to different R values - 25 (315.34 - 40)
 w. R ~~not~~ ~~choice~~ ~~sub~~ ~~dir~~ ~~for~~ ~~each~~ ~~problem~~. [well! not so obvious as to how to do this! - or if it would

03 (Work at all!)] One way is 317.36 : (compression for categories) \rightarrow 320.51
 (like L: Vinyt: Use Compression to Categorize Music.)
 04 **SN** On "definitions": T-way OOPS is KOZZ is AZ do defs looks about 4 - same!
 - Defint from Z-141; in which defs are made ~~after~~ ^{to} ~~to~~ ~~be~~ ~~useful~~! \leftarrow Actually, one could do TSA!
 Analyse & distance of t. systems: There may not be much quant. difference of t so, I can save lots of time! \rightarrow 17.11

07 [T. symbol that includes a data. will have a certain PC: ENBPM, This symbol will give its own ^{sub} corpus of successfull defns of t. past. Hvr, symbols before. definit symbol will influence t. subseq. symbol probab. I could disallow that, but it is a legit statistical dependance, perhaps I shouldn't complain?
 Compare t. details of t-way OOPS, KOZZ & AZ make defns. Also, see how recursive defns occure. They can be very inefficient in execution! If so, have special operator that finds them \rightarrow (Compiler) \rightarrow convert them to more efficient code. - T. PC of t. cond, hvr, is that of t. original, uncompiled code

16 I guess that one of t. objctives I had to OOPS was that t. redundancies in code were not apparent, but (perhaps) they were unlikely to be detected by PPM.
 19 A "good parser writer" is supposed to put all (?) of t. Regys into t. Definitions..
 20 Consider: We chose t. data symbol. This is auto analytically chosen non θ as a problem (i.e. things that have not been detected before). Any by pc sequence following θ will be a good candidate for "compression by definition". Once defined, this data could be used to reparse t. previous corpus! This would (probably) screw up PPM's WB xfm data object.

There may be some point (advantage) to having Defns not apply backwards. It's a way of "forgetting th. errors of th. past" & starting out on a "new life" \odot .
 - No it does seem to run more slowly. We could do evaln of t. pc of a new data. w.o. to parsing t. past!

30 O.k. : So .04--.30 does give a new perspective on Defining $\&$
 31 ~~Neat imp~~ **Neat imp** \rightarrow Neatly imp Q is .16-.19: what lang. Conventions/ Definitions might be poss. for PPM to recognize
 useful regys? [More generally, what lang. Regy data char system] works best? List + PPM searching method. What about Future PPM of Lisp w. OOPS - how defn changes? + PPM also means how $\&$ PPM
 Well, Lispish lang w. subtrac discovery seems fine (but expensive).
 Well, what about being able to recognize ^{sub} nearby things? Consider a certain subtrac, say;
 could it be poss. to index corpus so that all instances of that subtrac could be cheaply found?
 Maybe some kind of U., "holotape" Meme?

Perhaps Study **TRIE** data structures for prefixes on 25 - also to better understand discussion of PPM.

Spec
Q: 31.40: [I assume here, it's wrong to say "solve" a problem ... Some how we have to decide on a criterion for adequacy of a soln. - i.e. what pc for soln. is "adequate"]

So in the "solution pool", we are pushing toward getting solns containing a larger no. of probs.
Say we have 10 problem corpus. If we have a soln for probs 4 thru 10, then we will try very hard to get a soln for 4 thru 10 plus one or more of 1, 2, 3.
We will remove from 1. PPM corpus solns that are much worse than more recent solns. (But maybe not; we may still want some of them in abs.)
A big. Q is how much info to give to various solns.

Q: A slightly different tack: Say we have a prob. in T & Q; We run rounds on all 10 problems & using "edit distance" we obtain for each round a pc for a certain T & Q. At first, most probs will have very low costs, but we will put them into the PPM corpus anyway. We will then increase them depending on their costs for the T & Q corpus. Essentially, we only have one problem. Its an OZ problem. We can, say, elim from PPM corpus all solns w. cost > X: Then we select X to maximize gain - Mc.

We can assign a simpler (but could) pc to solns that take too long to find Edit Distance (2,2) → 320.31

Q: SN on the "edit distance" idea? Use a 31umc, try a few short R inputs, then use R w. smallest edit distance, or use all 4 or 6 ... R values & pick "closest" one to "True A"; Or, use 1. set of R values w. outputs "close" to A. Trying to get sort

Can move toward a useful "normal 31umc" w. R receives A exactly. 310mc 310.12

Q: SN on "edit distance". If 2 strings of length n & m are very close, it should be off by only

to find a short edit distance - or \ll time than \ll \ll min

Clearly we can use .09 - .15 when T & Q corpus is very large. Checking all previous QAs on a cond would take too long.

An idea related to (2,2): We may notice that a certain subset of problems gets by pc w. O_j ; Another subset of probs get by pc w. O_j . We then try to find a way to recognize (R system) which subset a problem is in, & use f. appropriate O_j

What I'm noticing: that using R input = A (i.e. 2 input umc) w. edit distance is essentially a using D-functions! So this is a path from D-funct to structure. The next expansion is

(.16 - .18) - A simplified version is to have each d-function have 4 or 5 ... outputs; respond to R having 2 or 3 poss. values, Next; perhaps have R w. values \ll w. d-funct w. corresponding to f. "Long R" of R.

Q: "R" device is equivl. to having a bunch of O_j d-functions a for each problem, picking them at random (w. each O_j func having its assoc pc (\equiv w.)). I discarded a model that was

Yes, some time ago, but I think the R choices for all QA's had to be the same. If the R choices are indep for each QA, then its OIR, i.e. \equiv normal 31umc. Go over this idea carefully & it seems very imp! Also Note 317.36 on 2 soft (probab) "R" function (316.00) → STM 4.23

313.90 : If σ^2 is true var., the model will have var $\sigma^2 \frac{m+n}{n}$. Go thru Calcus carefully using this fact.

Model PC: $\prod_{i=1}^n \frac{1}{\sqrt{\pi} \sigma} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2 (1 + \frac{m}{n})}} = \prod_{i=1}^n \frac{1}{\sqrt{\pi} \sigma \sqrt{\frac{m+n}{n}}} e^{-\frac{1}{2} \frac{m}{m+n}}$

True model PC: $\prod_{i=1}^n \frac{1}{\sqrt{\pi} \sigma} e^{-\frac{(x_i - \mu)^2}{2 \sigma^2}} = \prod_{i=1}^n \frac{1}{\sqrt{\pi} \sigma} e^{-\frac{m}{2}}$

Assuming m values of μ corpus are "Gauss" in both cases.

$$\frac{m}{m+n} = \sqrt{\prod_{i=1}^n \frac{i}{i+m}} = \sqrt{\frac{n! m!}{(n+m)!}} \left(\frac{n! m!}{(n+m)!} \right)^{1/2} = \left(\frac{n}{n+m} \right)^n \left(\frac{m}{n+m} \right)^m \sqrt{\frac{n \cdot m}{2\pi(n+m)}}$$

For $m \ll n$: $\left(\frac{n}{n+m} \right)^n \approx \left(\frac{1}{1 + \frac{m}{n}} \right)^n \approx e^{-m}$

e^{-m} could be part of the Gaussian factor, $\left(\frac{m}{n+m} \right)^m = \left(\frac{m}{n} \right)^m \cdot \left(\frac{n}{n+m} \right)^m \approx \left(\frac{m}{n} \right)^m \cdot e^{-\frac{m^2}{n}}$

But $\left(\frac{m}{n} \right)^m \cdot e^{-\frac{m^2}{n}} \approx \sqrt{\frac{m \cdot m}{2\pi(n+m)}}$

$\left(\frac{m}{n} \right)^m$ looks like a variation factor for $z = \frac{m}{n}$! (Remember: check de facts!)

Fits in $m (1 + \ln \frac{m}{n+m}) = m + m \ln \frac{m}{n+m} = m + m \ln m - m \ln(n+m)$

313.25 R: $C + A m \ln n$ $m + m \ln m$ could be part of C : in which case, we get $A=1$, not $\frac{1}{2}$.

Oh! I forgot to take square root of the whole expression so for $\frac{1}{2} \ln$ or $\frac{m}{2} + \frac{m}{2} \ln m + \frac{m}{2} \ln n$

plus $\frac{1}{2} \ln \sqrt{\frac{n \cdot m}{2\pi(n+m)}} = \frac{1}{4} (\ln n + \ln m - \ln(n+m) - \ln 2\pi)$

$\frac{n \cdot m}{2\pi(n+m)} = \left(\frac{n}{n+m} \right)^n \cdot \left(\frac{m}{n+m} \right)^m \cdot \frac{1}{2\pi} \ln \left(\frac{n \cdot m}{2\pi(n+m)} \right) = -\frac{m}{4n} + \left(\frac{\ln m}{4} - \frac{\ln 2\pi}{4} \right)$ part "C"

So the "C" part is ok, but maybe $\frac{m}{2} \ln n \approx \frac{m}{2n} = \frac{m}{2} \left(\ln n - \frac{1}{2n} \right)$ for the part that depends on n . The m depends on n .

I did assume $m \ll n$ in the approx. — so maybe this is acceptable for the $A m \ln n$ part.

$\frac{1}{4} \ln \left(\frac{n \cdot m}{2\pi(n+m)} \right)$ The part involving n is $\frac{1}{4} \ln \left(\frac{n}{n+m} \right) \approx \frac{1}{4} \ln \left(1 - \frac{m}{n} \right) \approx -\frac{m}{4n}$

So, its "in the ratio ball park", but I'm really not so sure its right, because my algebra is so buggy!

The + derivation of .00 ft is not so complicated!

311.40 : Consider 311.29 A (the "problem pool approach"): Suppose we had a way to list

S-factors is something like PC order (Note 312.00-01; 307 or a possibly better way to deal with).

T. Q is, how do we turn a bunch of QA's? In we find solns to single, then pairs, then triplets, etc.

What PC goals do we accept for singlets, pairs, triplets, etc? What wts do we give to solns in the say "PPM corpus"?

If we had one (QA) only in our TSQ, we'd just work on it until CB was reached.

If we had 2 in pool: Do one, Do the other then do both together. In 3 in pool: do all pairs first, then all 3.

Note: Working any subset of the corpus used to assist in solving the entire corpus

For a single corpus, we save by trying to do as many probs as poss. w. one Q:

NIPS
3TM 216.00
Redi: 4150
Korea Stack of
ED 611.00F
T < 2T

00:312.40: We do Lsrch using PC_0 as PC index & we use $\frac{CC_0 + CC_2}{PC_0}$ as a stopping criterion in Lsrch; In // Lsrch we work on 2 cand so as to keep $\frac{CC_1 + CC_2}{PC_1}$ as low as possible for all cand.

Result is that we are doing Lsrch on a 0-1 problem, w. $G = PC_1 \cdot PC_2$.
I think first for d-induction, we can use (INV problem) simple Lsrch, but for s-induction, we have to keep it as

a 0-1 problem. So before starting a problem of 0-1 Lsrch, we try to find best Gove w. available C . It's certain Gove, G_0 is known for a problem w. 2 cand whose $PC = P_0$; CC of soln is C_0 then we will find that Gove = G_0 & we use w. $CB = \frac{CC_0}{PC_0}$.

T. idea of T.S.Q. is to make PC_0 by T. T.S.Q. is "Adaptivity" to search PD will not modify CC_0 .

(The ~~TM~~ TM can affect CC_0 in "Phase II" TH's.)

HA! I found that proof that ϵ best normz function for same measures is

$\bar{x} \rightarrow \frac{x}{x+y}; \bar{y} \rightarrow \frac{y}{x+y}; \bar{x} + \bar{y} = 1$. Its on (NIPS) 3TM 207.10-19

T. problem is lines 12-13 and line 18. \rightarrow p.c. produced to the target corpus of. univ. d.f.

Work done $E \sum_{i=1}^n (C_{i,j}) \leq \frac{P_{sol}(n)}{u_{sol}(n)}$

for finite derivable u , $\frac{P_{sol}(n)}{u_{sol}(n)}$ will be bounded by the PC of the deriv of u .

Here, the minimum value of the normz constant occurs when one chooses a sequence of bits in the sequence

\Rightarrow This min is maximized when $f_1(x) = f_2(y)$ for all of the choices

I. any function for which $f_1(x) = f_2(y)$ and $f_1(x) + f_2(y) = 1$ is $f_1(x) = \frac{x}{x+y}$ & $f_2(y) = \frac{y}{x+y}$

Proof: $f_1(x) = \frac{x}{x+y}$ so $\frac{x}{y} f_2(x) + f_2(x) = 1$; $f_2(x) = \frac{1}{1 + \frac{x}{y}} = \frac{y}{x+y}$

So $f_2 = \frac{1}{1 + \frac{x}{y}} = \frac{y}{x+y}$. The Proof looks pretty reasonable

HMC (how many costs) In Sol 28 p 421 last part of col 1; I set ϵ in ϵ of ϵ derivs. \rightarrow structure of model. Consider linear predn w. m covs, n data pts. but $h(x)$ is replaced by $C + A b / \sqrt{n}$. \rightarrow no. of params (continuous) characteristic of accuracy of model

Using $G = \frac{m+m}{n-m}$ formula: True model will have $PC = \prod_{i=1}^m \frac{1}{\sqrt{G}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$

Subst $G = \frac{m+m}{n-m}$ \rightarrow $PC_{true} = \prod_{i=1}^m \frac{1}{\sqrt{G}} e^{-\frac{m}{2} \frac{(n-m)}{n}}$ \rightarrow $PC_{my} = \prod_{i=1}^m \frac{1}{\sqrt{G}} e^{-\frac{m}{2} \frac{(n-m)}{n+m}}$

ratio = $\exp(-\frac{1}{2} (n-m) (1 - \frac{n}{n+m})) = \exp(-\frac{1}{2} (n-m) \frac{m}{n+m}) = \exp(-\frac{1}{2} \frac{n-m}{n+m} \cdot m)$

m corresponds to b , the no of costs.

Woops! In the \prod , I didn't really do product of $n+1$ to M

It says $\exp(\frac{m}{2}) = \prod_{i=1}^m \frac{1}{2+m}$ start sum at $i=1$ $\frac{1}{1+m} + \frac{1}{2+m}$ It may be a $\ln(n)$ sum $\rightarrow \frac{1}{m} + \frac{1}{m+1} + \dots + \frac{1}{n+m}$

$\sum_{i=1}^m \frac{1}{2+m} = m + 2m \approx \frac{1}{2} = (m + 2m)(m+8)$

This approach may be correct, but the $\frac{1}{2}$ has to be given over carefully! \rightarrow 3K.00

0 : 31/01/23 : N.B. Int. 3 edit best idea (i.e. its Genset) of 31/01/21 : This really doesn't do traversals P. Cordero falls
is more like BU graph! Network! 30 30

ACS (Ant Colony System) : first use of $i \rightarrow$ ev. Compu. : (1997)

Full PDF in P3 : 12/15/04 (I also have hard copy of journal)

Authors say it was before Rein. Sim. based on "Evolutionary Compu." for perhaps Trav. Sls prob (TSP).
My impress. of how it may work: We have some ants that run a maze; making distances run by each (ed:)
after it finishes. Then each path is given pheromone wts $\propto \frac{1}{d_i}$ of the ants; that
thickens that path. Henceforth, ants having a choice will increase \propto to total pheromone of their
choice. I think paper is readable, but disregard part that says intuitive applic. to behaviour of Biological Ants!

This is similar to ALP: we have a set of alphabet letters. We have a set of codes

I don't think ants would ~~survive~~ survive using only they do eat!

But now strings of these letters - w. empirical store (G_i). If we regard each string

as a comp. code G_i \propto prob. of that string, then we want to assign wts to G_i alphabet
(\equiv pc's) \rightarrow t. total parallel pc off. solns is min. say L_k is the later in alphabet.

m is t. total no. If $N_{k,m}$ is s. no. of times (L_k) is used a max. freq. ($N_{k,m}$ often \propto)

Then val. of L_k is obtnd by maximizing \sum parallel pc of code by suitable
wt. assignment to L_k . It can be solved iteratively, by assigning initial wts & modifying
Reason basis of approx pc's of individual trials. (This is standard approx method for comb. prob. in

So next approx for L_i : \propto total wt. of L_i in successful trials

each trial is given wt. \propto product of pc's of its letters.

In t. TSP: two letters have "contact" (i.e. each city has only a few roads
coming out of it). For each contact, we have a separate alphabet of letters.

That (ACS) Rep should do so well (if indeed Rep did!) is encouraging for ALP!

Another possibly positive encouraging to ALP is 'The Wisdom of Crowds'

Possibly similar is Bad Prob ANN: Average wts of unsuccessful could solns.

Files
552/12/10/04
260 SP/12/10/04

10:00 Actually MDL has some diffy: we want a 2 part code to construct models in PC order, but
evaln of model can be surprising! While this is true of MDL, it is also true when we most
average our codes in \parallel . In fact, it is true for many, if not most (or even almost all) S-funct evaln
schemes!

If it is how do we compute cost (\equiv cc) of a particular, known. soln? Well, say we know part 1 of 2
MDL code has $pc = p_0^1$ & $cc = cc_0^1$; part 2 has $pc = p_0^2$ & part 2 has $pc = p_0^2$ & $cc = cc_0^2$

To part 3 soln, we assign $pc = p_0^1$ & $cc = cc_0^1 + cc_0^2$. p_0^2 is not regarded as a pc in this part of 2 soln -

we regard it as part of Gene; t. Gene is $(p_0^1 \cdot p_0^2)$. T. time to find it, using \parallel Level is $\frac{cc}{pc} = \frac{cc_0^1 + cc_0^2}{p_0^1}$ (3.00)

D 13 Oct, 4PM

S-functions (cont)

20:310:40 : finite no. of X_i , one into randomly sampled P from — but \emptyset P over an almost-always \emptyset X_2 .

② Even w. feedback of X_2 , one has no idea if X_1 ones have been included in P .

It ~~is~~ $[X_i]$ are simple scalar ones, no way be able to normal P from w. \emptyset randomly constructed $f(X_i)$. For vector X_i , not so easy.

Hvr, t. attracting ~~function~~ $f(X_i) = P$ is that it may be easy/fast to compute X_i — which would lots of time.

Conditional P.D.

So: Dropping $f(X_i) = P$ for t time being: Consider other forms of S-functions.

Consider t. \exists time: Starting w. $t=0$ time: One input defines O^j .

Another input is for target X_i (string or no.).

Another input is \exists random component, R .

The output = $U(O^j, X_i, R)$. Ideally t. machine reads all of O^j , then all of X_i .

Then read all (or part of) R and \neq then prints t. output. If it is allowed to print out part of

Looking at only part of R , then t. meaningful parts of local R 's will form a prefix set,

So their $\leq PC \leq 1$. If they are a complete prefix set their $\leq PC \leq 1$.

R need not form a prefix set: R can be a by finite strings of symbols w. an "end" symbol.

T. PC assigned to t , and symbol will be determined by a sequence layers of solns to problems

"t. present context". We assume symbols of R are uncorrelated. If they are binary

(0,1, end) then they will have $PC = (\frac{1-\epsilon}{2}, \frac{1-\epsilon}{2}, \epsilon)$ PC of a string of n binary and

will be $(\frac{1-\epsilon}{2})^n \cdot \epsilon$. ϵ is chosen to max PC of set of known R 's $\in \mathbb{R}^n$

We can actually make $1-\epsilon$

Hvr, I don't think I ever found a good way to search for solns: That older note in ~ Nov 2003?

It was on a yellow sheet of paper: first ~~at~~ Page 33 then \rightarrow (33)

That's it!! 3TM 133.00 \rightarrow bid 137. \rightarrow 138 (as) Buy on bid 138.285

3TM 319.10 ... $\frac{320.11}{23.21}$ supposed to be more serious approach — That then was Bugs in this

The it may be that w. recent discoveries, I could've found "BUGS"!

A diving back I'd consider for TSQ's is $\frac{2}{3}$ "problem pool" approach: We have an unordered

batch of QA's: TM tries to find individual solns then parts: then ... as many as possible by some O^j .

This builds up f. sets of cons in t. composition as usual way. We then give each of n the \emptyset O^j that can solve n of n QA's: \neq

Hvr, t. target is for d -induction. I had a previous discussion of this "poor" approach, but I don't know if it uses for S -predn. Well, even d -predn has a PC assoc w. each O^j .

Say we find a bunch of O^j 's: each assigns a PC to some subset of t. QA corpus. How to wt. this corpus of O^j 's?

Solns of time varying
 \emptyset 2 proofs
3TM 310.12-25
How Many Cuts,
3TM \approx NIPS
298-299; 303, 304
Case
3TM 266 on
Poles stack lang.
Conv. K. for
TM 2579.00
TM 65-30 Proof
that L. Schenker
only worse than words
in least order!

SPW
314.30

S. Funct Rev.

start w. discn. of why $P = f(x)$ $x = \begin{pmatrix} \text{number} \\ \text{string} \end{pmatrix}$ is not such a good form

We'd like formulae where pc would be > 0 is normed. (maybe actually > 0),

Basically, T. problem is to get a form $\in L$ such (or B such) is least expensive — In

particular — we are interested in the problem of minimizing PC of a OJ wrt

“corpus up to now” Since OJ is usually an S-funct (is we are most interested

in OJ 's that are S-functs because they enable phase 2) [Actually, even when OJ 's are d. functs, we will prefer by ^{to} appropriate to get

PC's of solns. ABDEFG 'a.b.c.d.e.f.g. ABCDEFG.]

We want S-functs in form P can be rapidly generated & tested to find PC of f. corpus wrt. P can.

For Zrch, we'd like to be able to generate S-functs in PC order — .30

Re: 3.10.04: I think I have a better understanding of it now so that t. \approx AZ lang

might not be bad (I do want to find way to do recursion or just do loops (to enable Prim. Rec. functs)).

One way to do 3.10.04: We use AZ to construct a function that uses O_i as ~~one~~ possible

input — whose output is $\approx A_i$. Actually, it's $f_i(O_i)$ a t. pc assigned to O_i 's output is its

“distance” from A_i (t. “correct” output). We need a good measure of distance (this is a clustering

approach). In many Numerical problems, it's simple $|x-y|$ norm — abs value or square or modulus squared

would be ok. But for strings it's not clear what measure to use. Ideally, we

would say: “How much info we need to turn $f(O_i)$ into A_i .” This usually takes time long to compute.

One approach is “Edit distance”. See Google; also (PS) editdistance 12/11/04: Also see Bookmarks on EditDist.

That paper says using dynamic it takes time with to compute edit dist between strings of length m & n resp.

The Edit dist can have a different cost, for Additional Deletion or Substitution. \rightarrow see $\begin{pmatrix} 15.069 \\ 15.21 \end{pmatrix} \rightarrow 312.00$

Consider S-functs of f. form $f(x) = \sum_{i=1}^k p_i x_i$: would like f to be > 0 for all x . if poss. normed \in

So we can compare various ~~candidate~~ f. forms. One way to \approx Norm: We want normed $f(x)$: Just try ~~various~~ k randomly chosen x_i . Then

$\left(\frac{1}{k} \sum_{i=1}^k f(x_i) \right)^{-1}$ will be ~~norm~~ a norm const. If ~~cannot~~ start w. small k ; if could be good,

use larger k for get more precision. **No!** True norm const. is $\left(\sum_{i=1}^k f(x_i) \right)^{-1}$.

By studying $\theta_k = O_k^2$ of $g(k) = \sum_{i=1}^k f(x_i)$ could we estimate $M_{k=200}$? — No so easy!

We have no idea as to whether large $f(x_i)$'s have been included in f. set. If ~~there were~~ $\rightarrow 311.00$

S-functions

0:308.40 : Even in BU mode we have to deal w. inventing "R" functs or eqn's. We can use the idea of 297.25 (AP) as a base on a fair no. of hard problems only. ~~IF~~ When we find a fair no. of ~~easy~~ ^{hard} cand's that work on that set, we begin focusing on random "easy" ~~or~~ problems. This may be good ex of to keep Lsroh & BU probs.

— The getting "R" functs is probly deservable — That may be able to generate t. eqn's of "R" functs by itself, w.o. a former telling it directly about them.

1: (308.40 spec) : T. idea of 308.37 is imp. ~~the~~ The technique descr'd would be fine for Lsroh of S-functs.

Somewhat reminds me of BU search — But BU search starts w. low G cand & tries to improve it ^{by mutation} ~~by~~ ^{GA, GP} ~~by~~ ^{fitting S-functs to data.} ~~by~~ ^{MEJ} w. assoc S-grammars.
Classified by mutation: less classically by crossover & ~~by~~ ^{MEJ} w. assoc S-grammars.

So maybe t. stuff I know now is not to do an acceptable "phase 1". ^{Like 302.20} (16)

SNAP poss. source of ideas for forms, representations of S-functs is to Duda, Hart, Stork Book!

6: (14) → T. 2. Main problems recently "Solved" ^{see 295.04 for some notes: 302 E.00-40:} How to represent a variety of S-functs to enable Lsroh & perhaps BU search. (2) How to deal w. problem of verifying O^j (in either Lsroh or BU search)

On t. whole previous couple's, when we are doing a larger $f_i S_i Q_i$.

I should perhaps write good summaries of Phase 2 "Sols". ^{Quick Abort} ^{(294.32-.40) is early form of "Quick abort"}

For (117) & Verifying, evaluating a O^j cond. (306.24-27) (297.25-40) (309.00-.05) ^{pp. 8, 9, 11} ^{part.} (272.37) a somewhat different approach!
The "R" function (Recognition Function) Discussed in Early part of FD31A ^{part.} (Also Note 250.11 for "problem sol" approach)

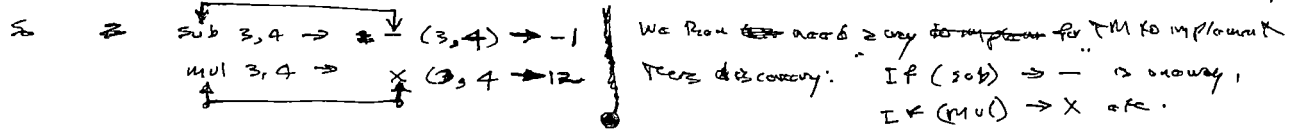
For (1) Representation of S-functs! (15) 274.04-40; from \approx 265 on: (266-267: on Universal Analog Machines)

- 270.13-19; 271.00-07
- 302 E.00-40; 308:00-40 (308-40); 309.00-12

0: 302.40: Originally I had that I'd start w. d-prada; then "go" to spread, but I didn't have clear ideas on how to do this - I write 302.40 to do it now... but what screens advise is:

To start w. s-functs: to first learn $(8, 37)$:
T. machine learns that soln. is $8, 37$
+ .25
x .25
/ .25
It then tries to correlate f.

empirically correct functional form, with correct soln. This is easy to do, because 4 functional forms & 4 solns. differ in only one place & those "out of phase" correlate 100%.



Here we have broken the 1 vup task into 2 parts: T. first part is boost by 2 (ob & f. search) & boost even more.

T. problem of ~~input~~ to evaluate complicated Arithmetic functions, reasons. T. reason its annoying is because I think I know how to do it, yet

I find it diff. to get TM to learn it - I don't know what primitives & observations to put into TM so it would get a my idea of "sub object", "f. eval in object".

I. idea of associating a number (value nor any a number) w. usable term common sub trees.
Some ways to deal w. this: ① find acceptable soln (by me) & put it into TM.

So TM has f. soln & 4 sub. abs. Go onto other problems, T. experience will perhaps enable me to get a usable soln. to this problem.

② in "SAARB TM" I had a stack machine & I think it solved this problem easily... but this was because notation for problem was match by "stack" structure of reference machine.

NB It may be easy for TM to use v. some methods that I use, to solve problems if I am to use Lsrb, which (usually?) means that solns. will be in v. form that I solve the probs.

27 **17** Possibly relevant: T. sub trees are "objects" w. "f. property lists" ("slots"). "Value" is one poss. property - A dirty hack is that we'd perhaps have "too many objects" for us to have property list for each!

0: 302.23 while 302.00-20 func P.D.'s on continuous receptors & at discrete analytical pts, finite m.no. or m. frunks. It doesn't give P.D.'s on ~~strings~~ (finite or infinite). ALP certainly does that! - But usually in a way by cc way.

302.27-37 Does ~~some~~ P.D.'s on strings. Actually, if we use a "proper" P.D. on code length ("3 input one" "3 input" isn't so bad!) for each (QA, say) try (A generated, ~~some~~ we compute (A - A correct) & add to list. code length of A. T. defn. of Norm can vary! We can start w. some simple, default, Norm, & modify it for each Q or Q class, when we get large and **SSZ** - (see discussion of "clusters" on 302.13f-37 its partly f. same thing)

37 **17** Hvr, as desc'd, these "clusters" aren't much good for Lsrb! WOOPS Unless we do an Lsrb in 1's by "cluster centers", & work our way is pc of cluster x pc of A or w. that cluster. This differs from standard Lsrb... 309.00 309.10

MEJX .00
May (winds)

20 : Variation on MEJ method ($\frac{\sigma_G}{G_{max} - \bar{G}}$): (MEJ winds)
 We have this complete population: In Bins for various G ranges. So say 100 bins of = width.
 (Or choose bins for Σ population). Then for each bin, we keep track of its population (Σ wt)
 its \bar{G} & σ_G^2 . Then, if we make a P.D. on the winds where $p \propto e^{-\text{Order of } G \text{ from top}} / \lambda$
 23 for each λ we have a different S-lag (Σ PD). For each λ we can obtain its
 $\sigma_G^2 \approx \bar{G}$ since they are both obtained by wtd means of the bins; wts are population
 times something like $e^{-\lambda G}$ or whatever. So we can ~~calculate~~ calculate $\frac{\sigma_G^2}{G_{max} - \mu_{G_{max}}}$
 for each λ & pick the one that is highest.

Using this "X winds" type of d.f. gives us benefit of occasionally low Σ trials —
 which is not true for sharp cut-off of "normal" MEJ ~~Call this MEJ winds~~
OR MEJX

HA! Some fun on Old MEJ! Say we have "Bins" of .00-.03. Now say "X" is t.
 G level of a Bin. Consider $\frac{\sigma_G}{G_{max} - \bar{G}}$. If we add a new bin's data to the existing P.D. (Σ G_{max})
 σ_G will change & \bar{G} will ~~change~~ change. If they both change in same proportion, we won't
 2 MEJ max. Say our candidate bin has a certain $\sigma^2 \approx \mu$. The change is down easy to compute.
 % change in Numerator is harder.
 $\sigma_G \rightarrow \sqrt{\sigma_G^2 + \sigma^2}$ ()
 $\text{new } \bar{G} = w_1 \bar{G}_1 + w_2 \bar{G}_2$
 $\text{new } \sigma_G^2 = w_1 \sigma_G^2 + w_2 \sigma^2 + \frac{w_1 w_2}{w_1 + w_2} (\bar{G}_1 - \bar{G}_2)^2$
 So I'm not clear on this! At first I thought it would be easy to find if numerator is down
 changed by same ρ ; so I'd know when I got ~~off~~ MEJ: But maybe not!

22 I also had that earlier "Complete Model of TM" ²⁰⁰² IP 510.09-10; 509.10-18
 509.36-510.08, 509.19-44
 This was the idea of a Grand PD used for TM₁ & TM₂ also. It was described in introduction of PSIA Report.
 So: I ~~was~~ want to (List, characters, Deriv) all of these & "Complete" Models of TM. I ~~was~~ was was Model
 had some 2 d. fits that I may have solved ... reference 309.16-18

474

Bibli. 100

20 : A MSR paper (2002) on improvements to PPM. Seems to have good discussions of various issues:

P: \PS \ PPM improvements - readable - Record-Process, PPM 11/18/04 date in "PS" file.

0: 305.21 (5) A fifth kind of thing is directing. - Say a set of RISC instructions; The objections to FOR (But it's not easy to find definitions, constraints) apply even stronger to RISC insts. Typically we have a small set of insts that get into: place ops (a) a logic of registers, (b) source stacks. (c) RAM (access). (d) Jump instructions. IT advantage is extreme speed. (Pro ~~PPM~~ = ~~slightly~~ compiler for PPM can do a better job.)

20 : Next: Suppose we have chosen a language: several choices of action:
1) we can do Lsearch to find ϵ candidates \rightarrow they fit the corpus $[Q_i A_i]_{i=1}^n$ as well as poss. As regarding P.O. we use a s-lang. w. corpus of $[Q_i A_i]_{i=1}^n$. We do this in \approx PC order.

24 : 2) If as in (1) we use Lsearch, this means each candidate tested on entire previous corpus. As corpus grows, this can become very expensive. One soln is to try worst (most diff) CRA's first. [Quick ABORT! see 297.25 ff.] see 297.35 for objection to P.O. - Suggests Prod "R" method of ID's reports in the direction.

27 : 3) We can choose not to do Lsearch (see 303.10 for dozen of Lsearch's. $B_{\text{new}} U_p$ search.) So in BU search, we start w. a not very good ("Default") code for corpus $[Q_i A_i]_{i=1}^n$ and we do "mutations" or other heuristic methods to try to improve it. GA/GP is one way to do this: we still have the "Test whole corpus" problem of (24-27), but. For search we can use \approx PPM and probably the MEJ ($\frac{O_c}{G_{max} - 2G}$) method (or mod of it),

30

4 TM

128
11 x 12 = 132
11 x 11 = 121

20: 302 27 **SN** in Every TSO₃W. s. func's: A early problem type could be to find good contexts for prodn. This is an OZ problem. → T. Problem of A.I. has been solved! → 302 27

22: 302 40 **A (Main) Main Part**: Make list of various Approaches to TM, à g. r.v. d. r. s. i. e. s. B. o. t. t. o. m. e. s. in each. // some refs: 294.20. 54.16: 2000 272.25--30

Bitwise I had been worrying about **Choice of a Lang.** for Lurch: Lisp time lags had ~~an~~ apparent dirty Past they normally weren't able to ask for more bits in it. This can be dealt w. in several ways:
① I did have an alternate system to ~~get them to ask for more bits.~~
② (new) One primitive to ask for another bit
③ Just use regular **Lisp** form, in which every string is a legal program. Assign pc's to reg's very some (post hoc optimized) function of pcn by pc — Not a universal pc of t. integers! **N.B. 294.00** **Note 294.14 ff** For: "primitive TM" one can use e. unv. d. f. on integers. [This may be about **Some as 1.065**]

④ Using Fortran: One problem was construction of data. Fourth (Lisp) gave definition facilities but they ~~didn't~~ were created as part of code for d. v. s. of c. n. d. s. A more useful kind of data is "fast hold": i.e. we just code could look for sub trees that have been used full — so that defining Recm compresses code. So generally, it seems that this process would be much more useful in Lisp than in Fortran. This is because

Any sub tree is a Lisp function. In fact, I don't know what a Lisp sub function looks like! **30** **SN** A probl. **BUG** in my view of functional layers: There can be a function ~~that~~ $S(x,y)$ which stores x at address y . $R(x,y)$ then retrieves content of memory y . The function $S(x,y)$ produces no immediate output; Also if " R, z " say appears in a formula, it doesn't usually mean something as " R, z " appearing elsewhere in the formula! $S(x,y)$ is a function w. "side effects".

IN Fortran, we can use stack as ~~is~~ not many (One is not such a good idea for!), so " z " puts z out. **stack**. In OOPS what sort of Memory Combinations are there? In OOPS the closest thing to **"Memory RAM"** is just that which is ~~the~~ **RAM** object for stack;

On next that copies from a place on one stack to a ~~new~~ place on another stack. **ON p 29**: "B. A stack of integer arrays ... that was in Rice's paper": This might be used as an ordinary RAM. Anyway: There are functions that use RAM in simple ways & I should try to characterize those functions. **See of phrase** is a way to break them down into "sub functions" (or sub... objects!). I can treat this messy stuff. 20-26 as ~~the~~ messy "logical reasoning": Try to write TSO's but don't use it, then when I used it, work back. **25-26** has an idea on how to do many. **N.B.** ~~the~~ function defus, via ~~inf~~ series don't use > 30 registers: They don't use RAM to generate the function?

21: **14** **Re Fortran**: Try writing some reasonable functions in Fortran (e.g. see what kind of "contexts" are useful for prodn. See if (2) PPM would be useful for prodn. in these things. → **306, 10**

22: **SN ON DEFINITIONS**: T. a advantage of Dofus over PPM is its "improvements", is that when we compress using Dofus, T. compressed code is in a form that is likely to suggest new defus (compressions). I don't think this is true (or "astute") for a "PPM". w. defus, we normally have parsing ambiguity, so we (a) may want to reparse code after each (or after several) new definition(s) (b) A/O use covered by PC parsers in 11 & look for new defus in Recm.

4 TM

20: 303.40: Normally in Lsrek, we have "Updating" in which the Guiding PD is updated in view of recent

MEJ
Max Error
Jump

Problems Solns. Updating is usually the most difficult part of TM's activities. Just how can PPM
be used in Updating of Guiding P.D. ? $MEJ = \text{Max Error Jump}$

Well, we can use previously "successful" P.D.'s (OJ's?) as corpora ... (possibly wtd):
PPM (either regular or using more complex contexts) can induce a P.D. on ~~the~~ OJ cards & we
can Lsrek them in PC order.

If we want a PD that involves Gore, we can use .005 ~~with~~ Max Error Jump ~~MEJ~~.

MEJ is one standard way to solve open problems - certainly not the best, but "not bad"
A better way would give a better PC of G as a funct. of dem of card, or would give
ways to construct ~~random~~ by ~~random~~ cards that led by pc of by G. ~~we~~ w. input
data of a set of cards & their G's: to make an S Grammar that assigns PC's to
(card, G) pairs. If we ask for a carried G, it will give a P.D. over cards w. that G.
Hr, perhaps there is a more direct "object" or "system" to find ~~optimal~~ options rapidly
using past experience.

Rothsats!
See
D:\SS2
But mostly
D:\PS
Also special
Blair in D\PS
in "Rothsats".
Is that some
my own idea for
"Approx. Lays"

16 In MEJ: instead of making a long ~~list~~ ^{sub matter} for each G level, I do an interpretation of a long ~~list~~ of all cards w.

G's > X. So each X has a long. But this MEJ was designed for GA: It had to do w.
M & G of ^{offspring set} pairs taken from a long. - well, that was our application. There was a major
direct application (I think). I did work out that application: it did not
use PPM

19 I've forgotten how I did it! - This uses PPM or a Grammar.
Say we use as corpora all cards w. $G > X$. This set has
a ~~mean~~ $\bar{G} = \bar{M}_G = \bar{G}_G$ & a G_{max} - all as functs of X: we pick an X_0 & $G_{00} / (G_{max} - M_G)$ is max.

25 As we evaluate cards, we ~~keep track of~~ add cards ~~that~~ w. $G > X_0$ to the long, &
it's $G_{max} M_G = G_G$ change. We also keep track of these points for the "new" by "longs"
w. $X = X_0 + \Delta$. ~~is~~ $X = X_0 + \Delta$: When one of these gets better MEJ is
we change X_0 to $X_0 + \Delta$ or $X_0 - \Delta$.

So .16 - .19 was for a normal GA w. crossover matrix.
20 - 25 was for "GA" w. PPM & no mixing or mutation.

30 One Ditty w. ^{Lsrek} Improving/Updating: If we use set of previously "successful" OJ for our
corpora: we'd like to w. different cards, but it's not clear how! I'd like to associate each
card w. some empirical Gore ... but unclear as to how (inf. Initial QA, T&Q)
If I give w. of 1 to all accepted cards, then I can't do MEJ .015, since
all cards have some "Gore".

What we want is a good compression method for the seq $O_1^j, O_2^j, \dots, O_n^j, \dots$ seq. of successful cards.
Actually, what we want is an optimum compression w. side info - To say of Q's is "side info".

0:502.13 / 2n' How many Coits.?' : $\sigma^2 = \frac{n+m}{n-m}$ v.s. $\sigma^2 = e^{\frac{2m}{n}}$ v.s. $\sigma^2 = \frac{n_0}{n-m} = \frac{n+1}{n+1-m}$

One way to decide/check on my algabre. would be to simulate $X_n = \sum_{i=1}^n X_{n-i} \cdot 2^i + noise$. and measure t. observed errors.

An easy way would be to use "consil compiler": it has matrix inversion; (Re inverting matrix is faster; n^2 rather than n^3 ... it's still easier to use pre-programmed matrix notation (?).

Also Looking my formula for k coits in Sol 78. — Also look at TM notes of next time to find out how I got these ideas: I. formula in Sol 78 says: $\frac{c}{n} \ln n < \frac{c}{n} \ln n$ $k =$ no. coits, n is no. data pts. $c =$ ln prob deriv of express containing k param. I got $\frac{1}{2}$ from Robinson, I think. T. idea in 78 was that in term 2 (coit) $k \rightarrow c \pm k \ln n$ (in \ln domain rather than binary logs). || Did I get the backward message $2k \ln n$

$\frac{P_c(n)}{P_c(k)} > e^c$
 $n = \frac{c}{k}$

But process & disc is $\frac{2k \ln n}{26 = (1+2k)}$
result of $\sigma^2 = e^{\frac{2k}{26}}$
so PC of 0.5018 \rightarrow
 $n = \frac{26}{26 - (1+2k)}$
 $\frac{n^2}{26 - (1+2k)}$

SN On LSuch v.s. BU search: LSuch is good when trainer knows (or suspects) what answer is a shorter code for the soln. that is within acceptable CJS.

When this is not true, then BU is preferred. It is very important that comes.

used for both methods to a reasonable to one another as much as possible.

In LSuch, if TM gets a shorter soln. to problem than trainer expects, TM should continue such to find Trainer's expected solns. This is so we can see that TM has acquired all of k.

coits that trainer expects it to have, so trainer can continue writing T52's w. trying to understand TMs "break thru". TM will always be able to solve trainer's problems in time that

Trainer expects, here. If TM is consistently finding significantly faster solns. than Trainer,

it is clear that TM has "gotten ahead of" trainer. Ideally, trainer should learn to

understand TMs (better) solns. — but this may be impractical & too time consuming.

The point of this "Trainered" phase of TM's education is to get to the point where a trainer can design training problems based on a preconceived needs of TM's mind — Just as a trainer does w. human students.

Can we Teach BU search using LSuch? I am thinking of BU search for OZ or IND problems as starting with a TM or maybe a set of "noisy" solns. soln that is not very good, — but it has a $G > -\infty$ (e.g. in

case of IND, $2pc > 0$). What TM has to Learn, is an Alg. for BU search, & this can be of a reasonable CJS. One trouble with this way is that such alg. levels (coits) are expensive to test.

Another (perhaps unnecessary) objection is that such algos are usually invented via logical analysis rather than simply doing algm. trials in a pc ordm. My guess at this pt. is because I haven't worked out details of how TM is taught to do logical reasoning. I had this idea that TM first learns academic, abstract logic, like it does Algebra. Later we teach it to relate these concepts or logic to problems in a M in which reasoning is needed. So TM is eventually a logic to reason logically about any thing in which such reasoning might be relevant.

4TM

S-Functs :oo

Also 4TM 20.00 &!
2.74.07-40

302-40 is imp

on S-functs: One good way to get S-functs: Take a common "primitive" set of S-functs
A take of functs of Perm: eg. Say one "primitive" S-funct is uniform d.f. Perm 0 to 1.

Take T. Gaussian D.F. - integrate it so it goes from 0 to 1. Then, if we have
a uniform D.F. from 0 to 1, the inverse of this will give us the usual Gaussian D.F.

Any P.F. from -oo to +oo can be expressed in a "fancy" way ... via its integral.

What about discrete d.f.'s? If we have a discrete d.f. of k pts uniform betw 0 & 1, we
can still use the inverse of the S of any d.f. from -oo to +oo. - I'm not sure that this is quite
(What I want)

I'd like a way to combine P.D.'s so I always get normed P.D.'s - Consider how P.D.'s are
combined, interact in R.W.!

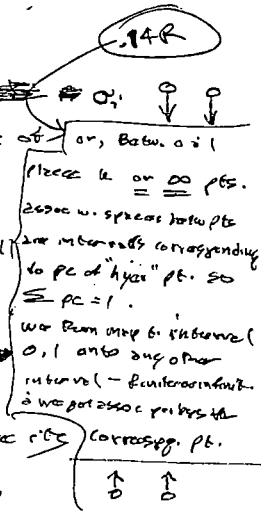
Given a uniform D.F. on k discrete pts, ~~the~~ uniformly placed betw 0 & 1. - we can
map Perm onto any other interval & preserve normality, since each ^{of the k} pts will always have
pc = 1/k.

How can I get non-uniform d.f. on k pts? Perhaps use \sum form, so each pt. $p_i \leq 1$
pc; of all pts of (assess) - So last $p_i = 1$. Any monotonic trend betw 0 & 1 will do this.

Then we map those k pts onto any ^{other} set of k pts. That we like. - So this gives us
poss. \rightarrow finite discrete D.F.'s.

We can have an ∞ of pts betw 0 & 1, & have a D.F. on Perm viz. 14-15
This ∞ of pts. can Perm be mapped onto any finite or infinite interval.

T. fancy, so nice! I want to develop it more, - more examples: Make sure it's
& separate ... can represent all I want/need. The "semi measure of ACP
shouldn't give trouble ... it is normed if we allow: "0" \in Yurbil \rightarrow 308 30



Another source of pc's is Tom Etter's "Join" function in EXEL. See my files on
Etter's stuff; I ~~very~~ have written a bit on it, "bit in Nolog on Etter": But much more elsewhere ... May be TM files??

Also note 305.00-01! This is a ^{Marked} ~~Perm~~ _{Perm} way, prodn; looking for useful contexts; making a

lang. to generate contexts: Lang consists of "parts" of various types and "combination rules" of various kinds
Both are expandable. PPM is a first-order tech way to do this. So a TM way of problems of
this kind could be to improve PPM! - A PPM is a component of TM2.

Another common typed start is a set of clusters: T. output of donee into box & search cluster center
& a distance function from Perm Center. The search could be easy (for a no. or a vector).

The distance function could be at first certain standard types, then more tuned to the functions
SSZ & also search results "make". We can also have >1 center: each with own pc

The most general case has many centers (or assoc. pc's) & simply pc's of non-center Perm
rapidly \downarrow w. distance from Centers, This is in direction of "3 input vnc" \rightarrow Note 308.30ff

A Methodology for S-funct invention: Give TM pairs in fsc. For each S-problem & types
devise an S-funct form that works for this kind of problem. After many S-funct been
types have been used, try to generalize to set to make it universal. \rightarrow (SPEC 308.00)

00: 301.29: In my other Maxm writings, I write $\frac{n}{n-m} G^2 \rightarrow \frac{n}{n-m} \frac{n+1}{n+1-m} G^2$

Unclear if this is meant to be for R window or X window.

If it's for R window, I expect the expression to be ∞ for $m=n$ (which it does). but

it is ≥ 0 for $m = n+1$ — which is not so reasonable! — The perhaps one should not consider $m \geq n$ at all!

Anyway, it would be relatively easy to get this experimentally, using n (to ~~test~~ PGM)

as $\mu(X_n) = \sum_{i=1}^m a_i X_{n-i} + \text{Gaussian noise}$.

I could test it PGM. in noise exp. Since the matrix is Toeplitz, it should take only $O(m^2)$ rather than m^3 steps. Num Recipes has PGM for Toeplitz inversion.

10) For, I may want a more general matrix inversion PGM: partly "in general":

partly for regression of the form ϵ simple non-linear.

perhaps write PGM in Causal computer (which I want to learn, anyway) — so ~~can~~

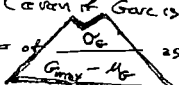
13) has a matrix inverter Col unknown Goodness (?), \rightarrow 303.00

Also Note 58.06-22 on "Grand Plan" for using GP for 28.6 TM

20: 297.40: So I. perhaps imp problems for a quick (maybe Phaset) TM!

1) T 3+1 areas of difficulty of GP on 272.25-29

22 1) Faster methods of finding good solns: Use of PGM w/ Lsearch: Better Models: as results of TM looking at problems suggesting model classes. Also Attempts to approximate Gove as a function of cond params, when ^{True, exact} Gove is very expensive

In general, try to ~~do~~ do this last, any way (even if Gove is cheap), so, with a direct form of Gove, may be able to more quickly find peak! Use of  as Gove for 2' (maybe for reference).

2) a) Profit from solns to previous problems; b) be able to do tsq's: Also note 22.6

3) Beyond Gen. Alg. Be able to watch Gove in action so as to understand why a cond was good/bad. ^{searchable}

4) Get good, useful, searchable forms for ϵ -funct.

From G&R S's P337:

$$(n-1)!! = (n-1)(n-3)(n-5)\dots$$

$$(1)!! = 1$$

$$\int_{-\infty}^{\infty} x^{2n} e^{-\frac{x^2}{2\sigma^2}} = (2n-1)!! \sigma^{2n+1} \sqrt{2\pi} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} x^{2n} e^{-\frac{x^2}{2\sigma^2}} = (2n-1)!! \cdot \sigma^{2n}$$

$x^0 \rightarrow \sigma \sqrt{2\pi}$
 $\int \frac{x^2 e^{-\frac{x^2}{2\sigma^2}}}{2} \rightarrow \frac{\sigma^3}{2} \sqrt{2\pi}$
 $\int \frac{x^4 e^{-\frac{x^2}{2\sigma^2}}}{4} \rightarrow \frac{3 \cdot \sigma^5}{4} \sqrt{2\pi}$
 $\int \frac{x^6 e^{-\frac{x^2}{2\sigma^2}}}{6} \rightarrow \frac{3 \cdot 5 \cdot \sigma^7}{6} \sqrt{2\pi}$
 $\int \frac{x^8 e^{-\frac{x^2}{2\sigma^2}}}{8} \rightarrow \frac{3 \cdot 5 \cdot 7 \cdot \sigma^9}{8} \sqrt{2\pi}$

$\exp\left(-\int \frac{-x}{\sigma^2} \left(\frac{x^2}{2} + \frac{x^4}{4}\right) \dots\right)$
 $\exp\left(-\left(\frac{\sigma^3}{2} \sqrt{2\pi} + \frac{3\sigma^5}{4} \sqrt{2\pi} \dots\right)\right)$
 $\approx \frac{e^{-\frac{\sigma^2}{2}}}{e^{-\frac{\sigma^2}{2} \sqrt{2\pi}}}$

Second Moment of numerator ($\bar{x} \cdot \bar{y}$) is $A_2 \cdot A_2 = \sigma^6 \cdot 2\pi \cdot n$ (299.35)

First moment of denom = $\bar{x} \cdot \bar{x} = A_2 \cdot n = \frac{\sigma^2 \sqrt{2\pi} \cdot n}{2}$
 Second " " " = $\bar{x} \cdot \bar{x} = A_4 \cdot n = \frac{3 \cdot 5 \cdot \sigma^4 \sqrt{2\pi} \cdot n}{8}$

First Moment Num = $A_2 \cdot n = \sigma^6 \cdot 2\pi$
 Sec Moment Num = $A_4 \cdot n = \sigma^6 \cdot 2\pi$

Value of denom = First Sec Mom - (First Mom)² - **Whoops!** (divided by 200th moment) $\ll ?$

Another whoops! - In .01 if I made a normalized d.f. for Gauss. So I should divide many nos. by $\sigma \sqrt{2\pi}$.

$\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} = \sigma \sqrt{2\pi}$	$x^2 \rightarrow \sigma^2$
	$x^4 \rightarrow 3\sigma^4$
	$x^6 \rightarrow 3 \cdot 5 \cdot \sigma^6$
	$x^8 \rightarrow 3 \cdot 5 \cdot 7 \cdot \sigma^8$

First mom. of denom = $\sigma^2 n$
 Sec " " denom = $3\sigma^4 n$
 Sec Mom of Num = $\sigma^4 n$ (here σ is the S.D. of the Normal Gaussian Noise)
 Sec Mom of denom = $(3-1)\sigma^4 n = 2\sigma^4 n$

$\frac{1}{\sqrt{2\pi} \cdot \sqrt{2\sigma^2}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} \cdot \frac{x^2}{2} = \frac{3\sigma^4}{2} = \sigma^4$

instead of the complicated multiplication of $200 \cdot 30 \pi$, consider the mean value of

$\frac{1}{1+x}$ in which x has d.f. $\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$
 $= \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{e^{-\frac{x^2}{2\sigma^2}}}{1+x} dx$
 $= \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} (1+x+x^2+x^3+\dots) dx$ odd terms = 0.
 $= 1 + \sigma^2 + 3\sigma^4 + 3 \cdot 5 \cdot \sigma^6 + 3 \cdot 5 \cdot 7 \cdot \sigma^8 \dots$ (see .01 R)

Somehow by W. Forgy: 1) $\frac{1}{1+x}$: when split by its numerator, is statistical noise with lower correlation

2) T. Num & denom. are highly correlated.

29: 2 imp. ideas: 1) How TM problem is already solved: 2) 299.00 W (who has worked well at past + Gauss to 302.00)

so: 2.99.40: $\int_0^\infty x^m e^{-\frac{x^2}{2}} dx = 2^{\frac{m-1}{2}} \cdot \left(\frac{m-1}{2}\right)!$ so for $m=2$: $\frac{2^{\frac{1}{2}} \cdot \frac{1}{2}!}{2^{\frac{1}{2}} \cdot \frac{1}{2}!}$
 $m=4$

$\frac{m=4}{m=2} = 2 \cdot \frac{3}{2} = 3$. Yitz!

In fact sum of (∞) ! for $2n=m$: $m=1 \frac{1}{2 \cdot 2} \sqrt{\pi} = \frac{1}{2} A_2$
 $m=2 \frac{3}{2 \cdot 2^2} \sqrt{\pi} = \frac{3}{4} A_4$

$\frac{A_4}{A_2} = \frac{3}{2}$ ← seems true!

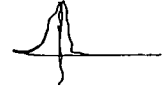
Thourcing an error = memo: top of p337 G & R table of S's

My representation of Gaussian D.F. by 2 pts at $\pm \sigma$ may not be good!



$\int_{-\infty}^{\infty} \frac{dx}{1 + e^{x^2}}$

$\int_{-\infty}^{\infty} \frac{dx}{1 + f(x)}$



$\sim \frac{1}{(1+f)} \sim 1 - f + f^2 - f^3 \dots$
 Integ of even terms only.

$= \int_{-\infty}^{\infty} 1 - f^2 + f^4 \dots$

$\int x^m e^{-\frac{x^2}{2\sigma^2}} \rightarrow \sqrt{2\pi\sigma^2} \left(\frac{m}{2}\right)! e^{-\frac{m}{2}}$ (30)

$\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx = \sqrt{2\pi}\sigma$

Paraphrase Σ ! try actually (linear regression w. known σ^2 of generator, Σ case

If σ^2 is correctly predicted & 'next case' is accurately predicted

try to see if Akhiezer's formula is better or worse than $\sigma^2 = \frac{u+m}{u-m} \cdot \sigma^2 \frac{(u+m)}{1-\frac{m}{u}}$

i.e. compare $\sigma^2 \frac{1+\frac{m}{u}}{1-\frac{m}{u}}$ w. $\sigma^2 \cdot e^{\frac{2m}{u}}$

$\frac{1+\epsilon}{1-\epsilon} = (1+\epsilon)(1+\epsilon+\epsilon^2+\dots)$
 $= 1 + 2\epsilon + 3\epsilon^2 + 4\epsilon^3 + \dots$

$\frac{1+\frac{m}{u}}{1-\frac{m}{u}} \sim \frac{1+\epsilon}{1-\epsilon} \text{ v.s. } e^{2\epsilon}$

$e^{2\epsilon} = 1 + 2\epsilon + \frac{4\epsilon^2}{2!} + \frac{8\epsilon^3}{3!} + \dots$
 $1 + 2\epsilon + 2\epsilon^2 + \frac{8}{6}\epsilon^3 + \frac{16\epsilon^4}{24} + \dots$

say $\epsilon = \frac{1}{2}, \frac{1}{3}, \frac{1}{4}$
 $\epsilon^3 \frac{1}{8}, \frac{1}{27}, \frac{1}{64}$
 $\frac{2}{3}\epsilon^3 \frac{1}{12}, \frac{1}{40}, \frac{1}{100}$

So they differ at 3rd order! & much at 4th order.

Thus, m can be close to u so $\epsilon \rightarrow 1$ say, and $\frac{1+\epsilon}{1-\epsilon} = 9$ $e^{1.8} = 6$ (not far off)

So actually $\frac{1+\epsilon}{1-\epsilon} \approx$ good approx of $e^{2\epsilon}$ $\epsilon = .5 \left| \frac{1.5}{.5} = 3 \right. e^1 = 2.7$

$\epsilon = .9 \left| \frac{1.9}{.1} = 19 \right. e^{1.8} = 6$ for ϵ close to 1, the 2 methods give very different results.

0.12

try $\prod_{z=i\Delta}^{\infty} \frac{1}{1+f(\Delta \cdot i)}$ $f(x) = e^{-\frac{x^2}{2\sigma^2}}$ $= \exp \sum -\ln(1+f(\Delta \cdot i))$

$\ln(1+\epsilon) = \epsilon - \frac{\epsilon^2}{2} + \frac{\epsilon^3}{3} \dots \approx 49\%$

$\frac{d}{dx} \ln(x) = \frac{1}{x}$, $\frac{d^2}{dx^2} = -\frac{1}{x^2}$ at $x=1$ $\ln' = 1$, $\ln'' = -1$

No! it's $\prod \frac{f(i\Delta)}{1+i\Delta} = \exp \left(\sum_{z=i\Delta}^{\infty} (\ln f(z) - \ln(1+i\Delta)) \right)$

No! to give 2 values ≥ 2 small w $m \geq 2$ produce, mult by $\geq \epsilon$, w. small ϵ .

So maybe $\exp \sum_{z=i\Delta}^{\infty} f(z) \cdot \ln(1+i\Delta) = \sum e^{-\frac{z^2}{2\sigma^2}} \left(i\Delta - \frac{i^2 \Delta^2}{2} + \frac{i^3 \Delta^3}{3} \dots \right)$ rather messy

$\exp \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} (-\ln(1+x)) dx \left(-\frac{x^2}{2} + \frac{x^4}{4} - \frac{x^6}{6} \dots \right)$

10:27-89.40 :

Consider 289.30: given y_i a model $z_i = \mu$ constant to find μ , $\mu \in \mathbb{R} \subseteq \sum_{i=1}^n (y_i - z_i)^2 = \min$.
This gives $\mu = \frac{1}{n} \sum_{i=1}^n y_i$. If we assume $[y_i]$ was generated by a model of mean μ & $\text{Var} = \sigma^2$.

1033 What is good estimate of σ^2 ? Working Backward: If we know μ & $\text{Var} = \sigma^2$,
What is most likely (max) σ^2_{obs} for n -cases? For true $\mu=0$, M_{obs} will be $\mu \pm \frac{\sigma}{\sqrt{n}}$

σ^2_{obs} will be $\sim \frac{\sum X_i^2}{n} = \left(\frac{\sum X_i}{n}\right)^2$. T. first term is $\sim \sigma^2$, 1. second $\frac{\sigma^2}{n}$ so $\sigma^2_{obs} \approx \sigma^2(1 - \frac{1}{n}) = \sigma^2 \frac{n-1}{n}$
so $\sigma^2 \approx \sigma^2_{obs} \frac{n}{n-1}$.

Try for regression eqn for 289.32: True model is $X\alpha = y + \epsilon$ (noise about zero).
First try simplified version of 1033 say $\alpha = 0$. What is α_{obs} & σ^2_{obs} ?

Find $\alpha \Rightarrow \sum (\alpha X_i - y_i)^2 = \min$

~~$\frac{\partial}{\partial \alpha} \sum (\alpha X_i - y_i)^2 = 2 \sum (\alpha X_i - y_i) X_i = 0$~~
leads to center around zero. $\alpha \sum X_i^2 = \sum y_i X_i \Rightarrow \alpha = \frac{\sum y_i X_i}{\sum X_i^2}$
So, for normalized \tilde{X} , $\alpha = \tilde{y} \cdot \tilde{X}$ = projection of \tilde{y} on \tilde{X} . $\tilde{X} = \frac{X}{\sqrt{\sum X_i^2}}$

To get estimates of α & of σ^2_{obs} we need more info about \tilde{X} .

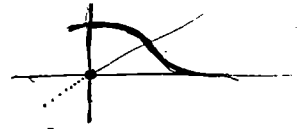
One kind of constraint: y_i & X_i are from stationary Time Series, say 1. seq. is "stable" so $|\alpha| < 1$.
In IIR α_{obs} would be ratio of 2 autocorrelations, $\Delta t = 0$ for numerator; $\Delta t = 1$ for denominator.

say $\alpha = \rho$ is seq. is stationary. $\tilde{X} \cdot \tilde{X}$ will not be \propto proportional to n if $\rho \neq 0$. Since y_i random X_i is random
w/ X_i , $\tilde{X} \cdot \tilde{y}$ will be \propto product of 2 random Gaussians: X will be symmetric about 0 & random w.r.t. y ; y will be random Gaussian if 1. seq. is stationary.

so $\tilde{X} \cdot \tilde{y}$ will be \propto normal D.F. w. $\sigma^2 \propto n$.
oops! X & y will both be gaussian d.f.s of $\sigma_x^2 = \sigma_y^2 = \sigma^2$ of n noise, so we know what d.f. of $X \cdot X$ and d.f. of $\tilde{X} \cdot \tilde{y}$ are. for IIR, we can (a.s. by?) find first 2 sec. moments of $\tilde{X} \cdot \tilde{X}$ is known.
to derive, $\tilde{X} \cdot \tilde{X}$ is first mean is $\propto n$; say $\alpha = \int_{-\infty}^{\infty} z^2 \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} dz$ from $\tilde{X} \cdot \tilde{X} = n \alpha$.

ed. Fraction

C (Coulby) / Boon Chesnut hill \rightarrow route 9 - Hystop. 166 :



All that I need is first and second moments of $X \cdot X$ & $X \cdot Y$ for $n=1$.

first mean for $X \cdot X$: $\int_{-\infty}^{\infty} z^2 \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} dz$ (A2) second moment: $\int_{-\infty}^{\infty} z^4 \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} dz$ (A4)

$n=1$: $X \cdot Y = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy e^{-\frac{x^2}{2}} e^{-\frac{y^2}{2}} dx dy = \frac{1}{2\pi} \int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2}} dx \int_{-\infty}^{\infty} y^2 e^{-\frac{y^2}{2}} dy = \frac{1}{2\pi} A_2 \cdot A_2$

from 301.01: $A_2 = \sigma^2 \sqrt{2\pi}$
 $A_4 = \frac{3}{2} \sigma^4 \sqrt{2\pi}$

for $\alpha=1$ first moment, $\sigma A_2 \cdot \sigma = 0$ for $\alpha=2$ it's $\int x^2 e^{-\frac{x^2}{2}} \cdot \int y^2 e^{-\frac{y^2}{2}}$

so first mean of $X \cdot X = A_2 \cdot n$; second moment $X \cdot X = A_4 \cdot n$
" " " $X \cdot Y = 0$ " " " $Y \cdot Y = A_2 \cdot n$

Auff approx: T. 2 second moment of numerator is divided by

$\frac{1 \pm \frac{\text{some constant}}{n}}{\sqrt{n \cdot \text{width of Baffor.}}} = \frac{1}{\sqrt{(1+\frac{\epsilon}{n})(1-\frac{\epsilon}{n})}}$
 $\approx \frac{1}{\sqrt{1 - \frac{\epsilon^2}{n^2}}} \approx \frac{1}{1 - \frac{\epsilon^2}{2n^2}}$

If $\geq 4\sigma$ var V , then $k \geq 4\sigma$ has var $k^2 \rho$

The unknown decim is $X \cdot X = A_2 n \pm \sqrt{A_4 n} = A_2 n \pm A_4 \sqrt{n}$
 $= A_2 n (1 \pm \frac{A_4}{A_2 \sqrt{n}})$ (can't use for $n=1$, or perhaps ≥ 2)
 $\frac{X \cdot Y}{X \cdot X}$ has var $\frac{A_4^2 n}{A_2^2 n^2 (1 \pm \frac{A_4}{A_2 \sqrt{n}})^2} = \frac{A_4^2}{A_2^2 n} (1 \pm \frac{A_4}{A_2 \sqrt{n}})^{-2}$
 $\approx \frac{1}{n} - \frac{A_4}{A_2^2}$
It would be nice if $A_4 = A_2^2$!

PS: Ratschinski Occ Razor in Pp 97. PS. 19/06/04 : For discussion of Source of April

Akaike AIC 8/28/04

$$AIC = \left(\ln \left(\frac{1}{n} \sum_{i=1}^n \sigma_m^2 \right) + 2m/T \right) \quad \text{or} \quad \ln \left(\frac{1}{n} \sum_{i=1}^n \sigma_m^2 \right) + 2m$$

or $\ln \left(\frac{1}{n} \sum_{i=1}^n \sigma_m^2 \right) + 2m$ (Total sq. error)

T = no cases = n; m = no. of coeffs in (linear) regression. $\approx \min \sigma_m^2 (1 + \frac{2m}{T})$

Since T is constant over n

on the other hand, I was very

$$\sigma_m^2 = \frac{T+m}{T-m}$$

single parameter model = $\sigma_m^2 \frac{1 + \frac{2m}{T}}{1 - \frac{2m}{T}}$

We get: taking ln we get:

$$\ln \left(\sigma_m^2 + \ln \left(\frac{1 + \frac{2m}{T}}{1 - \frac{2m}{T}} \right) \right)$$

$$\ln \left(\sigma_m^2 + \frac{2m}{T} \right)$$

So the AIC is using k to expected error for next product is min.

So: How to consider

$$\left(\frac{n+m}{n-m} \right) \cdot \sigma^2$$

with

$$\left(\frac{1}{\sigma} \right)^2 \cdot \left\{ \begin{array}{l} \text{volume in } m+1 \text{ space} \\ \text{coeffs plus } \sigma \end{array} \right.$$

$$\text{or } AIC \rightarrow \sigma^n \cdot e^{2m}$$

$$\rightarrow (\sigma^2)^n \cdot e^{2m}$$

$$\sigma^{-1} - \frac{m-1}{2} \text{ necessary increase in}$$

$$\left(\sigma^{-1} \cdot e^{\frac{2m}{T}} \right)^n \gg 1 \text{ & } \text{volumes to be compared.}$$

could not get σ^2 at x_i "same place" I got

"spurious" e^{-1} in ZIHI analysis

NB

Perhaps m.c. decomposition $\frac{n+m}{n-m} \sigma^2$, we make some assumption about a priori's?

result is $\sigma^2 / 2$ function of $\frac{m}{n}$, a dimensionless constant.

Actual m is $\frac{1}{2}$ may be dimensionless constant.

Consider coeffs: for m coeffs, there is a region in n space where σ^2 is within σ_{max}^2

$\sigma_{max}^2 \propto (n+m)$. For linear regression, the coeffs are all dimensionless. So in theory,

one could compare k volumes of different dimensional spaces. Also note that coeffs tend to be "u" (so maybe $\frac{1}{m}$?) for linear regression.

This may be ok for linear regression, in which all coeffs are of about the same significance. But in linear Curfit, they are not so. Then we can curvilinearize the basis vectors

so perhaps coeffs tend to be $\sim \frac{1}{m}$.

A big Q is: Does this / m dim space approach give same results as $\sigma^2 \frac{n+m}{n-m}$?

It works done: for given m, there will be a vector \vec{v} that give min σ^2 . Given new by values of \vec{v} will give smaller σ^2 .

A single case is a set of coeffs plus a σ^2 . for each set of coeffs there will be a σ^2 that gives max pc of corpus.

We can assume a uniform prior for σ or for σ^2 (I don't know if it makes much difference).

We have $n+1$ data pts: t cases y is f our predicted. To get E.D.f for f now p , integrate over all $m+1$ space.

Another way to define product is to consider f to pc of corpus. all pts in $m+1$ space. Each point gets a certain W.B. of f pc. if gives to corpus.

Each pt. gives a product for every possible value of f (and p) & to get a product of f .

evaluate set w/ average over points product. w. prior w/ σ^2 . we can also use sums of pts to give product. for each set of coeffs, we can assume got with Σ of all σ^2 values in predicting n past data y , a set σ^2 is a w.t. for that set of coeffs... which can then be used to get w.t.d. product.

CRITICISM of

Akaike formula!

It does not blow up

for $m = \frac{n}{2}$.

$\sigma^2 \frac{T+m}{T-m}$ does blow up

so it should.

FTDY

QUICK ABOUT: .25

10 In Sol 78: I had proof that Cover's method of induction was worse than ALP by an ∞ factor (a "very small" ∞ , hvr @) — This was for a version of Cover's method that used Gao's method of Norman. I recently chvd ~~over~~ for ALP, & Norman constant is always bounded by 2 ~~acounting for true~~ model of data. ~~It~~ Id .02 .3 for us for Cover's method, then ALP must be better than Cover by this "small" factor. (?)

Check this out!

On my web page, I could write "Comments" on Sol 64 re: Levin's analysis: on "non-convergence". Also his "discovery" that prefix sets must be used to get good PC values. Also remarks on Sol 85: That equation for when singularity occurs: ~~is not correct!~~ new eq. was derived & that its exact closed form soln is. And more recent stuff on "Modre's" law. A (Bon Neuman's Remarks 2 book "experiments" — Explain what it means by Modre's law i.e. holding of cc per bit every 1 year 1 1/2 yr. & no limiting size or needs by host.

Also how singularity occurs if amt of money invested each yr. is bounded.

0 :294.80: pc's of a "best" can't thus far (when they had gotten as far as that particular problem) "Satisfactorily" LET'S look at conventional QATM: say it has "solved" 10 probs "Satisfactorily": so it has at least one funct that can do all 10 probs w/ "reasonable pc" for f. entire set of 10. we give it Q.A. —

23 we want O^J for all n problems. We just try total soln for (QA)_n along with list of satisfactory pc. (25)

24 **SN** for each QA problem in a young TM: Train gives a pc level that it MUST achieve ... as minimum.

25: (23) After solving (QA)_n we test it again previous problems: f. lowest ~~of (QA)_n~~ ^{previous} ~~max~~ ^{min} problem first

27 and then we test it on all other probs in order of diffy of problem. We fail if ~~it~~ ^{product} of probs of "solus plus far" is "too low". ← (was ~~some~~ this "too low" level may be set by Warner also, or ~~be~~ ^{be} deemed from ~~some~~ thresholds given by Sumner (n.24)).

28 So f. target is 24 L such soln, we will also use R (recognition function) Mechanism of IDSA report.

0 My impressn is that f. "Quick about" idea of .25-.27 would quickly reject most cands that would eventually "fail". Hvr, in General, we don't want to have to test a new O^J on a

to entire page corpus. Usually f. modifc in O^J will only affect recent problems. ———

33 We should be able to take ad vantage of this fact (i.e., indeed, it is true!). Clearly would be the "R" method. Could we somehow arrange for it to Discover f. "R" method? ... Suitable T.S.Q.

Also note that usually L such will not descr. Quick About, because method of f. P.B. does not (normally) enable Quick about. A Phase II Machine will look for good QPZM Methods & Phase must be able to consider all over such time & .i. ~~not~~ be in position to discover Quick Abouts. → 302.20

Notes
309.00
302.20
302.20

4TM

ANN .00

SM: Median Filters : 80

293.00

10. ~~293.00~~ : Hvr. one could get second derivative (Array/Matrx/Tensor) & solve ~~matrix~~ eq. to get best jump size.
 T. Costs of Tensor can be obtained partly by chain rule, partly by making several trials in various directions.
 Also, one might assume all off diagonal elements are zero, or, that only f_i elements in Dim 1, 2, 3, ...
 indices of diag. are non-zero; Or other assumptions, methods, to reduce no. of elements handled.
 In a standard Non-linear Curvature technique, I think they assume all diag. elements are zero.

0 : SM Median Filters : used "Min Median error" should be good for SD: Perhaps start w. linear Models; Then use ANN. Linear models are not really linear because $e^{-|x|}$ error criterion (or $|x|$ error criterion), but we can start w. ~~Median~~ error costs to fit for Mean sq. error criterion as first approx.

In ANN, using Back prop, or - Marquardt method; $e^{-|x|}$ error produces no big change in methodology.

For "x window" a "median" of $x_{n-1}, x_n, x_{n+1}, x_{n+2}, \dots, x_{n+m}$
 would find $\frac{1}{m} \sum_{i=0}^m |x_{n+i} - \bar{x}| e^{-\frac{|x_{n+i} - \bar{x}|}{\sigma}}$ was minimized.

To obtain this $\frac{1}{m}$, use successive approx, starting w. \bar{x} . It might be possible to show that in successive approx, f_i errors always \downarrow , so i with most bias has largest error. This would enable one to know which of past data one had to retain (remember) in order to make corrections (over!).

Try using $\frac{1}{m} \sum_{i=0}^m |x_{n+i} - \bar{x}|$ (non-saturated) v.s. $e^{-\frac{|x_{n+i} - \bar{x}|}{\sigma}}$ (non-saturated) [A trials] for Median prediction. — Empirical tests. Use generator of data in which b "noise" part has $\geq e^{-|x|}$ distribn.

4TM Rev:

20 : This is a preliminary backhaul review from ~ 288 backward.

277.20 to 288.20 + variance on 291.00-.18

275.00 - 277.19 : is a "Go" at the 2141 paradox; Not carefully done, but may be useful if I want to do it more carefully.

04 274.07 - .40 Mainly on Problem of finding Good ways to represent S-functions; Some good suggestions: Also 274.25: Links this problem to TSC writing in a very useful way! This seems very important (particularly the link to TSC's). It's an important part of the general problem of the form of TM & its TSC.

- 273.26: Lists 3 Big Problems: ① How to use solns to past probs to help solve present problem.
- ② Lsearch v.s. BU (Bottom Up) search: It was to be simultaneous, how to find concs that may both can use.
- ③ how to represent S-functions: (see 104) is 274.07 - .40 for good ideas on this.

→ 272.30: If I had 2^k ways to describe TM, for $k=1, 2, 3, \dots$ how would I do it... say for small k ?
 I think this remark started off the stuff that followed thru 2141 to 288.20
 The "3 Big Probs" of 104 were an important part of that thread - but I guess I forgot that I was trying to solve ~~272.30~~ 272.30! I eventually got to the "Lsearch v.s. BU search" question.

Note 25706-22
 A "Grand Plan"
 for TM using GP.

Good! The stuff leading to 272.30 is nice!
 272.21 reviews the 6 needed improvements of GA - Note "3 areas" of 272.29 1/2

From 265 onwards, there was discussion of ways to represent S-functions:
 Before 265, Discussion of Lsearch v.s. BU (Bottom Up & GA) search.
 (266-267) is on Universal Analy. Machines

"What has worked well in Past" : 1.00 - 1.19

FTM

What has Worked Well in Past

This is a general method of getting & updating ones' exps.
 We have 2 particular algms, say, ^{task} _{assignment} that calls how to get exps for present problem based on
 success/failure of ~~param~~ ^{param} & prop/meth in previous problems. (Getting 2 of d_2 of
 -1.36 would be an example). Each ^{Poss.} algm. will \downarrow bc of corpus, until one adds in bc of }
 f. d. cor. of: "worked well in past algm". So one must first measure just how much use of
 that algm \downarrow bc of corpus.

Particular important v. score of .00 ff! : Geometric Probability : how to compare
 n volumes in different nos. of dimension of param. space (e.g. -1.34)

I want to make sure that I have an general, unambiguous, & race, & effective, formal, practical

Way to do this: (14)

Perhaps Give ^{Institute} Paper on .00 ff! "T. Univ. Dir and Appr'd." : Agent in Theory and Practice.

The Computation of A.P.P.S
 Computing A.P.P.S in Theory and Practice. In Search of A.P.P.S: - Theory and Practice

Also Discuss Occ razor: When it works & when it doesn't work (reference papers) to when it doesn't work).

See PS: 19/06/04: "Occ Raz in Phys":

T. problem of .00 is closely related to problem of creating trials for "test tasks", using concs that
 "have worked well in the past". T. usual problem in .00 is that we know for go often just when how
 we created the models if we are using & we don't remember the param assignment algms that should
 have been "Bundled" w. the same models. Normally when a ~~new~~ model is created, it's tested on 8. post \rightarrow so
 see if it's a good model. If it is a good model (i.e. we decide to accept it), then we must have had
 techniques within the model description just to id how to test - which is equiv. to knowing how
 to use that algm for production - either on prevalent old data (verify the model), or on new
 data (not yet verifiable model)

quick
 "A. quick take" on "T. problem of A.I. has been solved".

Consider a problem of relating new ^{G.P.P.} ~~prob.~~ prob. to old: T. demote latest problems compared in terms of previous
 prob. that have been solved. T. problems that are "close" to present problem (all past problems have their
 terminal
 & final pop. data stored on disc) have their terminal population mixed to form a initial population for
 new problem. There are certainly other, possibly better ways to do this. Weighing them &

assigning param. to them should be done by .00 - 1.0

For a closeness metric we can use PPM compression, or better (usually slower) compression

methods.

Main by Q's ~~First~~ Notes: 272.25-29 (More like 5 general areas of improvement \rightarrow see 272.29) "improves" in modern G.P.

Also, can it do this w/o trouble? < (This is tricky! "Can it do S-Pracs well?")

I think a main ^{Bottleneck} problem was representation of S-functs: As is my model (now) not good enough

to do LSrch for QA problems or do much ~~(BU)~~ (BU) search in GP.

Say I had a way to list S-functs in pc order (with a ^{Swy} factor of 2" of correct order):

still, I'd have problem of speed testing new cards on Enkro Computer. (Testing is not so hard for S-functs
 because one can order cases in order of pc of failure, so cards ^{fail} and quickly (retest). - Could we do anything
 - Being tried

analogous for S-functs? - Well, YES! Put problems in pc order (or LSrch order?) w/ lowest pc first.

"pc order" will be from more pc over past cards. Repeat search & card when t.

pc for entire corpus (Plus far along list of probs by pc order) \rightarrow become \leftarrow extreme threshold (perhaps to $P \leftarrow \frac{1}{2}$

for pc's ... like $T \leftarrow 2T$ LSrch). - Actually, this ~~is~~ doesn't seem to be such a hot way to do it!

A perhaps, better way: Discard card when pc for ~~problems~~ Plus far is sufficiently less than \rightarrow 297.20

Key to decision
 whether we're doing
 LSrch or BU Srch
 or Q's!

47
SPEC
290.40

Paper: "Notes On Normalization"

Given results on 1) convergence of $\sum^2 \sigma(U)$ to ? 2.1.2.k.

2) The norm converges, see 290.20-40

3) Fast Normed P.d. is always better than unnormed: $\sum (\sigma_i)^2$ is $<$ by min of σ in Normed case.

4) The "usual" norm function is best

5) Discuss "by constant factor" is important in ~~practice~~ practice (indeed because the choice of ~~reference unit~~ reference unit is not arbitrary. (No mis is dis. Consider some default) in my book 'table' (paper).

0 267.06

ANN w/o cross val. 1 Pick random set of initial cond for a net.

2 Use back prop. for k "seconds" (or other count). Do prod. w. resultant state.

3 Repeat 2 1) 2) l times.

4 Take average of ~~prod. results~~ prod. results: perhaps with average.

Q's: Total cc \approx kil. Given a total cc (2k.1) yes! How to allocate 'kil. or perhaps: to spend ^{more} cc on certain sub-runs?

Other param Q's: How many params (Nodes) to use in net? What topology -- how many hidden layers, what population of hidden layers? Perhaps in (2.2) \in (1), chose topology at random. also. No. of params could be unknown, but

small ~~no. of para~~ in (2.3) (averaging) small no. of params get more wt. (32)

These params, as well as various cones that enable one to correlate the params with initial problem form, can be solved by "what has worked well in past" (See 294.06)

28

T. advantage of ~~to~~ to ~~log~~ log ~~med~~ med over "max flatness" is that larger space of

"hypoths" is available. One could 1 no. of params in max flat, but notes "wildly" (numerous).

32

24 Another utility trick: when utility func of no. params, use error on known Corpus,

if error = 0, use "ratio" of ~~no~~ no of params space to no of prod. space of constr

34

"constr" prod. A Big Q is how to compare lvls of diffrent no. of dimens

36

One way: To make sum of lvls converge: each dim dimension has factor, α of loss ext ("hvl")

α is found by optimization of prod prod. Conceivable one can use a diffrent.

α : for each successive dimension, if one has out of "past data".

On ANN: (unrelated to immediate preceding stuff): Normaliz "backprop" they use various tricks (momentum, imp rate) to decide on stop size. T. direction of jump is known via derivatives of wts wrt. Gove. (Spec \rightarrow 296.00)

4TM

Notes on ZFC:

1) One way to explain the problem: We know this corpus: (say Alphabet is a, b, c): We have Press 2 grammars: Which is ~~more likely to be~~ 1250 of pc of each generating corpus? A sort of Bayesian problem: Th. 2 grammars: (1) simple Bern w. a, b, c .

2) is ~~simple Bern w. a, b, c, γ~~ ; ~~or $\gamma \equiv a b c a$~~ .

Th. 2 grammars do not get here pc's for each of their choices. We assume uniform prob's \in over all possys. Certain pc values will be most likely - but we are mainly interested in total pc's of the 2 grammars in view of their priors.

Essentially, we want to know the pc of detecting γ .

Re: Re "1/6 factor": 277.20 - 279.02 is a simple way to get pc of γ , but it doesn't have 1/6 factor: This treatment is incomplete because it does not consider γ in pc of other ($\neq \gamma$) symbols when γ appears pc > 0 . As soon as we introduce that ψ of pc, we get the 1/6 factor. 278.26 - 279.04 explains how 1/6 arises.

4TM REV

List of Problems: unsolved & "independently" solved.

Note
99TM 114!
2 V. p. review of work
on 2-141.

1) 2-141: The 2 versions of behavior of system wrt. varying n! 285.31, 32
Also, the general feeling that Z^R is not so reasonable.

$$Z \equiv (\alpha e^{\frac{1}{n}} - 1): \alpha \equiv \frac{R}{\rho n} \text{ where } \alpha \text{ constant as } n \text{ grows.}$$

But in my letter to Wolff I said that I had some understanding of why Z should have that form: (i.e. $\frac{1}{n}$ factor)
See 292.10-13 for discussion of how $\frac{1}{n}$ factor arises
I. main ~~the~~ recent problem was that 285.31 & 32 differed by a factor of \sqrt{n} a 'serious discrepancy'.
It seems unlikely that if approx. used in 283.37 should be off by a serious factor!

2) "How Many Cords": 286.00; 287.00 - 289.20 I got close to a soln. to this problem, but it does need work. The stuff I did in 1962 in that city is probably relevant/helpful.

3) f. problems of .01 areas when I felt that "coding in derivatives" did have serious advantages over the method used in PPM with ~~frequency~~ ^{some} 278.20 it was an attempt to put in more exact form the idea of 278.23-25 (i.e. that 277.20 - 279.22 (278.00 - 22 in particular) solved the "1/n factor" problem in 2-141 (apparently it didn't) & the problem of making deltas of news w. unusually low (or zero) frequency. I think this second idea didn't have very much to do with the 1/n factor! See 292.10-13 for discussion. ^{280.18-19 is imp!}
Wolff's: see 279.10 for the beginning of PPM. Continued to 280.19 (SN on this: would parallel coding of "unusually low freq news" contribute further to the difficulty of this idea?)
Even unusually high or low freq. of any news must lead to a shorter code; 275.15-16

4) In my Summary over all passl phases of X (in 2-141), I used the factor $\frac{R!}{(R-m)!m!}$ as the no. of passys. Should it be $\frac{R!}{(R-m)!}$? (R is the no. of lines)

X occurs m times; m is the no. of times it is considered to be X, $0 \leq m \leq R$.
Also, I ^{only} summed codes w. m before

20
21

10

18

20

20

MEM (Max Entropy) .00

Normalized : 20

10 : T: 2 versions of MEM ^{computer} Max Ent. & Max Cross Ent. may both be known!
 Terms of ALP.

For ^{simple} Max Ent. : Statement: we have a seq. w/ $P(X_i)$ being ^{unknown} function of part of X.

Also constraints on $P(X_i)$. MEM ^{is in fact} in this situation, most likely $P(X_i)$ is ^{specified by} part for which $\sum \ln P(X_i) \rightarrow \max$ This is seq. Part ^{is} known to a part of seq.

04

How .04 is not same as .045 (for history!). While min of .04 is correct, it ^{is} not Max (clearly not); .045 is not ALP, but is less clearly stated. - But it implies .04 is unknown.

At any rates, if ALP trivially implies .045, then if $P(X_i)$ is a function both to part of X_i & to part of X_i (a seq) then .045 is not Max Cross Ent. Version of MEM.

For both versions of MEM (T: norm (MEM is not special case of cross ent)) in which $X_i = \text{constant or null}$) if .045 implies Max of .04, then perhaps to conclude of ALP convergence theorem (Sol 28) holds! E.g. min-max dot. is min MS error via Parse methods! - I don't know if Parse facts are known. If they are known, then it may perhaps give new proof of Parse theorem.

Corollary of Conv. Thm (Sol 28): If $\frac{P(X(n))}{\mu(X(n))} > \alpha(n)$ for all sequences of length n.
 Then μ ~~is~~ (Expected value ~~of~~ $(\ln \mu)$) of \sum sq. error of cond.

pc. of successive bits, $\mu < -(\ln \alpha(n))$.

If μ is a generating D.F. & P is a Universal D.F. (Normal or not normal),

Then is that if μ had finite bits, then possible

$P > \mu \cdot e^{-D}$ for all poss. strings $X(n)$. The Normal constant (if P is normalized)

can never be $> e^D$, because if it were if normalized P would be $> e^D \cdot P$ and

$\frac{P}{\mu}$ would be > 1 , which is impossible. since μ is already normalized.

Still, I must remember that each string $X(n)$ usually has a different Normal constant.

I know, in line with .20 if $N(n)$ is the smallest normalized constant for all strings of length n (for given unnormalized D.F. $P(\cdot)$), then \sum sq. error in P_{norm} will be $<$ that for

P_{norm} by $\ln N(n)$. (In considering a "normal factor" is ≥ 1 (by data)).

032 i.p. $\frac{P_{norm}(X(n))}{\mu(X(n))} > \alpha(n) \cdot N(n)$ since $\frac{P_{norm}(X(n))}{P_{normalized}(X(n))} \geq N(n)$ for all strings of length n.

obtained by multiplying 030 by 030.

Spec
293.00

4PM

SM.00 (continued)

Spec

20: 2.99.40: On the other hand, if one considers price schemes w.o. slippage one can have very good daily yields, and still have a gain of α to justify the price scheme.

For a system that with consideration of slippage removes the slippage, computes a mean subtract the cost of the down (in units) from $\ln \alpha$ to obtain true yield: Then "subtract" slippage from that result.

The idea is that cost of down of stock composed of the slippage-free yield.

It may, however, be possible to do trades w. very low slippage - e.g. if one knows expected ~~price~~

gain of ~~very~~ very many stocks, one can "look for bargains" and buy/sell them.

At any rates, that the ~~unadjusted~~ slip-free yield $\ln \alpha$ & cost of stock down should be compared may be an important insight (return!)

It is clear that a compression in sm. data of b bits is equivalent to a $x \cdot 2^b$ slip-free yield, and that compressor decr. (any) must be subtracted from b .

For example: Consider a stock w. 2 returns: say 1 had k values, the other z values. One might run all $k \cdot z = 12$ strategies & make each bet w. α out of mult. of k & z for. This wts are not updated simultaneously is in each variant of the stock bets at different times (usually). But at each point, each stock bet (cost) contains mult. amt. of money, and that's ~~the~~ the fraction of money one bets on it. The money actually bet ~~is~~ on a strategy is a separate account from the mult. yield of the stock.

Now, note that w. $k \cdot z$ variants, if one is much better than the rest, we still have to

divide ~~its~~ its yield by $k \cdot z$. Actually, we need to keep track of individual multi. yields of each individual (strat): final yield will be $\frac{1}{k \cdot z}$ x sum of multi. yields.

Is there a way to use the info that certain params are "close to" a μ or σ ? Well our mult. divide each param. into 100 divisions & follow them as ~~best~~ best: there would be many strategies by yield near peak.

21: (2.99.20) spec: I think the true part of the problem is: for a given \bar{M} and \bar{Y}_0 How does one choose N (w. \bar{Y}_0)

vary w. size of money, N . I.e. True variance of date \bar{Y}_0 is var of \bar{N} .

But bet \bar{C}_i are ~~smaller~~ smaller over time & smaller var. for its prodng.

works! I've been working on Corfit! ~~Now is added to prodns only.~~ Now is added to prodns only.

the linear Regression noise is added to ~~the~~ data also. I'm not sure it's my problem, but.

It seems like a reference to have \bar{M} values indig. N , but if not sure it's really needed for the analysis.

Notes: Perhaps try as "study problem" & proving that $\sigma^2 \rightarrow \frac{n}{n-1} \sigma^2$ (by ~~the~~ $\frac{1}{n-1} \sum (x_i - \bar{x})^2$)

Remember that the other ~~factor~~ factor; ~~it~~ follows, perhaps from first factor: $\sigma^2 \rightarrow \frac{n}{n-1} \cdot \frac{n+1}{n} \cdot \sigma^2$.

perhaps .27 can be found on any statistics book! May be look at Regression.

I ~~try~~ try to do .27: First do T.S. w. 1 param. to approx it: so we want n of eq.

given X_1, \dots, X_n to find $\mu \Rightarrow \sigma^2 = \sum_{i=1}^n (x_i - \mu)^2 \Rightarrow \min$. Show how σ^2 of Gaussian = $\sigma^2 \cdot \frac{n}{n-1}$

do same in which ~~the~~ model is $\sum_{i=1}^n (x_i - y_i)^2 = \min$. Note that this is same as ~~the~~ 31, if

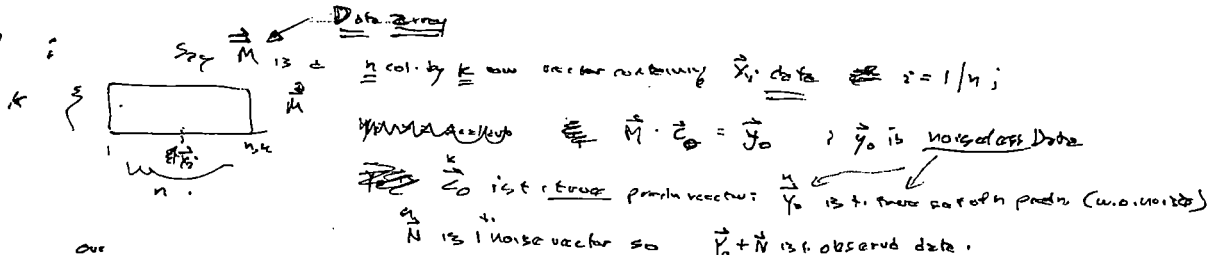
$$\sum_{i=1}^n (x_i - y_i)^2 = \min$$

Finally show it for $\sum_{i=1}^n (\alpha x_i + \beta z_i - y_i)^2 = \min$ This is presumably easily proved to

and only loss of no. of $(x_i, z_i, w_i = \dots)$ vars.

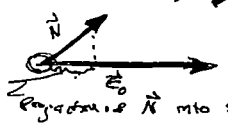
Spec 2.90.20

20:287.40



our estimator \vec{c}_0 will be $\vec{c}_0 + \vec{E}$ which $\vec{M} \cdot \vec{E} \approx -N$ as closely as possl.
 i.e. $\vec{M} \cdot (\vec{c}_0 + \vec{E}) = \vec{y}_0 + N$ w/min total sq error so $\left(\vec{M} \cdot (\vec{c}_0 + \vec{E}) - (\vec{y}_0 + N) \right)^2 = \min$ so $\left\| \vec{M} \cdot \vec{E} - N \right\|^2 = \min$.

Consider N to be an a vector in an n space: we want to approximate N by a linear combination of k vectors that are rows of M . The k rows of M form a k dim subspace; \vec{E} will be in that subspace and \vec{E} will be a projection of N into that subspace.



Getting to Volume of our set of row vectors is "Murray's problem". - Solvable by Gram-Schmidt. So to get a projection, we first make an orthonormal basis for the k , n vectors of M . From these bases, we get the projection of N . We express this projection as a linear combination.

oft, rows of M . (There may be a more direct way). We may be able to run Gr. Schm backwards on N to get the expansion of N in terms of the k row vectors of M .
 A Gram-Schmidt inspired way to get projection of N onto M : For i -th row vector M_i , remove from N : $N_i \rightarrow N_i - \frac{N_i \cdot M_i}{|M_i|^2} M_i = \tilde{N}_i$
 The scalar $\frac{N_i \cdot M_i}{|M_i|^2}$ defines projection of N_i onto M_i .
 $\tilde{N}_i = N_i - \frac{N_i \cdot M_i}{|M_i|^2} M_i$ (Normal factor)

Actually, 17-18 is a way to find \vec{c}_0 for $\vec{y}_0 + N$ 289.20

SM In auto prediction (1 stock) select model on basis of min $\sum_{obs} \frac{(n+1)}{n-1}$
 Vary both n & m (choice of model); Probly best to choose n and m which has somewhat different formula: See EMW (1960, etc):
 In cross prediction (choosing m sample size to auto predict) an additional param is "size of population" that one selected "best" of. This is related to, but not identical to t .
 SOY problem. (Spurious Optimization Yield). Its maximum in cross predn because $n \rightarrow t$.
 No. of pairs in a population is \gg + population itself.
 In using \parallel (parallel) models: does SOY still occur? - There is no optzn selection, I think.

SOY is for runs of 2 probs only. If we do SM by predicted diff, then we can avoid SOY. (Probably ... I simulate sure). It does mean we cut ourselves off from many interesting SM strategies!

T. O'Leary: that any market legit strategy can profit a pd out. corpus \rightarrow Next strat would work. Any pd. out corpus is a coding technique; it's subject to "cost of code" constraints. I think both slippage & slippage free systems can be treated.
 Any way, SOY is strategy is able to grow multi. of capital by factor K in one year using Kelly betting or whatever. Then it should be enough to compound or $\ln K$ in n yrs. Then time series away have 10k bits/yr, so 2 bits/yr that we get from strategy is quite small, so we can't have much into in f. Strategy darn.
289.00

47M

How Many Coeffs (Count) SEEMS like V.G. Fine (Soln. for linear/non-linear: Coeffs & regression)

EXAMPLE: "Needs work"!

OP: 286.96; Spec Q. Linear regression: A complete design of corpus is $\text{Dist } k \text{ values, } k \text{ coeffs, value of } \sigma^2$.
From this info, one can assign 2 pc to each value (~~for~~ $k \leq 2h$) and to any proposed $n+1$ values.
So we have a statistically set of $k+1$ params. How much info is in them?
Actually, no doubt here! We just use all sets of $k+1$ params with codes: i.e. from top to bottom

cat This may give $\sigma^2_{observed} \times \frac{n+h}{n-k}$ var about the non-predicted "counter" cases (\equiv "best" code).
The d.f. in σ^2 accounts in σ^2 spec are rather short — so reductions to apply d for coeffs &
 σ^2 are not varying rapidly. Each code has same pc but is mult by its accuracy (\equiv "it.")

07 to give final distrib. Since we sum over all codes, the cross of part part is ϕ !
073 is true for various values of k , so we verify don't care about it. We use k !
T. widom $k \leq \sigma^2$ is also unimportant. T. only impk thing is control σ^2 its width —
which is (perhaps) $\sigma^2_{obs} \frac{n+h}{n-k}$ (w. modif. for x window).

Should we give = wt. to ~~all~~ σ^2 predn for each value of k ? (for $k < n$)
∴ $(k \neq n$ means σ^2 exists & we have model w. that var.)

With each value of k has a predn of its own wt., we still have to give σ^2 value of h ~~to~~
specify each set of codes. Since $1 < k < n$ we can use current d.f.
Unless we have found in x space that for seqs of k possible types, some record of
out σ^2 values of k gives better predn.

00-17 seems to solve the VMC problem for linear regression. It will probably
work for coeffs & non-linear regn. & coeffs if σ^2 is sufficient to predn
the system is locally linear.

T. VM predn prob is $\sigma^2 = \sigma^2 \frac{n+k}{n-k}$; I have proved it for $k=1$ & its prob
true for $k > 1$. One way to do it is to consider 11 codes as us. 00 00 00

Another un proved prob (possibly not exactly true). I assume that the predn (for
 σ^2 given k dim space & σ^2 space) is rather short, so that diverging over these codes, is
very close to using the best predn, σ^2 codes. Actually, it's more complicated.
for each k vector we get a d.f. on σ^2 's — i.e. if σ^2 using ridge σ^2 's. It would seem that σ^2
 k th dim. "with σ^2 " of that d.f. should be wt. of that k vector σ^2 prob.

As I said now, for each value of k , we have a d.f. on coeffs, σ^2 spec
Say k^2 set of coeffs is σ^2 is assoc var. For each value of k^2 σ^2 we have
2. σ^2 pc for the corpus is a wt. σ^2 ; we also have a predn of the coeffs of the d.f. will have wt. w.
- pred d.f. on coeffs x wt. of the coeffs, will have wt. σ^2

So, each k^2, σ^2 will have a wt. (\equiv pc of corpus) & a predn mean σ^2 for product x wt.
we obtain final d.f. on σ^2 predn by the wt. mean of all k^2, σ^2 sets (32)
The wt. σ^2 of the bunch of normal D.F.'s is a normal D.F.? No! But in many cases, yes.
i.e. we can get $\mu \pm \sigma^2$ of any bunch of σ^2 predn of the coeffs of the d.f.

If μ_i, σ_i^2 are wt. w., prob $\bar{\mu} = \sum w_i \mu_i / \sum w_i$; $\bar{\sigma}^2 = \sum w_i \sigma_i^2 / \sum w_i$?
Maybe we have to do: for each k^2 , first give a wt. sum of all of the possible σ^2 values. Then give a single
 σ^2 (a wt.) for each k^2 . We then have a wt. of all k^2 's & assoc sum of all σ^2 's.

Needs work!

N/2/04

4 TM

How Many Coeffs? $COS \equiv HMC$: out of bot 2 67.00 - 17 seconds
 to Do **FINAL SOCN** **APPARENT**
 (or non-linear regu) \leq e.g. ANN \leftarrow (in conflict w. Barron 1993)

20: HMC: For linear regression: (Non-linear reg. is locally linear when σ^2 is small (continuous large SSZ)). T-formula $\sigma^2 = \frac{n+K}{n-K}$

Probably looks. To do a prediction \hat{y} LDL would choose k w. min of σ^2 .

ALP will wt. products all $\neq 0$ $0 < K < n$.

To show ok for $k=1$ is easy. To show it for $k > 1$, use k -dim space w. t. "correct" set of coeffs of ϵ -optimum. Any deviation from t. set of coeffs at y -optimum ($\equiv \vec{b}$) will give you a non observed error of $|\vec{b}|^2$. If σ^2 is true σ^2 for true coeffs, then

using a finite SSZ : $SSZ < \infty$ ($n < \infty$) will give a smaller σ^2 due to overfitting.

Then, since for from this observed error one can find σ^2 true parameter. What one needs, first,

σ^2 for chosen set of coeffs which is an additional cost > 0 since t. chosen set of coeffs

$\neq \vec{0}$.

$\sigma^2_{obs} \cdot \frac{n}{n-K}$ gives σ^2 of true parameter. This has to be k by factor $\frac{n+K}{n}$ to give

compensation for fact that t. set of coeffs chosen is not t. ~~best~~ set of coeffs of ϵ . True parameter

I may be able to do things more exactly: 2 ways:

1) simply vectors in n & k space & Gaussian D.P.'s.

2) In k space, it's harder to keep all poss. sets of coeffs (if correct set not ϵ optimum).

using wt. or $e^{-\sigma^2}$ not set of coeffs. I perhaps see notes of 1960, McLeod

If t. thing is a really correct ALP solve, how does it fit in with t. idea of (about some problem)

Does each extra coeff have a certain (standard) extra cost? Also how does it deal w. σ^2

extra cost assoc w. using ~~an~~ overparameterized set in k dim Σ space?

Perhaps "another" approach: linear regression as a special case of curve fitting.

For each k , there will be some ~~truly~~ optimum set of coeffs - a k -dim vector $\vec{0}$. If observed y

output is $\vec{0} \cdot \vec{X}_i + N_i$. \vec{X}_i is i 's set of k data pts.

say $\vec{A} = \vec{0} + \vec{R}_A$ (\vec{A} is an approx to $\vec{0}$). Then $\vec{A} \cdot \vec{X}_i = \vec{0} \cdot \vec{X}_i + \vec{R}_A \cdot \vec{X}_i$

say ~~random~~ random k -dimensional basis

Then we find Σ over $\vec{0}$ vs. Part of \vec{A} . We want \vec{R}_A at $\vec{A} \cdot \vec{X}_i - (\vec{0} \cdot \vec{X}_i + N_i) = \vec{R}_A \cdot \vec{X}_i - N_i$

$(\vec{A} - \vec{R}_A) \cdot \vec{X}_i = -N_i$. (I'm confused!)

On derivs of n point corpus w. k point models: T. realization wts are of k -density values in k & n space. A good code maps a good distribution (low info) in to a narrow D.P. The σ^2 's of t. D.P.'s have to be derived, but since we are only interested in the σ^2 's of \vec{p} 's of σ^2 's, it's a bit easier. We want or less k quant. "factual size" of σ^2 at output.

Still, comparing densities in 3 v.s. 4 dim Σ spaces (k spaces):

Measuring is not clear!

67.00

10 : $\left(\frac{m}{n-m\ell}\right)^m = m^m \cdot (n-m\ell)^{-m} = m^m n^{-m} \left(1-\frac{m\ell}{n}\right)^{-m} \approx \frac{m^m}{n^m} e^{+\frac{m\ell}{n}}$

Let's look at $e^{\frac{m\ell}{n}}$: on 284.20, we have an expression for $R^{\text{th}} \text{ year}$.
 we want to see how $\left(\frac{m\ell}{n}\right)^m$ compares w. it. with $\left(\frac{e^{\frac{m\ell}{n}} - 1}{\frac{m\ell}{n}}\right)^R \approx \frac{R}{\frac{m\ell}{n}}$

03: $0.2 m \in R - \varphi n$ | $z = 1 - \frac{R - \varphi n}{\varphi n}$ | $e^{\frac{m\ell}{n}} = \left(\frac{e^{\frac{m\ell}{n}} - 1}{\frac{m\ell}{n}}\right)^m$ $0 < \frac{m}{n} < \frac{R}{n} - \varphi$
 $e^{-\frac{m\ell}{n}}$ | $\frac{m}{n} < \frac{R}{n} - \varphi$; which tends to be small as φ is small; hvr $\ell \left(\frac{R}{n} - \varphi\right)$?
 How does $\ell \left(\frac{R}{n} - \varphi\right)$ compare w. $\left(\frac{e^{\frac{m\ell}{n}} - 1}{\frac{m\ell}{n}}\right)^m$? See (284.20R) for sizes of z .

If ℓ is very long, maybe $\ell \left(\frac{R}{n} - \varphi\right)$ could be significant i.e. ~ 1 .

But look at 282.36 $\left(\frac{m}{n}\right)^m \left(1 - \frac{m}{n}\right)^{n-m\ell}$: where do these factors come from?
 $\left(\frac{m}{n}\right)^m$ is hvr of δ insertions; $\left(1 - \frac{m}{n}\right)^{n-m\ell}$ is the \downarrow of pc's of ~~former~~ former symbols of corpus.

Further factor; more exactly: m insertion is a total corpus long ℓ of $m - m\ell + m = n - m\ell(R-1)$

12 $\Rightarrow \left(\frac{m}{n - m\ell(R-1)}\right)^m$ for $\left(1 - \frac{m}{n}\right)^{n-m\ell}$, \downarrow correction is more complex: T. necessary was φ
 have a new symbol in pc = $\frac{m}{n - m\ell(R-1)}$ - substitution of $n - m\ell$ old symbols.

So we assume all old symbols have their pc's \downarrow by factor $\left(1 - \frac{m}{n - m\ell(R-1)}\right)$

So f. second factor is $\left(1 - \frac{m}{n - m\ell(R-1)}\right)^{n-m\ell}$

Note that ℓ expresses of ℓR \downarrow ℓR more symbols. T. some of exponents is $n - m\ell + m$ - as desired.

T. product of z is a common function of $m, n - m\ell$: $\left(\frac{x}{x+y}\right)^x \cdot \left(\frac{y}{x+y}\right)^y$

~~...~~ $\frac{x!y!}{(x+y)!} \approx \frac{x^x y^y}{(x+y)^x (x+y)^y} \cdot \sqrt{\frac{2\pi x y}{x+y}} = \sqrt{\frac{2\pi}{\frac{1}{x} + \frac{1}{y}}}$

One error in .12, .14 is that when δ is defined, not all old symbols have their pc's decreased by same factor. The symbols in δ have their pc's \downarrow more than others. Symbols not in δ may have their pc's \uparrow , since corpus is shorter yet their cases counts are larger. The "factorial" treatment of ≥ 14 respects this fact better.

Looking at 99 TM $\sqrt[122.02]{z(14)}$ ratio of \uparrow of pc due to δ in:

$\frac{\beta+1}{\beta+2\ell+1} \frac{\beta}{n} \sqrt{\ell m} \sqrt{2\pi} \left(\frac{\frac{\beta}{\gamma} e^{\frac{\beta}{\gamma}} - 1}{\frac{\beta}{\gamma}}\right)^R$

R is cysa count of δ . (δ has no repeated symbols, can't over-represent β) $\frac{\beta}{\gamma} \approx \frac{R}{n}$, $\gamma \approx \varphi$

for $\frac{\beta}{\gamma}$ close to 1 $\frac{\beta}{\gamma} \approx 1 + \frac{(\frac{\beta}{\gamma} - 1)^2}{2}$ (for small $\frac{\beta}{\gamma} - 1$) $\frac{\beta}{\gamma} \rightarrow "z"$ of 284.20

or ibid 122.10 $\frac{\sqrt{2\pi} \beta+1}{\beta+2\ell+1} \gamma \in \frac{z^R}{\sqrt{R}}$ If you assume $\frac{\beta}{\gamma} \geq 1$, hvr.

32 4 TM 284.21 $\frac{1}{\sqrt{n}}$ Equiv of $\frac{1}{\sqrt{n}}$ $z^R \sqrt{\frac{\beta}{\gamma} - 1}$ $\sqrt{\frac{\beta}{\gamma} - 1}$

The big difference betw. .31 & .32 is $\frac{1}{\sqrt{n}}$ factor: its quite large $\frac{\beta}{\gamma}$, but presumably $\left(\frac{z}{\gamma}\right)^R$ is dominant. Compare behavior of $\frac{1}{\sqrt{n}}$ approaches w. $\frac{1}{\sqrt{n}} \in z^R$ as $n \uparrow$.
 .31 is $\propto \frac{1}{\sqrt{n}}$, .32 is not $\propto \frac{1}{\sqrt{n}}$

Since m of interest is usually by m , we can use $\frac{1}{\sqrt{2\pi m}} \rightarrow \frac{1}{\sqrt{2\pi A}} = \frac{1}{\sqrt{2\pi(R-qn)}}$ as a correction factor.

We need to find largest term

$$\frac{R!}{(qn)^R e^{-qn}} = \frac{R!}{(qn)! (qn)^{R-qn}} \text{ mult by } \sqrt{2\pi} \text{ and } \frac{1}{\sqrt{2\pi(R-qn)}}$$

$$Z = \frac{R!}{2 (qn)! (qn)^{R-qn-\frac{1}{2}} \sqrt{\pi(R-qn)}} \text{ which is } \approx \text{ sum of } 282.36$$

Next, I want to search look at e in terms of $\left(\frac{R}{qn}\right)^R = e$

$$Z = \frac{R^R}{e^R \sqrt{2\pi R} (qn)^R \frac{e^{-qn}}{e^{qn}} \sqrt{\pi(R-qn)}} \dots$$

$$= \frac{R^R e^{qn}}{e^R (qn)^R (qn)^R \sqrt{2\pi R} \sqrt{\pi(R-qn)}} \dots$$

$$= \left(\frac{R}{qn \cdot e}\right)^R \cdot e^{qn} \cdot \sqrt{\frac{R}{\pi(R-qn)}} \dots$$

Expected e

First write Brook renews Poisson analysis of convolution $Z(n)$ using factorials.

$$\frac{R}{qn} = Z \left(\frac{e^Z - 1}{Z} \right)^R \cdot \sqrt{\frac{R}{R-qn}}$$

$$\sqrt{\frac{1}{\pi}} \left(\frac{e^Z - 1}{Z} \right)^R \cdot \sqrt{\frac{R}{Z-1}}$$

This looks like some debt as before! i.e. $\frac{1}{e}$ penalty!

The $\sqrt{\frac{Z}{Z-1}}$ factor is peculiar! For Z close to 1, it becomes very large!

It can be used as (∞) $\frac{\sqrt{2\pi R}}{\sqrt{2\pi A}}$

The $\frac{1}{\sqrt{2\pi m}}$ factor is from 283.01 $\frac{m^m}{m!} \rightarrow \frac{e^m}{\sqrt{2\pi m}}$ factor

$\frac{R}{qn} e^{\frac{qn}{R}}$	$\frac{e^{qn}}{e^{qn}}$
2	$\approx e^{\frac{1}{2}} - 1$
1	1
1.1	1.0044
1.5	1.074
2	1.21
3	1.54
5	2.246
7	2.97
20	7.73
∞	$\frac{e}{e}$

Q: $m \approx 283.37$ $\int \text{ say } L(1+\frac{1}{2}), L(1+\frac{1}{2}), \dots$

Maybe should be $L(1+\frac{1}{2})(1+\frac{1}{2}) \dots$ No

Check formulas for 283.22 $Z \approx 31$ simple recursive formula for $\frac{R!}{(R-n)! e^{-n}}$:

$f(n) = 1$ $T(n) = \frac{R-n}{n} \cdot T(n-1)$

$P = \varphi$

$n! = \dots$ [later, include "2" as mpt. variable]

$B = (R-nf)$

$X = 1/\text{sqr}(2p)$ - type α .

$S = 0$

for $m = \dots$ to B :

$X = X \cdot (R-n) \cdot n! : \dots$

$S = S + X \cdot \text{sqr}(m)$

Print S.

9/29/04

4TM

for $m \gg 1$, lots of cancellations in 282.36 terms! If $m \ll R$ which is common

$$\frac{m^m}{m!} \approx \frac{e^m}{\sqrt{2\pi m}} \left(\frac{R!}{(R-m)!} \right)^{1/n} \approx \left(\frac{R^m}{n^m} \right) \left(1 - \frac{m}{n} \right)^{n-m} \approx e^{-m + \frac{m^2}{n}}$$

(30) gets rid of this

So we get $\frac{1}{\sqrt{2\pi m}} \left(\frac{R}{n} \right)^m e^{-m + \frac{m^2}{n}}$
 This factor is the factor: as m + we multiply preceding factor by $\frac{R-1}{n}$.

until $R-1 = n$: i.e. 1. factor is 1. This $\frac{1}{\sqrt{2\pi m}}$ factor

Re: $e^{-m + \frac{m^2}{n}}$ $\frac{m}{n}$ has a max value of $\frac{R}{n} - 1$: so $e^{\frac{m}{n} \cdot mR}$

Why of all factors, why is most unreasonable! for very large m , $\frac{m}{n}$ will have a certain max size, but m can be very large, so $\frac{m}{n} \cdot mR$ can be very large!

This trouble is due to $\left(1 - \frac{m}{n} \right)^{n-m}$ factor (the $\left(1 - \frac{m}{n} \right)^n$ part seems ok, it's only e^m factor)

but $\left(1 - \frac{m}{n} \right)^{n-m}$ decreases. Take it to see could occur:

$$\ln \left(1 - \frac{m}{n} \right) \approx -\frac{m}{n} + \frac{m^2}{2n^2} \quad \left(1 - \frac{m}{n} \right)^{n-m} = e^{\left(-\frac{m}{n} + \frac{m^2}{2n^2} \right) (n-m)}$$

$$= \exp \left(\frac{m^2}{n} - \frac{m^3}{2n^2} \right) = \exp \left(\frac{m}{n} \cdot mR - \frac{m}{n} \cdot \frac{m^2}{2} \right) = \exp \left(\frac{m}{n} - \frac{1}{2} \left(\frac{m}{n} \right)^2 \right) mR$$

so second (probably by) order, doesn't help:

m 's values available values $\propto R$: if m very large, m can be very large.

Perhaps in 282.36, $\left(\frac{m}{n-mR} \right)^m \left(1 - \frac{m}{n-mR} \right)^{n-mR}$ would be more correct! It would

correctly give a near result! - No! $\left(\frac{m}{n-mR} \right)^m$ still gives $e^{\frac{m^2}{n}}$ factor. $\rightarrow 285.00$

$$\frac{1}{\sqrt{2\pi m}} \frac{R!}{(R-m)! (n\varphi)^m} \quad \text{This second factor is } \frac{R}{n\varphi} \cdot \frac{R-1}{n\varphi} \cdot \frac{R-2}{n\varphi} \dots \frac{n\varphi}{n\varphi}$$

More exactly $\frac{R}{n\varphi} \frac{R-1}{(n-1)\varphi} \frac{R-2}{(n-2)\varphi} \dots$ **NO!**

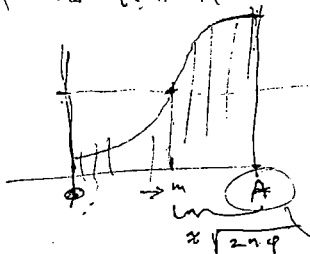
So this begins to look reasonable again!

its $\frac{R}{n\varphi} \cdot \frac{R-1}{(n-1)\varphi} \cdot \frac{R-2}{(n-2)\varphi} \dots \frac{R}{(n-2\ell)\varphi} \cdot \frac{R-1}{(n-2\ell)\varphi} \cdot \frac{R-2}{(n-2\ell)\varphi} \left(\frac{R}{(n-2\ell)\varphi} \cdot \frac{R-1}{(n-2\ell)\varphi} \cdot \frac{R-2}{(n-2\ell)\varphi} \cdot \frac{R-3}{(n-2\ell)\varphi} \right)$

I'm not sure it's ok, but it's beginning to look like a good bet!

$\frac{A}{n=0} \approx \frac{(n\varphi)^{n-m}}{(R-m)!}$ is a very well known (i.e. no $\frac{1}{\sqrt{m}}$ factor),

The terms in $n=31$ increase exponentially at first, but then flatten out as the limit is reached. So plot of "half value points" show how large m is.



Working back words: largest terms L ; next is $L \frac{n\varphi-1}{n\varphi} = L \left(1 - \frac{1}{n\varphi} \right)$

Next is $L \left(1 - \frac{1}{2n\varphi} \right) \cdot \left(1 - \frac{1}{2n\varphi} \right) \dots \left(1 - \frac{R-\alpha}{\alpha} \right)$

$\left(1 + \frac{1}{2\alpha} \right) \approx \exp \frac{1}{2\alpha} \rightarrow \exp \frac{r(r+1)}{2\alpha} \approx \dots \approx r(r+1) = 2\alpha = 2n\varphi$

$A = R - \varphi n = m_{max}$

0/29 04
4TH

2045B

00

Try $m=2$ in 280.22: $\epsilon = \frac{2}{n(1-\phi R)}$

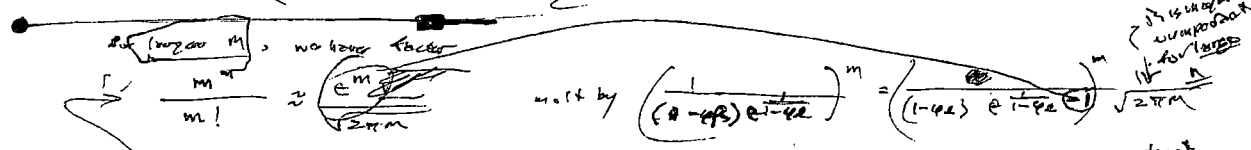
in 280.22 we had $\epsilon = \frac{1}{n(1-\phi R)}$ for $m=1$

if $m=2$ $\epsilon_{max} = \frac{2}{n(1-\phi R)}$ 280.22 gives $\left(\frac{2}{n(1-\phi R)}\right)^2 (1-\frac{2}{n(1-\phi R)})^{n-2} \cdot \frac{R(R-1)}{2} \cdot \frac{1}{\phi^2}$

$= \frac{2 \cdot 2}{2} \cdot \frac{1}{(1-\phi R)^2} \cdot e^{-2 \cdot \frac{1}{1-\phi R}} \left(\frac{R}{n} \frac{R-1}{n} \cdot \frac{1}{\phi^2}\right)^{\pm 1} \approx 2 \cdot \frac{1}{(1-\phi R)^2} \cdot e^{-2 \cdot \frac{1}{1-\phi R}}$

$\approx 2 \cdot \left(\frac{1}{(1-\phi R) e^{1-\phi R}}\right)^2$

2.82	1
3.69	$\frac{1}{2}$
6.70	$\frac{1}{3}$



This is a little better than the sum product; but now, the possibility of poor use of ϵ is dependent on ϕR which seems strange!

$\frac{m^m}{m!} \approx e^{m/\sqrt{2\pi m}}$

T. source of ϕR : 279.34: $\phi(1-\epsilon)^2 + \epsilon = \frac{R}{n}$

To second order: 279.34 $\phi(1-\epsilon)^2 + \frac{2\epsilon(1-\epsilon)}{2} \epsilon^2 = \frac{R}{n}$

$1 - 2\epsilon + \frac{\epsilon}{\phi} + \frac{2\epsilon(1-\epsilon)}{2} \epsilon^2 = \frac{R}{n\phi}$

$\frac{2\epsilon(1-\epsilon)}{2} \epsilon^2 + (1-2\epsilon)\epsilon + 1 - \frac{R}{n\phi} = 0$

$2-1 \pm \sqrt{(2-1)^2 - 2\epsilon(2-1) \cdot \left(1 - \frac{R}{n\phi}\right)}$

$(2-1)^2 - 2\epsilon(2-1) \cdot \left(1 - \frac{R}{n\phi}\right)$

$2^2 - 2\epsilon + 1 - 2\epsilon^2 + 2\epsilon$

$= -2\epsilon^2 + 1 = 2\epsilon^2 - 1$

$x^2 - 2x + 1 = 0$

$+ 2 \pm 4 - 4$

$= 1$

For small $1-\phi R$ (0) looks like its not far off from reasonable. Perhaps using full 280.18

miter to ok.

$2-1 \pm \sqrt{(2-1)^2 + 2\epsilon(2-1) \cdot \delta}$

$\frac{2\epsilon(1-\epsilon)}{2} \epsilon^2 + \epsilon \left(\frac{1}{\phi} - R\right) + \left(1 - \frac{R}{n\phi}\right)$

$\sqrt{\left(\frac{R-1}{\phi}\right)^2 + 2\epsilon(2-1)\delta} \approx \left(\frac{R-1}{\phi}\right) \sqrt{1 + \frac{2\epsilon(2-1)\delta}{\left(\frac{R-1}{\phi}\right)^2}}$

$2-1 \pm \left(\frac{R-1}{\phi} + \frac{2\epsilon(2-1)\delta}{2-\phi}\right) / (2-\phi)$

for $\delta = \frac{R}{n\phi} - 1 > 0$

$\frac{2\epsilon(2-1)\delta}{2-\phi} = 1 + \frac{2\epsilon(2-1)\delta}{(2-\phi)^2}$

There is another small factor: $(1-\epsilon)^{-m} \approx e^{\epsilon m}$ in (10)

Weight $\left(\frac{1}{(1-\phi R)} e^{1-\phi R} - 1 - \epsilon m R\right) = \frac{e^{1+\frac{m^2}{2n}}}{(1-\phi R) e^{1-\phi R}}$

$\frac{m^2}{n} = \epsilon m R$ which can be very large if m is large!

Back to 280.22: if $\epsilon m = \frac{m}{n}$ $\epsilon m = \frac{m}{n-m}$

$\sum_{m=0}^{\infty} \left(\frac{m}{n}\right)^m (1-\frac{m}{n})^{n-m} \cdot \phi^{-m} \frac{R!}{m!(R-m)!}$

then $\epsilon = 0$ form is 1 for $m < 1$ form

$\frac{1}{n} \left(1-\frac{1}{n}\right)^{n-1} \phi^{-1} \cdot R = \frac{R}{n} \frac{e^{\frac{1}{n}}}{e^{\frac{1}{n}}}$

cancel $\frac{e^{\frac{1}{n}}}{e^{\frac{1}{n}}}$!!

inversion of bar $\epsilon > 0$.

36

See 284.06 for approx. summation.

4TH

IAP | Schedules office | Post ::::: IAP url
253 4788 | http://mit.edu/iap/

organize : non-credit
post and review. 1/1/10.

0 : Actually, a messy way to check formula -1.22 or -1.25, would be numerical of M Back or Maple or Mathematica. I can start in first order Approx. of $E_m(-1.22R)$ (279.34 R) for second order. The previous was R, ϕ , n, λ . To start, input pc as a function of $\frac{R}{n} \approx \phi$ (say n) proceed

For factorial, use $\ln n!$ precision formula.

In Basic; I have forgotten how to do subroutines (for X! say).
formula -1.22 doesn't look bad!

I want to compare it in pc before coding: $500 - 1.07$; so -1.22 gives a mpc due to ϕ definition. It should be? if $\frac{R}{n} = \phi$; but using $\phi \approx \phi$ in Rect case (?).

0.08 try to be $\frac{R}{n}$ very slightly $> \phi$. Anyway, this gives $\epsilon \approx 0$. No!

0.09 we sum m from R to R . one term $\phi^{-R} = \left(\frac{R}{n}\right)^{-R} \cdot \frac{R!}{R!} \approx \left(\frac{R}{n}\right)^{-R} = \left(\frac{n}{R}\right)^R$

0 Some have only ϕ^{-R} factor: -1.07 is supposed to explain it.

1: (-1.07)! The ϕ^{-m} factor does not occur; was ϕ^{-R} still

No! - (0.08) recurring. $\epsilon \approx \phi$. I have only gotten ϵ smaller values confused!

I really have to go over formula -1.22! - one -1.07!

Re: -1.07 ϕ^{-} At certain points ϵ may be confusing in such

Poss. confusion: I would want to go from ϕ to $R - \phi n$ (I was confused before R and n)

Since $\epsilon \approx \frac{R}{n} - \phi$, $\epsilon \approx 0$ if $\frac{R}{n}$ is only marginally $> \phi$.

Exactly what is m ? I think m is the no of times λ appears in the code for the corpus. ϵ can be $> \phi$, but ϵ is just ϕ code for corpus. So $m \approx \epsilon n$ at most $\approx \phi$ at minimum. If $m = 0$ then pc of code = $(1 - \epsilon)^n$

would it be adding extra correction for λ units on M (-1.22) to be $0 \approx R - n\phi$?

View of reasoning - 1.22 begins to look reasonable! $e^m (1 - e)^{n-m}$ for $m=0$

1: ϵ , $\phi \approx \frac{R}{n}$ is $1^m \approx 1$, but if would be more reasonable to have $m=1$ & n large.

O.K., say $m=1$; $m > 1$: $\frac{R}{n} = \phi n + 1$ $\frac{R}{n} = \phi + \frac{1}{n}$ $\frac{R}{n} - \phi = \frac{1}{n}$;
So for -1.22 $\left(\frac{1}{n(1-\phi)}\right)^1 \left(1 - \frac{1}{n(1-\phi)}\right)^{n-1} = \frac{1}{n(1-\phi)}$ $e^{-\frac{n-1}{n(1-\phi)}} \approx \frac{\exp(-\frac{1}{1-\phi})}{n(1-\phi)}$

$1 - \phi > 0$ because $\frac{1}{\phi} > 1 > 0$ because $\frac{1}{\phi} > 1$ because $\frac{1}{\phi}$ is expected bits below occurrences of λ . Reason for $\frac{1}{\phi} > 1$

$\frac{1}{\phi} - 1 = \frac{1}{n(1-\phi)}$
 $\frac{1}{\phi} = \frac{1}{n(1-\phi)}$
 $\frac{1}{\phi} - 1 = \frac{1}{n(1-\phi)}$
2 approximations

In -1.22 we have $\frac{R!}{m!(R-m)!} = \frac{R!}{1!(R-1)!} = R$ (if we sum ϵ from $0 \leq 1$)

$\frac{R!}{0!(R-0)!} = 1 \left(\left(\frac{1}{n(1-\phi)}\right)^0 \cdot \left(1 - \frac{1}{n(1-\phi)}\right)^n \approx e^{-\frac{1}{1-\phi}}$

So it looks like -1.22 gives $(R-1) e^{-\frac{1}{1-\phi}}$
for $m=1$; 1 term is $\frac{R}{n(1-\phi)\phi} e^{-\frac{1}{1-\phi}} \approx \frac{R}{\phi n(1-\phi)} e^{-\frac{1}{1-\phi}} =$

for $m=0$: $e^{-\frac{1}{1-\phi}}$ so this is worse than $e^{-1+\phi} \approx e^{-1} \cdot \phi \approx \frac{1}{e} (1 + \phi)$.

$\frac{1}{\phi} - \phi = \frac{1}{n}$
 $R - \phi n = 1 \approx n$
 $R - \phi n = 2$
 $\frac{R}{n} - \phi = \frac{2}{n}$

Can we get how 0 $n \left(\frac{R}{n} - \phi\right) = R - n\phi$?

we have R poss. positions \approx pc = $\frac{1}{\phi}$ so $\frac{R}{n(1-\phi)}$ we can't by $\frac{1}{\phi}$ to get $\approx \frac{1}{1-\phi}$

There is still 1 term $(1 - \epsilon)^{n-1} = \left(1 - \frac{1}{n}\right)^{n-1} \approx \frac{1}{e}$. So I still seem to have this $\frac{1}{1-\phi} \cdot \frac{1}{e}$ for coding (1 occurrence of λ).

0: : Actually, for PC of codes using ϵ & defn. are much higher due to ϵ coding: Of all occurences of ϵ in δ ; any no. can be coded by ϵ defn, & this can be done in many ways, so we get aH Diverse codes - many of different PC.

1: Here, the definition of δ has PC ϵ & will δ PC of letters comes by δ $(1-\epsilon)^n$ (actually, ϵ exponent is n no. of symbols other than ϵ , no. for ϵ .) in δ - no. codes of ϵ comes. $\approx \epsilon^{-n}$

2: Say ϵ is at the R occurences of δ , we code ϵ of them as δ defn. ϵ (R-m) of them as alphabet symbols.

3: T. PC of coding is $\epsilon^m \cdot (1-\epsilon)^{n-m}$. ϵ is different for each m . $\epsilon_m = \frac{m}{n} - \phi$ (where $m \leq \phi n$ then $\epsilon_m \geq 0$.)

4: for a given m there are $\frac{R!}{m!(R-m)!}$ ways to do this. ϵ is different for each m . $\epsilon_m = \frac{m}{n} - \phi$ (which goes to zero as m from ϕn to R - $\phi \cdot 20$ - ϕ is low freq.)

5: 2.79.14 : EN Re: -1.10 To make this less. \approx A.H., apply it to all occurrences of δ : i.e.

6: If δ has been defined, whenever δ is created by a prefix: i.e. destroyed, the way way δ can be created by its choice as a symbol (\approx "to know").

7: An objection to .10: say δ is ϵ symbols (say: I'm looking at a corpus which other last-2 symbols of δ have just occurred. There are 2 ways that ϵ symbol can occur: 1) By alphabet 2) by δ selection so ϵ is of that symbol is very high; However, if we read ϵ symbols of δ in a corpus then ϵ is for ϵ .

8: next symbols ϵ - even this complex δ , it can only occur by δ selection - so there is this discrimination

9: at just the point at which δ is very relevant!

10: Perhaps drop ϵ - 1.10 is (280, 40) for white: The "unusually low freq. of an asen" Does it occur?

11: Very often. Until I can think of a possl cause of it in RW, it seems not so heavy!

12: we know redundancy of factor $\frac{R!}{m!(R-m)!}$

13: So we want $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m} \frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

14: $\frac{R!}{m!(R-m)!}$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

15: $\frac{R!}{m!(R-m)!}$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

16: $\frac{R!}{m!(R-m)!}$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

17: $\frac{R!}{m!(R-m)!}$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

18: To each phrase of Gauss. form: find phrase's secret derivative.

19: $(1+b-2am)^{n-m}$ may be expressed as $\sum_{k=0}^{n-m} \binom{n-m}{k} (1+b-2am)^k$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

20: $(2am+b)^m = 2^m m^m - b 2^{m-1} m^{m-1} + \dots$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

21: Looks like ϵ is $\frac{\phi m}{2!}$ which is $\frac{\phi m}{2}$ $\epsilon^m (1-\epsilon)^{n-m} \phi^{-m}$ $\frac{R!}{m!(R-m)!}$ $\epsilon_m = \frac{m}{n} - \phi$ $\approx 2\phi m$ ϕ is low freq.

22: with ϕ of $\frac{\phi m}{2!}$. It's also quite so simple! I have m^{m-2} and $\frac{m!}{(m-2)!}$

0: 2 2840 : So: when one defines a new grammar; It may be that the resultant grammar never parses some PCs as previous grammar: even if the new defined had same previous frequency as was expected from its symbols!

0 4: ① 275: 15-16 still seems true ② Definitives VERY implementation Sci discoveries
③ -1.36-37 ^{of} Penalty factor of ϵ for each occurrence \rightarrow of interest!

BALDWIN effect
Baldwin effect,
Lacovick,
Feed back from generation of children Fitzpatrick

Perhaps one can define Contexts (for produ) w. impunity.

0 6: ① Some other ideas: γ can occur 2 ways: by defined symbols, or by creation by a reverse in itself.
Try code in which $\varphi =$ twice expected from other symbols. Frequencies of all symbols will change:
See how coding cost of corpus varies as I change PC of δ & ϵ of other symbols.

② for PC of γ in lower than expected & δ have operator Not look for δ & destroys it w. PC ϵ . Since this can be done in many ways, to produce observed corpus = we have || codes δ & of PC! Since δ is less likely, to resist to fossil outputs & to grammar have to \uparrow something, to preserve "Normality".

504
26.10
230.10

③ way this is done: any symbol in the corpus can be created (parsed) 2 ways; δ solves leads to PC of δ symbols.
④ By direct creation δ , & ϵ is δ (PC δ) Next create β .

06 - If See has reasonable solutions to ϵ . 2 14 problem of -1.36-37

If we have no constraint equals other than their frequencies, it is clear that ϵ code of δ -1.26-35 gives us 2 ways to create δ & ϵ , total PC of δ is much greater than (about twice) what we expected!

On the other hand, δ PC of ϵ code ϵ occurs ϵ ϵ times as often as δ .

The result of this: If $\frac{R}{N} = \varphi$, then all instances of δ are accounted for as generated ^{current} alphabet. If $\frac{R}{N}$ is only slightly $> \varphi$, but corpus is long ($n \gg \frac{R}{N}$) then it will pay to do δ , but only a small fraction will be generated by δ code. δ & we will know which

if $\frac{R}{N} \gg \varphi$ say 10 times as large than most δ 's will be generated by δ code, but we will know which; furthermore, there will be many || codes because there will be many permissible combinations of many ways to create δ ... by definition, δ by original alphabet symbols.

If $\frac{R}{N} = \varphi + \epsilon$ then defining δ so it has PC of exactly μ , will not work because

Proz decln. \downarrow PC of other symbols so PC of creation of δ by alphabet φ .

If we define δ to have PC = ϵ , then $\varphi \leftarrow \varphi(1-\epsilon)^2$ (since all other alphabet symbols keep their PC \downarrow by factor $(1-\epsilon)$ as δ is 2 symbols long).

So we want $\epsilon \Rightarrow \varphi(1-\epsilon)^2 + \epsilon = \frac{R}{N} \dots \approx \varphi \epsilon^{-2} + \epsilon = \frac{R}{N}$

for small ϵ ; $\approx \varphi(1-2\epsilon) + \epsilon = \frac{R}{N}$ $\varphi = \epsilon \varphi + \epsilon = \frac{R}{N}$ | $\varphi + (1-\varphi)\epsilon = \frac{R}{N}$

$(1-\varphi)\epsilon = \left(\frac{R}{N} - \varphi\right)$ (this is ϵ ans. we have to \uparrow PC of δ to get to $\frac{R}{N}$)

$\epsilon = \frac{\frac{R}{N} - \varphi}{1-\varphi}$ usually $\varphi \ll 1$, so ϵ will $\approx \frac{R}{N} - \varphi$, which is reasonable

More exactly, $\varphi(1-\epsilon)^2 + \epsilon = \frac{R}{N}$ to second order $(1-\epsilon)^2 = 1 - 2\epsilon + \frac{2\epsilon^2}{2} \ll 2 \dots$
So $\varphi(1 - 2\epsilon + \frac{2\epsilon^2}{2}) + \epsilon = \frac{R}{N}$

This can easily be repeatedly solved by successive approach, $\epsilon_i = \frac{R}{N} - \varphi(1-\epsilon_{i-1})^2$ or if it doesn't converge do it

by $(1-\epsilon) = \left(\frac{R}{N} - \epsilon + 1\right) / \varphi$

00 : T. result of -1.38-.39 is simply obtained & very reasonable. T. previous stuff w. re factor of e & re have various diffys in them & may be wrong.

-02 Also this result is: idea of 275.15-.18 should make it clearer to how to use unusually high or low freqs of ngrams to d code length.

T. way it works: If an ngram has "hydraulic lowcut freq." & to say sm & cost will pay for its definition. [I suspect low freq. ngrams will not often pay for their defs, unless fr. corpus is very long!]. The way we code using the definitions:

< Say for ngrams w. high freqs: > We make definition of the high ngram $\equiv \delta$. ~~It has a "normal" freq. for its component~~ (by concatenating its components) coding using δ would not change fr. b cost of corpus. The other ($= \delta$) symbols will have their freqs changed, but in such a way to preserve Pair ratios; R. all $= \delta$ symbols will have their freqs \downarrow from before, however we now have a new symbol, δ & it has a certain freq.:

~~... T. foregoing analysis holds whether δ has unexpectedly high or low frequency.~~

My guess is that we will not often make defs of ngrams as very low freq. because

they will usually have only a small ~~freq~~ & mb cost - not an outlay for definition. ^{directly/Leffcus} ^(Dcr's) ^(Laps/role)

To foreg. ~~...~~ That's worse for coding using probabilities: If we code using δ we should get about same results, ~~but there may be "e" factors~~, but they should cancel out

when we compare diffnt methods of coding (i.e. with v.s. w.o. v.s. a defn).

I should calculate cost of corpus in which I define a δ that has same pc as product of pc's of its components. To gain ^{this cost} "cost of definition" just add a symbol, δ for prevars.

Do analysis 2 ways: 1) simply using probs relative to same ^{for that symbol} ~~each time~~ that symbol is used.

2) using pair coding (≥ 14)

(26) $(1-p)^{\frac{1}{p}-1} = (1-p)^{\frac{1-p}{p}}$

I think 1. foreg. stuff (277.20 - 278.22) (278.00 - .22 in particular) may be a very simple break thru in the ≥ 14 problem! It seems to solve a Paradox "factor of e" Paradox 2) How to handle ngrams of unusually low frequency.

For (2), use idea of -1.38-.39 of coding a "typical" section of T. code, containing one δ . ^{many} ^{probab} of δ ~~...~~ defnt symbols ^{No! This is not a paradox, but a cost factor!} ^{See 272.16 and 278.22 - 278.27 for how $\frac{1}{e}$ arises}

Consider case in which ϕ is the normal pc assigned to δ & consider δ symbols. in the $a+1$ symbol corpus, one symbol has pc = ϕ (?). The other symbols have their pc's \downarrow by factor $(1-\phi)$.

So, before coding a new symbol, the first a symbols had pc's A , the last 2 had pc's ϕ . Total pc of corpus = $A \cdot \phi$. With new codes, T. first a symbols have pc's $A(1-\phi)^a$ and the last symbol has pc's ϕ . So total pc = $A(1-\phi)^a \cdot \phi$. ANALOGY $(1-\phi)^2 = 1 - \phi^2 = (1-\phi)^2$

This looks like old ≥ 14 paradox: we save parity of about $\frac{1}{e}$ factor of e for each occurrence of δ , when we define δ . AGAIN constraint we

There may be more constraint on the δ symbols, since linear that they can't even define δ . The resultant Bernoulli grammar w. δ does not produce

1. some corpus that originally had δ & unmodified, δ occurs w. greater freq. than before to defn of δ .

20 : $m = (0.01) \cdot 2 \cdot 0.8$ then $(1-q)^m = (1-q)^{n-LR}$ has to be under the \leq sign.

$$\sum_{R=0}^{\infty} R \frac{(qR)^R}{R!} (1-q)^{n-LR} \quad [m \equiv n-LR]$$

$$\approx \sum_{R=0}^{\infty} \frac{(qR)^{R-1}}{(R-1)!} e^{-mq} = \frac{(qR)^{R-1}}{(R-1)!} e^{-mq}$$

constant $R!$ so this is an approx. - division.

$$= qm \cdot e^{-n} \sum_{R=1}^{\infty} \dots$$

oops! $\sum_{R=0}^{\infty} R \dots \neq e^{qm}$, because m is a function of R .

$$(qR)^R \approx \left(nq \frac{(n-LR)}{n} \right)^R \approx (nq)^R e^{-\frac{LR^2}{n}}$$

$R_{max} = \frac{n}{2}$
 $e^{-\frac{LR^2}{n}} = e^{-\frac{L}{2}}$

$$e^{-\frac{LR^2}{n}} = e^{-1} \text{ when } \frac{L}{n} R^2 = 1 : R^2 = \frac{n}{L} : R = \sqrt{\frac{n}{L}}$$

for $n \gg L$ (i.e. a long corpus), R will be large before this occurs.

$\frac{(qR)^R}{R!}$ will perhaps be quite small: say $qR \approx R$ $\frac{R^R}{R!} \approx \frac{R^R}{R^R} \sqrt{2\pi R}$

so not so small! $\approx e^R \cdot \sqrt{2\pi R}$ - This is mult by e^{-mq} to get e^{R-LR} ...

which is not so huge compared to $R!$ when n is large \approx sum of series.

T. func. stuff is very "fuzzy" : I should check this stuff numerically, I don't, say a Monte Carlo

run to see how frequently an new approx.

Look at it this way: we are going into sequences incrementally ϵ . we have $p = q$ of emitting δ .

If we do emit δ , then L consecutive δ 's, but then all we can emit is δ again.

Let's look at d.f. of distances betw. δ 's! Study of ϕ is ϕ prob of (1) $(\phi)^L$

prob of n is $(1-q)^n q$. Expected distance is ϕ (1) $(\phi + 1 \cdot (1-\phi) + 2 \cdot (1-\phi)^2$

say $1-q \equiv \alpha$ ($\Rightarrow \phi \equiv 1-\alpha^2$) $= \frac{q\alpha}{1-\alpha^2} = \frac{q\alpha}{2(1-\alpha^2)}$ $\equiv \frac{\alpha}{1-\alpha^2}$

$$\sum_{n=0}^{\infty} n \alpha^n = \frac{1}{1-\alpha^2} \quad \frac{d}{d\alpha} \frac{1}{1-\alpha^2} = \frac{2\alpha}{(1-\alpha^2)^2} = \frac{2\alpha}{(1-\alpha^2)^2}$$

$$\frac{d}{d\alpha} \frac{1}{1-\alpha^2} = \frac{2\alpha}{(1-\alpha^2)^2} = \frac{2\alpha}{(1-\alpha^2)^2}$$

expected distance = $\phi \frac{2}{(1-\alpha^2)^2} = \frac{1-q}{q^2} = \frac{1-\phi}{\phi^2} = \frac{1-\phi}{\phi} - 1$ \leftarrow I made some mistake with this one.

Expected distance = $\phi \cdot \phi + \phi \cdot 1 + \phi \cdot 2 + \phi \cdot 3 = \phi(\alpha + 1 + 2\alpha + 3\alpha^2 + \dots)$

$$1 + \alpha + 2\alpha^2 + 3\alpha^3 = \frac{1}{1-\alpha} \quad \frac{d}{d\alpha} \frac{1}{1-\alpha} = \frac{1}{(1-\alpha)^2} = \frac{1}{(1-\alpha)^2}$$

$$\alpha \cdot \frac{1}{1-\alpha} = \frac{\alpha}{(1-\alpha)^2} = \alpha + 2\alpha^2 + 3\alpha^3 + \dots$$

$$\phi(\alpha + 2\alpha^2 + 3\alpha^3 + \dots) = \frac{\alpha^2 \phi}{(1-\alpha)^2} = \frac{1-\phi}{\phi} - 1 = \text{expected no of symbols betw } \delta$$

so we have a "corpus" of $\frac{1}{\phi} - 1$ symbols followed by L symbols. - which has 1 α init.

so for n relative probability so in a corpus of length n , we expect δ to occur

$$\frac{n}{\frac{1}{\phi} - 1} \text{ times. It is } \frac{n \phi}{1 - \phi + \phi} = \frac{n \phi}{1 + (L-1)\phi}$$

if $(L-1)\phi \gg 1$ (which is usual); # δ 's $\approx n \phi$.

0:25:40:

$$x \frac{d}{dx} \sum_{R=0}^{\infty} \frac{x^R}{R!} = \sum_{R=0}^{\infty} \frac{R \cdot x^{R-1}}{R!} = \sum_{R=1}^{\infty} \frac{x^{R-1}}{(R-1)!} = x e^x$$

so $\sum_{R=0}^{\infty} \frac{R (e^x)^R}{R!} = \sum_{R=0}^{\infty} \frac{R e^{Rx}}{R!} = x e^{e^x}$ which is ~~more~~ ~~unreasonable~~ ~~to~~ e^{e^x}

01

My calculations of P of just R occurrences may be wrong!

A (perhaps) minor point. $\sum_{R=0}^{\infty} \frac{(e^x)^R}{R!} = e^{e^x}$. The largest term is $\frac{(e^x)^m}{m!} = e^{xm} \cdot \frac{m^m e^m}{m!} \approx \frac{m^m e^m}{\sqrt{2\pi m}}$

$(e^x)^m$ is usually quite small. mult by $\frac{m}{2\pi}$ still very small. So e. approxn of α is $\frac{m}{2\pi} e^m \cdot \frac{m^m e^m}{\sqrt{2\pi m}}$

How I really enter req of + duplication of cases whenever hour > 1 & at same point?

number $\frac{m^R}{R!} \cdot e^R (1-e)^{n-R}$ $(e^x)^m \approx e^{-pm}$ when n is small & multi factor for m. $\frac{m^m e^m}{m!}$ when n quite reasonable. $\frac{m^m e^m}{m!} \approx \frac{m^m e^m}{\sqrt{2\pi m}}$

09

T. d. work over prob of various R values as $\frac{(e^x)^R}{R!} \cdot \frac{e^{-n-R}}{e^R} = \frac{(e^x)^R}{R!} \cdot e^{-n-2R}$

$= e^{-n-2R} \left(\frac{e^x}{R} \right)^R \cdot \sqrt{R}$ $\left(\frac{e^x}{R} \right)^R \cdot \sqrt{R}$ say $\frac{e^x}{R} = z = 10$. $2^5 \cdot \sqrt{5} = 32 \cdot \sqrt{5}$

$z = \frac{e^x}{R} = 10, f = \sqrt{R} = \sqrt{10} : \Rightarrow R \approx 10, f \approx 3.16$ say $R = \frac{z^2}{2} = 5$

$R = 1, z = 10, 10^1 \cdot \sqrt{1} = 10$ so peak here. $R = 1 \approx R = 10$ ($R = 10$)

$R = 1; F = z$
 $R = 2; F = \left(\frac{z}{2}\right)^2 \cdot \sqrt{2} = \frac{z^2}{2} \cdot \sqrt{2}$
 $R = 3; F = \frac{z^3}{3} \cdot \sqrt{3}$

so peak is peak for R before z .
 100 Hz $\left(\frac{z}{R}\right)^R$ only peak $z = 10$ mod: $\left(\frac{z}{R}\right)^{\frac{R}{z}}$ $\frac{R}{z} = \alpha$

for fixed R $\frac{z^R}{R!}$

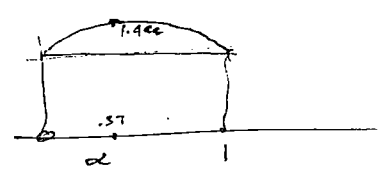
α^{-x}	
1.047	.1
1.25	.1
1.38	.2
1.44	.4
1.432	.45
1.414	.5
1.36	.6
1.29	.7
1.20	.8
1.10	.9

So peak $\alpha \approx .37 = \alpha^{-\alpha}$

Arb. peak. $\alpha \approx 1.444$ (of $\alpha = 1$)
 $\alpha \approx .57$

1.3	1.435
.35	1.44659
.37	1.44659
.4	1.442
.45	1.432

.38 m. add 38



\Rightarrow trying to z favor of f will increase R . 1.44 peak \approx (of $\alpha = 1$)
 $\Rightarrow z$ is as much ≈ 100 , \sqrt{R} is of computer - R peak will
 var of $\alpha = .37 \approx$ width of peak for z

$x^{-x} = -x \ln x$ $\frac{-x}{x} \ln x = -\ln x$ $-1 = \ln x$ $x = \frac{1}{e} = .36787...$

$e^{\frac{1}{e}} = 1.444667 =$ peak. get some div. to calculate.

$x^{-x} = e^{-x \ln x} (d) = (1 - \ln x) e^{-x \ln x}$ $\frac{d}{dx} = -\frac{1}{x} e^{-x \ln x} + (1 - \ln x)^2 e^{-x \ln x}$

so sec deriv = $\left((1 - \ln x)^2 - \frac{1}{x} \right) x^{-x}$ $x = \frac{1}{e}$ $\left((1+1)^2 - e \right) e^{-1} = 1.444667 \approx 1.851657049$
 = second deriv of peak of x^{-x}
 peak = $x = 1.444667$

from, note that probability is $\frac{z^R}{R!} = \frac{z^R}{R!} e^{-z} = \frac{z^R e^{-z}}{R!}$
 operator $\frac{z^R e^{-z}}{R!} = 1, \frac{z^R e^{-z}}{R!} = 1, \frac{z^R}{R!} = \frac{z^R e^{-z}}{R!}$

$\left(\frac{z}{R}\right)^R \rightarrow \left(\frac{R}{z}\right)^{-R}$ peak at $\frac{R}{z} = \frac{1}{e} = .36787$
 even at $R = 1$ ($\frac{1}{1} = 1$) we still have peak $\frac{1}{e}$ factor

0/25/04
4PM

The Z141 Paradox: 10

A "Final Sobolov" seems to be 277.20 - 90? Then 278.00 - 22
This also solves problem of defining n/m's w. unusually low frequency. - This search is not to let

Not sure how best to call
say 285.31 - 32!
say 280.10 - 19
say 282.10 - 13
for how $\frac{1}{E}$
log. timeliness
at 300.

0: 27.40 ; One characteristic of the s-branch of -1.34-90: While it is "relatively easy" to dig over
Grammar of this type from a corpus, they don't easily assign a pc to an arbitrary object.
for grammar models, R. object has to be parsed! (The R's are R's of Matrix parsing would work
o.k. for s-grammar. It takes n^3 operations to mult 2 $n \times n$ matrices. (only n^2 for d-matrices
for $n \leq 64$, say). These models can give Monte Carlo d.f.'s easily, however. These problems
is in making a prediction i.e. a continuation of a string

11.14.04

The n^3 is disastrously large! Most of elements of matrices of interest are zero - particularly
if + start of + computation. All non-zero mults in addition are necessary important into final result.
We may use d-matrices to keep track of where non-zero elements occur in s-matrices, so as to
speed up calcn. To cost of keeping the d-matrices will be \ll cost of multiplying s-matrices.
So this parsing of s-grammars may not be as bad as it initially sounds ($M \geq C$ since).

1: On the Z141 paradox, say φ is the pc of n symbols obtained via pc's of its symbols.
say δ occurs R times in corpus of n symbols. \rightarrow so $\frac{R}{n}$ is its observed frequency.
 φ is the pc that δ would occur at an arbitrary, 2, in the corpus. To estimate freq. of occurrence of δ in
corpora in which the only constraints were frequencies of subsymbols of δ : - Also some trouble in overlap.

If δ is of length l , there are only n^{l-1} opportunities for δ to occur (??).
At any rate, it deviates much from expected no. of occurrences of δ from what we should have

2: ϵ way to use this fact for prediction - i.e. a way to shorten codes. This argument is
3: true of R as too small or too large: I imagine it's usually easier to find discrepancy for
4: every large R, because "small R" can only be ≤ 0 ; but large R can be many times $>$ expected.

5: First consider case of δ being "unnumberable" (i.e. sequence "aaa" is over letters,
6: ϵ sequence "abc" is not.) [I think I did work this out in detail.]

7: A way to do it: To get pc of just R cases of δ in corpus of n symbols: Consider l n -R
8: symbols that are not part of δ 's: Any of these symbols can be followed by a δ [we consider these inserted δ 's to
9: be of zero length, so some of them can fit into δ 's & δ 's following any one of the n -R symbols] (29)

10: Anyway, for very large n , it would seem that δ would occur w. frequency $R = n\varphi$. a unit
11: R deviates from $n\varphi$ much (i.e. $\frac{R}{n}$ deviates from φ much) from what we ~~should~~ pretty have & regularly,
12: In my Z141 analysis, $\frac{R}{n\varphi}$ had to deviate from φ for very large n ... so this was paradox.

13: Also note .15-16, proposed as to how to use a relatively low R for coding

14: $R = (22-25)$ to avoid many possible orderings of δ 's "at ϵ same point": say $M \geq n$ -R ϵ out
15: insertion pts. M^R doesn't work w. redundancy at each pt. where we had ≥ 1 δ .
16: M^R - which is hard to compute! another way: $da \cdot m(m-1) \dots (m-R+1)$ - different way to count
17: δ - we only allow one δ at a point, ϵ - but I'm not sure of details, already
18: On second thought M^R may be correct ϵ - if we insert a view δ just after ϵ inserted δ ,
19: or between 2 two symbols (adjacent to an inserted δ) so $\frac{M^R}{R!}$ ($R!$ different orderings of insertings)

20: $= \frac{(n-R)!^R}{R!}$ ways. If the pc of any one way is φ then pc of R occurrences is $\frac{(\varphi)^R (n-R)!^R}{R!} = \frac{(\varphi m)^R}{R!}$
21: Expected no. of occurrences: $\sum_{R=0}^{\infty} \frac{(\varphi m)^R}{R!} = e^{\varphi m}$ exactly which seems crazy!
22: φm would seem like a more reasonable expected value of R

23: $\sum_{R=0}^{\infty} \frac{(\varphi m)^R}{R!} = \sum_{R=0}^{\infty} \frac{R \cdot X^R}{R!} \quad \left| \quad X \frac{d}{dx} \sum_{R=0}^{\infty} \frac{X^R}{R!} = \sum_{R=1}^{\infty} \frac{R \cdot X^{R-1}}{R!} \right. = 276.00$
24: woops - a paradox value

4TM

S-Functs. 25

SN see 4TM 26 (4) for w. Paul's Grammar for ANH.

0:27340 + QAs as possl. - i as short as possl. Th. Gave at Ric's point is a bit vague, but perhaps acceptable. (Note 272.37). I'm thinking of the BU method of soln. {I don't know how Lsrah would work here.} After TM has done a lot of Ric's, we got it to do prediction once locally... but TM's main job is usually just finding short rules dir-Rules - finding a kind of "key" language for a growing corpus. From previous codes, TM gets lots of useful cones that can be used in practice in Mark problems. TM might be able to do Lsrah well, after Pius kind of training.

27 Re: Lsrah My impression is that Lsrah w. a suitable TSO is a good, reasonable way to do a TSO: T. main problem may be in presentation of S-functionz.

So, review recent work on S-functions: (essentially functional up to 272-20) Burroughs - Wheeler. etc.

One recent () idea was use of Bernoulli: Essentially 2 qst. use of PPM / Bur. Wh. etc. to find cones rapidly: Also look for subcones (≡ Entails). Also expansion of 2-1qt into CFG & CSG. This would give jobs of cones & not bad prodn. Could I then introduce more factorial forms of S-functions? One more general format "Context" is a function strings: It amounts to a function that tells if a string is a member of set of contexts. If an input. This is, func. is one way to dec. a context. (See 4TM 20.00 for discussion of how to create rules from old. Also has ideas on S-functions. An "S-gam" or S-objects is a "S-quest" which is a s-parameter & a p.d. in all strings. 4TM 20.08 ff discusses cross & comb. strings

SN on Sorting for PPM: One fast (cheap) way: use quick sort. Then, when new info coming, make a new file w. ^{row} Pz/dz sorted by itself. When creating new cones, use both files. The first (larger) file will usually contain all info used, but the second will occasionally be used also. For the second file we can use 2-3 lines or some expensive way to insert info, because of smallness of file it will not be very expensive in time or space.

Quick sort v.s. Heap sort.

25 SN One way to decide on how TM should represent & search for S-functions: Write a TSO in which includes a sequence of S-functions as solns. Then, for these particular S-functions list different ways they can be represented, & give heuristics for finding solns. This will be initially done for Lsrah, but perhaps similar functional forms & search methods can be used in BU method. (34)

PCP "end of input" part "end of input" is determined by experience w. N problems in the past.

Try listing all solns (to date) of S function generation / search problem. Also note first things I wrote on soln. of 3 input case. "P 33 or 133" - yellow paper. I think I found a bug in it - but I'm not sure.

How, I have Modda of 3 input case. form of S-functions that may not P. I use a functional form for PCP in being as small as possible, based on past experience. I know it is a kind of "special" kind of problem.

4.1.70 One general kind of TSO would be to kind of data used for 2-1qt: use 2-qt models, Pz, d, n, m, c, f, g, e, s, c, o, c, t. This gives rise to particular kind of model for S-functioning. How general is this? - How can they be generalized? These things models how calculating process, Pz, d, n, m, c, f, g, e, s, c, o, c, t. are automatically calculated, not "sah for". In (I know) process is most correct thing, usually such is always over-educate speed & certain patterns are calculated automatically (not "sah for")

See 302.00 for Good ideas on S-Functions

10:27:40 EN in -1.35-40, its not clear if it was binary of INV, IND or OZ probs!

02 In -1.35-40: The idea is to decide what population to use ~~problem~~ when starting a problem;

past populations consist of individuals and their scores) w.r.t a particular goal. (10)

In line w.r.t. idea of 1.00 Δ INV probs seem to be not so good, since they don't fit into an "hour misses" (i.e. Gene has only 2 values).

Perhaps Don't try to solve INV probs, but discover ways to convert them to OZ probs like "GPS" - w.r.t. vector Gene.

On the other hand IND probs using Lscho are a kind of INV problem, because of particular character of Lscho. We do get continuous Gene values, but are many of them, a piece

are usually not easy to compare because they are for different corpora... (Nrv. Reason objection)

may fade away after we work out a way to do GP on TSO's!

03 27:33 One way is to look at present problem: decide what probs in past are similar/related...

Developing useful variations of this sort is done by analyzing TM. In high TSO's

"redundancy" is a mut. relation. The "Compression trick" of Li-vit (others) is also useful, a

can be improved, when induction (=compression) methods improve - so it can be self-improving

technique.

15 SN In -1.37-39: Say we were (not) able to find a single OJ that did f. entire corpus, then to do a prediction, we could use the "closeness" of Q relation to decide

which OJ would be best for a particular new Q (prediction). This idea seems close

to see much older idea of having a function "looked a new Q"

to decide which functions to use to (help) advise.

10 Alternatively, for induction (or other OZ probs), one could assign a "default" value to Q A's for which one has found no code. - So one always has a code that complete corpus

E.g. smooth for some Q's f. OJ may give a pc to some A's. We assign a certain

a merit of f. pc of all "non OJ" outputs that don't include all A's - i.e. we use this p.c. for

for "default codes" for all otherwise "unresolvable" A's.

So-called Apparent BIG problems

6 1 How to use info from previous problems a problem solnt. (.10 is an imp way) (10)

Def 2 The Lscho v.s. BU method. - If we do both simlty, how can we show info best?

Solu. to 1) may be direct for Lscho v.s. BU.

3 (I hadn't noted this first time around) How to represent, such, S-functs.

Another approach to .10 if: "update" is not done until a (new Qnt) comes in! We then look

at Q's (22.11) that look related to Qnt, i.e. we pass grammar "Per their (Q's, A's)S.

32 (26) This can be done gradually; First we learn to solve individual problems - using, to some extent, info from previous problems. But at first, TM has no way to "recognize" a problem is decision on what comes to use.

The worst possi. interaction is that of usual GA, which: Even after he solves a problem (perhaps in talk, refer to try Michel or other account) creation inherent decision any kind.

In (33), the problems in the TSO are related, but a "solution" doesn't have to be able to solve all of

the previous problems - so we don't do the "conventional" QA TSO at first. At first, our goal

is not prediction! We are mainly interested in getting some common solutions for many of 274.00

combine? i?
(50) on SM...
could be...
to list of...
strategies
or a power of
type of
prob?
lower of...
> 1 times "best"
stat. process
< 1 moves forward
= w.r.t. all steps.

9/19/04

47M

20. (Spec): What are H costs for deciding when to put a D.F. at a particular pt. in the genome from the cond?

(Vagooly): That is essentially 2 or more Taken choices at that pt. seemed to be of high "G".

In Lsach, we will tend to search each branch pt. in a tree exhaustively.

In QA: Wouldn't it be reasonable to spend fraction of genome "working backwards" from A toward Q? This is certainly an O(n) search technique.

Say we have found, via PPM in Lsach, good codes for (Q_1, Q_2, \dots, Q_n) . PPM provides a D.F. over the set of $\geq n$ objects. We use this D.F. on Q^j 's to put a D.F. on $Q_{n+1} \rightarrow A$ But ... Press looks like a nice way to go.

Some Wo. Backtrack!

Since it's all possible in Lsach, we have a universal backtracking! How a D.F. is a Deterministic O(n) func, based on a corpus of $D \cdot O^j$ that have been successful for $\sqrt{z} \approx 1/n$.

I'm not sure that this is what I want! Consider a set of O^j 's that have been able to do $Q_i \rightarrow A_i$.

This set of O^j programs, induces via PPM a P.D. on O^j items.

If we use this test on ANL strings w. +, -, x, /. we would expect for e. ppm. (w. input $O^j, 1, 7$ {ops: + x - /} output $\begin{pmatrix} + & .25 \\ - & .25 \\ x & .25 \\ / & .25 \end{pmatrix} A_2, A_3$

We could continue to search for better pc's for corpus. If we allow "non-brink" even

T.M. could look at pairs of probes & successful sites. (which saves memory like BU method)

21 To forget stuff is involved w. very general search models & general search methods.

Lets go back to GP: it's a or 5 different & how I was going to fix them!

258.06 - 22 is a relatively recent peak. (0/5/04)

270.11 - 29; (low 292, 10 - 14, 20 - 40)

235.06 - 40; "improvements of GP"

205.341 Not many other low peaks could be fasted.

Using faster approx. of Gac routine, in early stages of GP. Also (correction) that is not of Gac routine, in early stages of GP.

- 1) Speed: Better models & model fitting
- 2) Profit from past problems & solns.
- 3) TSP
- 4) Lsach? or More optimized search routine like $M_{ij} = \frac{O_i}{G_{max} - G_i}$ i.e. other ways to model G as a function of Gac params.
- 5) Being able to watch phenotypes in action / fail / succeed. Being able to watch Gac being applied.
- 6) When Gac is expensive, use fast ruff approx. Mostly Generated & related (E1)!

29 So (1, 4, 6) (2, 3) 5. So More like 3 general cases of improver.

30 Try to outline what I would do if I had to do a TM in 1 wk / 2 wks / 4 wks / 6 wks ...

31 One simple way was to use GP (w. or w/o. various augmentations): Th. most interesting / imp. problem is TSP / learning from previous problems.

33 1) Use Lsach-like TSP: ~~learn~~ each problem in order, keep pop in / populations cards that scored well in past problems. We can do this by scoring to "the best" of each soly population or take "10k best" & select k from them at random. Pick & evolve target. Variance in Gac.

Note +1.00

37 2) ~~Evolve~~ solve Q_1, A_1 to Q_n, A_n individually, using previous populations & - but a soln to Q_i, A_i need not try to find common solns. After we have seen individual solns., we try to find common solns to pairs of Q_i, A_i 's, then triplets ... ect until we have a soln for them all.

273.32

-0.01 SPEC

0.270.19!

04

07

In the case of 2 level: If I have a uniform d.f. over 0,1 as source of f. elements of f. d can I will start to search w. 0 & 1 as possible trial cond. ~~etc~~ - best = for each pass, then we try the new 2 bit code of $\frac{1}{8}, \frac{2}{8}, \frac{3}{8}$. (I'm not sure how we start 0 & 1). if we next $\frac{1}{8}$, use $T \leftarrow 2T$ we do $0, \frac{1}{8}, \frac{2}{8}, \frac{3}{8}$ on next round - then $0, \frac{1}{8}, \frac{2}{8}, \frac{3}{8}$ on next round.

I'm not entirely of ~~course~~ $0.00 \dots 01$ - there are other symbols outside! - but this is search, & we do by physics PL order.

Note that Hertz did have a way of using GP to compute numerical problems to his precision. I just know did he do it? ... I think there is a Mt Carlo version of (0.00704)

If its not search, then we can just use a / s-grammar (can be GP population based s-grammar) to generate Mt Carlo code - using something like M_{ij} - ("Max Expected jump" $\frac{O_{ij}}{G_{ij} - B_{ij}}$)

A possible serious diffy w. GP: That one could offer genes to branches & one "climbs up the landscape". Normally, GP makes "decisions" at each branch & continues. To some extent this is "fixed" (or made less bad) by doing many GP runs out. Some problem - so a alternative branches until any be chosen.

Lsearch does have advantage of looking at all branches (up to current CB level). - but its looking for a somewhat different thing.

There are "Multiple Goal" GA schemes, it may be that we can use them to investigate "Branching paths" better.

I really think (1) branching should be addressed by a "s-grammar" for each of branches: sub-grammars would, of course, have many common cons, but would correspond to ~~closed~~ branching & ~~generating~~ generating trees.

In a search system, we'd want parents to be (mainly) from same branch, if we wanted to do a branching exploration off. "Fitness landscape" - How one identifies "which branch" the parent is on, is unclear! Would use of a PPM type "grammar" be able to deal w. "branching" of this type?

(0.00 - 0.10) (Note 0.09 - 0.10 in particular) does seem like a general way to generate search s-branches.

I'd like to link it to the Bernoulli ideas of $(\frac{269.34}{269.34 \pm 270, 27-29})$

In (0.00-0.10) we chose p.c. values & try them out on the corpus. In 270, 27-29 we chose elements to do the laps rule on, & we compare the p.c.'s. We should get about the same PC values.

So first, I want to be sure I know how to use 0.00 - 0.07 for Lsearch! Then ~~etc~~ generalize to 0.00-0.10; then see if it can ~~find~~ fit in (1.22-1.30).

In line (-0.01) I noticed that the discrete D.F. should differ from Gaussian ones: i.e. ~~etc~~

$P_i (i=1/n)$ p.c. of taking a certain feature. (I perhaps uniform d.f. modified by function) continuous is uniform uniform D.F. over 0,1 mapped by some function onto a different range so it gets a different (usually non-uniform) d.f. In the first cases the D version is a single token at first context. In 6. other it is a single real (or complex) no. at a context.

In both cases there are ways to find "best fit" for phrases that depend on D.F. In 6. Discrete case, Bernoulli (Laps rule) is in the continuous case, maybe something like: (say divide 0,1 into a few divisions - then use finer divisions (more bits) if necessary / warranted. (SPEC 272.00)

4-TM

00:269.40: f) A variant on 2.69. ³⁴ We have a Universal (with set), Σ so strings of its symbols are Σ^* "Lisp" forms
Grammar G is a (perhaps) CFG whose leaves are strings of Σ symbols \equiv words. PEWS.

G_2 generates CFG's. By looking at a corpus of Σ PEWS, we can deduce G_2 grammars for them. We get a density d.f. of Σ PEWS. I had to idea some time ago that this could not

~~give a~~ give a s-funct that was universal (in the sense - I want @ updated) a universal d.f. on s-functs. The idea being that I had a way of generating

So far, Peil was universal - so it would design $p.c. > 0$ to a $p.c.$ finitely decidable s-funct. (function universal s-funct d.f. as s-functs)
Well, Φ (.00th) seems to do this. - NO, it doesn't! It is a (univ. d.f. on s-funct) or D-functioning.

07 \rightarrow Any (d.f. Universal (or not)) on s-functs is itself an s-funct.

Any s-funct or d-funct is an s-funct.

Are .07-.09 deceptive? I.e., a s-funct (is. Pd.) on a function is a functional. f_1 of f_1 ,
is a kind of f_1 argument of f_1 - I'm not talking about a composition of functions.

13 In learning d-functions: I use defns \approx /o (Context) ^{original} for induction: to create
new kinds. If any old / primitives \approx are s-functs, then my old induction methods of
dedn \approx contexts will give defns that are s-functs \approx contexts that are s-functs. - we will end up
w/ s-functs of (possibly) \approx very complexity - probly all s-functs could be generated, ^{possibly} ~~some~~ ^{so the method}
is univ. How to do it? ^{3 PPM}

17 While .13 is surely reasonable, I don't want to ask for any examples of how
to use it! Well: could it use it for f. Bern problem of 269.28?

19 In APL is other similar (pass from is su instruction that complies 2 1 dim arrays
by $t - x \pm$ $a \equiv (a_1, \dots, a_n)$ $b \equiv (b_1, \dots, b_n)$: $a \circ p b \equiv (a_1 \circ p b_1, \dots, a_n \circ p b_n)$

I may want to use this sort of thing ^{or just as in PPM or augmented PPM} ~~rather~~ to assign p.c.'s to a list of "tokens" or whatever.
NB Primitives a common ~~base~~ ^{base} type ~~copy~~ in which one defines a context \approx P.c. automatically
defines a p.d. and talking tokens. In this sense, the defn of \pm context is a complete defn
of \pm s-funct. A context also can be used to define a p.d. on things lower than a simple symbol -

26 ~~A~~ token can be a name of any size - (or perhaps be a s-string \equiv Pd. on strings.)

27 In \approx ~~26~~ Consider a common ³ ~~method~~ ^{method} of induction: One does this solely by defining
names, negats, functions, contexts. Each of these operations can take a pc of corpus \approx

29 Can we promissory deduce objects or ^{names} ~~parts~~ ^{parts} of them in its ~~conclusion~~ ^{conclusion} functions.
 \rightarrow An arby complex OJ function can be part of as a special case of .27-.29 ^{Later! How?} \rightarrow ³⁵

30 Heuristics (for "suspecting" that a new case might be constructed in a class of ways)
are also constructed of old cases \approx "parts" of them.

Initially, all of this \pm can be done in PPM style, using BW x PPM.

Later we look for sub-graphs.

35 For .30: Perhaps this: That OJ operators on \pm strings of $Q_i(N_i)$ parts into \pm sq \equiv copy.

0/16/04

FTM

Spec. 20: 268.40 : A transition. 268.36 : Its extra expensive to say which symbols to change and so what: A better kind of S-function will have few symbols to be changed (marked) - which is like 268.28; but we could use a default diff for the symbols or use the d.f. implied by a recent context.

Just how to assign pc to forgo. it's unclear! Say I use a special symbol/int. code to denote the d.f. of several values. Say I have 4 of such symbols in my code for a cond. Then, to evaluate the pc of a "A" value that I need, I have to scan over all combinations of all 4 "special symbols" to find a comb. That gives the value A. The pc's of "special symbols" is 1, (best 20). The d.f. is implied by the "special symbols" gives the final pc of first ~~value~~.

code that produces the value A.

One might have a standard "Correction Code" that is used to "store an uncorrected A, into the correct A value.

How does forgo. relate to Maximize of O's functions in QATM?? Do the pc's of special symbols correspond to "R" values in the "simple model of S-functs?"

I really, I'd like to be able to take any O's function & find a pc (>0) for a corpus w/ that O's function.

T. problem seems to be: I have a few ways to represent S-functs: Now O's is one

of S-functs. What are good ways to search over a space of S-functs, w. goal of obtaining

O's that assigns (w. its own) by pc's to a corpus. I have a few such methods: 1) look 2) M.C. search 3) M.C. search with cross

Whoops! In general, one can't simply replace any symbol in a function, w. any other symbol: we have to replace them w. symbols of same "arity".

To study S-funct forms: Give many examples of S-functs. Plus could be solved

QA probs.

1) In early ANL, Tom has learned about +, -, ÷, x, but doesn't know which to use! - so for binary operators there are 4 posys. in some known express like α, β, γ (in PN) in which α can be 'add, sub, mul or div' w. probly. T. probs. transducers could be obtained by Bernseq. analysis of the corpus.

2) Consider predicting a simple Bern Seq. w. known letters & alphabet. We have a Null (= constant) φ.

So we want an uncond. pd. This seems like a very fundamental problem.

How best to express a k-way probabilistic jump? Perhaps have a "primitive concept" one

We have a random variable 0 ≤ t ≤ 1. Pick out k intervals on [0,1] so 1. random we choose one of t. k: w.p.c. (depending on) interval chosen. For the search process, we also start out w. k equal intervals.

3) One method of representation of S-functs and similar from is used of S-CFG's: start out w. a simple Bern model (approximate needed) - then hunt for production rules. T. "post" probs are computed by parsing & counting up of trees made choices made. This is usually reasonable approx. I'd like to formalize this-idea to universality. Sensitive Grammars can be universal, - but can not solve this is a good, useful representation.

4TH

so: ^{space} 265.90 : Done both Lsrch's BU method seems like a good idea, but it would like to link them - so they help one another! Common confusions, is one way. Both are guided by PD's!
 (265.90) ~~space~~ Now similar are these PD's?

Well, maybe to 2 PD's aren't so different! I have been thinking of D IAD probs for Lsrch, not (Stochastic)! If we do S-INDuction, then Lsrch examines causes in order of their QAs apart of causes. The causes, of courses, are subsets of the individual in the corpus.

In BU method, we have a PD on s-functs also. These s-functs (as do those for Lsrch) assigning a PD to each QA. Equivalently, for each Q there are PD on all poss. Assoc. #'s.

The 2 PD's don't seem so different! The Lsrch P.D. looks at "successful" causes (s-functs) & induces a D.P. on s-functs... The P.D. itself each one will be the "successful". Note that "success" is a soft quality.

The BU method D.P. also tries to rank causes on expected future "success" (soft quality).

The 2 methods seem to be looking for the same thing - a cause w. much success w. current corpus. They could, indeed, use the same PD on causes!

Ho. BU using MEJ, would use a certain subset of causes to generate a s-ling (≡ PD on causes). And do Mt. Carlo trials on this way Lsrch would simply try to list causes in that lang. in PC order. [Also note, it's not clear how the Mt. Carlo search of BU deals w. non-terminating causes]

Also in general, I've not dealt w. the softness of the search. Even for both methods, we just search & keep track of the "best k" causes. But for (one way sort) all of the causes found w.r.t. the score?

In BU, there are at least 2 approaches: One, classical G.P., uses Markov/counts to generate Mt Carlo causes. Another makes a stoch S-grammar to generate the causes. This last is closer to the Lsrch method, since one can (sometimes) use that S-grammar to generate causes in PC order.

T. forgo. suggests a unity of Lsrch's BU method: The generation of s-functs may be very easy: All you do is include 1 or more stunts as "primitives". Some useful stunts to use:

28 D Uniform d.f. of X on interval 0,1.; for discrete d.f. 2 mult, div have PC's = $\frac{1}{2}$.

This second is not so hot! The first, I am relatively d.f. I like w. suitable probs.

2 I don't see how to refer to "second"

[One A.H. way: Include parameters gives val. PC's & choices: The param. evaluation may not be precise - precision, but PC ↓ w. precision (costs more bits),

Another way to deal w. PD's on discrete params is to use 11 codes, - one code for each value of the discrete parameter.

A related approach: 2 (Practical idea! We have a d-gram) D that does Q → A. To get a s-distrib for A, 2 make changes in the D to get it to create other A's. The "info" needed for these changes is then added to the base of D, to give PC of the new A. A simple way to modify D, is to change individual symbols, words, pairs of symbols, etc. → 269100

4TM

ANALOG: Rite way to do ANN (Also RANN, but finding peak takes longer)

.03-.06 is rite way to do ANN (i.e. maybe RANN).

In normal feed forward ANN: what is normally done: The pick random initialization pts & recurse hill climbing until they get a peak on "test data". ~~Then~~ This random choices run many times, then average y. results

~~ANALOG~~. Unhappy to have "test data": Do as above but don't stop until ~~test~~-final peak is found. (Pts will usually be grossly over fitted). At each peak find "area" of peak in (initial) param. space (possibly by local M. Carlo trials). Do this for many initial pts & average, using

as WFS, i.e. local size of each peak, mult by it. Give opt. "training data" → 5 @ 293 20 on his

(diffy express)
A poss. objection to .03-.06: That almost all ^{final} peaks will be narrow & 2/0 have poor gene for trng. data. This show will be only a few "good peaks" & they will give quite different ~~and~~ products:
As a result, one will have to run a plum a very large no. of times to get an average value

Sample for good product

This To speed up convergence use that v.g. non-linear optimization method. (--- Mergart method)
Method, also gives "width" of each peak it finds.

In terms of no. of trials needed to find peak opt, .03-.06 looks (short it might be some speed (same no. of trials needed) as LSOCH(?)

→ RE .03-.06: Just how is it better than choosing parameters randomly, then using ~~some~~ product for

Best set of params by trial & error versus with that set of params? I think .03-.06 is more efficient, but it really has to be shown. Each peak has a "basin of attraction" There are probly many fewer peaks than the no. of pts we would need to get count accuracy via .16-.17. Her, this is true only if i. model is good enuf to give only a few "Basins of Attraction".

Anyway, forv. discuss, strongly suggests that ANN could be very useful for learning real functions. For learning Discrete functions, I'm not so sure!

SN I had this idea of using a ^{large?} (size?) Nucleus of an Atom as a "very fast" dynamic system, so we could get it to do computation very rapidly. One by diffy was getting ~~in~~ in for m & out of Nucleus. 10⁻¹⁴cm λ = 2.4 x 10⁻¹⁰cm
E = 2 x 10¹⁰ MeV
= 3 x 10¹⁹ eV
= 3 x 10¹⁵ GHz
One way into by NMR. — This/2 computation could change the Magn or electric (or whatever) Moment of a nucleus which could be detected by NMR. It may be that almost all nuc/reactions are very short lived, so diffit to detect — (i.e. life is << time of passage in ~~of~~ of nuc. in H field of accessible strength (say 3T = 30k Gauss)). Normal NMR (MRI) is ~ 1.5T = 15k Gauss. There are some nuc states w. very long lives (Years) ... but I think they are rare → ? But
Actually any nucleus, stable or not, can be regarded as a "nuclear state" ... because there certainly lots of nuc. states w. all ~~the~~ sizes of 1/2 (spin)!

ANALOG Compu: .00ff

DA: primitive operations w/ Add, Mult, integraten. Very reasonable

Analog Compu.: Differential Analysis, Digital D.A. "D.A.": Shannon paper on DA.
Paper by Chris Moore (PS: 10/13/02) | Also Moore's paper in 1996 on how to map dynamic system into a TMC.
Could NMR be used to probe state of nucleus, structures doing some computations? ~~Also NMR~~
If so, it should be poss. to use NMR as input to ~~simulate~~ nucleus.
Just (how) is a D.A. anal. "universal"? Can it (at least) do all Prim rec. funcs? ^(to what extent)

All DA's are of limited precision. Can we use a finite no. of them to get arby good precision
like Shannon's channel capacity theorem. (Actually Shannon was this book by Conrad & Winograd
on how to do this w/ ANN. I think, but I never figured out how it worked. McLaughlin was
involved in this to some extent. Shannon did some work on this, but I don't know if
he got any very useful results (He called them "crumbly" components)

One idea of "universality" that we have a fixed ~~set~~ set of interconnections; but
a set of ^{continuously} variable inputs. That any configuration can be simulated by a suitable
set of analog initializations!

To compare an analog value/integ a switching signal: Integrator \times (v.s. time).
It will $\rightarrow +\infty$ or $-\infty$ (depending on its sign). If we have "threshold" functions,
 $T(x,y) = \max(x,y)$ Then we can decide if $x > 0$ or < 0 . i.e. if $\int_0^{+\infty} x > 1$ then...
or if $\max(x,0) = 0$ then $x \leq 0$. $\max(x,0) = -x$ can be done by multiplying x by -1 .

To simulate a ^{non-connection} "connection" betw 2 operators, we multiply the "wire" by 0 or 1.
2 operators by 0 or +1.

Also, Moore is (perhaps universal) "circuit computing".

For induction it may not be necessary to have exact computations. Say we have slightly
"crumbly" components: we scan thru input config/param values until we get output close to
desired, i.e. we observe "next output" (\exists predn). We then repeat with same & nearby input
values, i.e. get D.P. of "next output". We may then continue scanning over inputs to get better
match of reproduction of outputs: (closeness to outputs can be RMS, ~~or~~ R window or x window).

How could we do ~~analog~~ Time series predn? Cross-predn?
^{subseries}
^{corpus}

There are 2 inputs to "Net": ① config of net ② input data. We also have to have a way to
"search" for "output" (perhaps have several poss. outputs & decide on "best one"
post-hoc). ~~for each config. input~~, we run several ~~net~~ ex. corpora & see how outputs
correlate w/ "correct" predn. Approp of an input confign would be related to (identical to)
the density of inputs or ^{input} Realization, but fine outputs "in a region"

ANN's are an interesting example: we can have connections fixed & vary inputs; or vary both
connections & inputs. Another large expansion is RANN (Recurrent ANN). In both cases (ANN, RANN)
most of the remarks of 20-23 hold.

[SN] In normal ANN (or maybe RANN as well?) we normally use first partial derivs to guess at
next trials. If we also used second order derivs & solved resultant eqns for extremum, we would
get much faster convergence! How to adapt RANN to this idea (or even first derivatives)
is not immediately clear.

00:264.40: Note that after 6. TSC of "Phase 1" TM itself has a choice of what induction methods to use: **So it should itself be able to decide** ^{on} Lsrech v.s. BU method - or how much (relatively) time to spend on each.

03 During normal "Adaptive Lsrech" TM develops an assoc GPD. Could ^{4. Same} GPD be used somehow, by BU method? - Or would BU method have to develop its own techniques independently of Lsrech's "GPD"? If another Method uses "data" or "status" they would be able to use one another. On trying to mix data used by 2 methods

06 Vary Imp. Q! Compare Updating process in Lsrech v.s. BU method. Also, compare such procedures.

10 In Lsrech, t search is always "standard" (say, $T \leq T$) - ~~that~~
In BU method, there are many poss. ways. - Tho In GAT @ MEJ, I had been considering populations in which we could find deterministic or stochastic comparable, approximatible con structures on t population to control G & G_c in a (known way).

X in strings for UMCS. It may use partially ordered trees as the data & as the BU from algebra. Is this useful? Some times cluster linearly ordered can UMCS be linearly ordered? Maybe, but produces not "effectively"

19 In Lsrech we have population ordered in probability that a individual will meet criterion; T criterion or INV being (for INV) problem - ~~that is~~ an acceptable code. (Yes/No). So induction becomes an INV problem. T. updating of GPD involves adding restrictions, control (conditional) pc of new symbol (it we can say generally). Updating can occur by adding positive or negative cases to corpus, or by statistical biological analysis of "trajectories" of system. (130)

20:242.35: Bifurcations Another approach to Bifurcations in SRT. If there is "boominess" of T then the cardinality of the set of UMCS will be too large: (i.e. uncountable) - sense of "Planus" (continuous).

19 Consider branching factor as a function (curve) of T in T .
The no. of T mc is countable & t . UMCS is a subset of T mc so its cardinality can't be part of t . Continuum.
We do then conclude by the "proof" of 242.35-40, which involved randomness.
So maybe a relationship betw. randomness and countal cardinality or ... ?

21:19 For $SS2=0$ we have standard ^{set} ~~of~~ exists for Lsrech: ~~that~~ in loop rule this zero set
Case is deal'w. by t . "precomp" Dist grows > 0 pc to a universal set of symbols/insts/relations.
For $SS2 > 0$, we look for various reps in corpus: Bernoulli firsts, then perhaps ngms, ngms, cut mult. rules for agmts, etc. Content dependence, Common Subtrees,
DNA type analysis of corpus, etc.

37 I did try to soften Lsrech Success Criterion by using sample sets of t . set of problems
TM is to solve (ie. from t . TSC). This way, each cond would get a score > 0 if
it solved $> \phi$ problems. Actually, it is complicated: its t . list of successful codes with pc of each code
Just 25 in .32-34; I BU method, updating can look for all kinds of reps in t corpus (which methods
@ soft'g' into (not just $\alpha = +\infty$ or $-\infty$). (.35-37) trk Lsrech does similar analyses.

A mean boost could occur if the set of problems is finite

0/11/04
4TM

try (mutation) a mouse w/ a tail to be solve the cut-off commonly were still being done.
(Licht Circumcision)

Anyway, GPS does look at testing of conds. to suggest Modifiers in "Next Round" of trials.
Lsuch can do this a little by looking at common characteristics of failures, or "Next Moves" by but not very by & conds. I think it does this in Phase 2.

Phase II uses (TM): Phase I does not. So in any Lsuch TM...
We do tm. to work for kinds of paths that TM2 would solve, but we don't yet apply this capability to TM itself.

Another feature of Phase II could be observation of Evaln. of Conds: so it can understand their success or failures. (It may not be necessary to put this in Phase 2).

SN perhaps look up { Elitism, Genetic Algorithms } Google. Search engines has d. covered optimization degrees of Elitism.

In Jaeger's OOBs: It may end up w. paths that have max pc of jumps to Gmax which usually will also be Mej!

14 2001 T. Method

Out of 2 basic methods of solving problems: Lsuch finds cond. $\frac{PC}{CC}$ taking time $\frac{PC}{CC}$. If it doesn't find it, it finds Nothing!

2) Common method starts w. A.H or Pruned search. But as a rule, it has a best found rules per. BU method any start w. parallelism of conds: Both A.H & Pruned search. BU method is "random" default. No! A.H & Pruned search are very quite different. YA. Greedy Pruned search may have some PC, but

TD method. Lsuch method. Top down method. TDM = Lsuch. Bottom up Method. BU method. BU M = Greedy Pruned search.

One Arg for Lsuch was that if one could find a soln. T.SQ was inadequate. In such cases, trainer has to find PC of system (more time, money) or if problem diff by "Hints". In real world, the trainer isn't around to give hints! However, in Phase 1, trainer is present - and any thing that gets thru Phase 1 (i.e. capability of doing Phase II problems) is o.k.

So, perhaps pursue both strategies sim. if. (Time share, perhaps) ← easily parallelized.

T. T.SQ is designed so that Lsuch works w. it. Whether TM can continue w. Probs that are not structured so that Lsuch will work, is unclear. If an Lsuch soln. is actually poss. w. 6. CC available, I'd certainly want TM to acquire it, because it is indeed (almost always) a v.g. soln. (i.e. vary by PC).

In RW TM will perhaps often be asked to solve probs in which Lsuch would not be able to get soln in available "time". However, A priori, one doesn't know this. If one starts out w. BU method, one will get some feedback on how things are going, by way of a in ~~phase~~ boost of corpus per cc. However, Lsuch seems to afford no such feedback! It either codes the corpus with doesn't code the corpus.

In T.SQ, if the trainer thinks that problems should all be solvable by methods of trainer or aware of, he should be able to write a T.SQ for it, or TM should be able to find soln within maybe 10 X ~~desired~~ "solncc" constructed by Trainer

0/9/04
4PM

203

20:262.40 : Optimality of LSrch is also tied in w. Subscript (2) in Sol 89. Perhaps to effect that

any conceivable law can be expressed as a modification of G.P.S. for LSrch. [I think one can realize

→ Prin: Any Legal law is ~~not~~ logically &/o statistically derived from observations on past.

→ As such, it should be possible to express results as modification of G.P.D. (same Q of this list)

Hvr. This picking up of laws will occur in Phase II, when TM is looking at probabilities of 6 Post

I recently in effort to write 2nd volume to comment on "PS Facts".

→ trying to use P's 1980 to speed up present's Srch. One Conclusion: G.A. is Adaptive LSrch can't look at Phenotype directly ~~or~~

1) So; write review of 589 PS ~~in~~ in include it in Commentary 589 folder. [Mention "unsolved" problem of

2) Note area of superiority of ~~the~~ (augmented) G.P. ~~over~~ LSrch. q.v. 262.37-38

3) Another impl pt. in G.P.'s LSrch is ~~to~~ distance in approach to IBD: (261.30-39) also

4) After doing 1), write a review of G.P.'s LSrch problem: this should be ~~not~~ get me back

"on track"!

N.B. MEJ is one way of modeling G as a function of Cond. down. Tho it seems to work

well for certain down methods, it is decisively not the only way to model G as a function of

Cond down! Another (not fully developed) way way to try to model some monofunctional $\delta(C)$

function of G as to pc of some S-lang. (S-Grammar). I haven't (yet) done much w

this last.

Also Note: Any expression of G d.f. as a function of cond. params. is an S-lang.

(Any ~~simple~~ p.d. is an "S-lang".) on any arguments

Actually, MEJ does model G as a function of cond. down. It finds a way of ~~the~~ MEJ,

Plan McCarb gives its population. This is a way of getting by G cond. Its similar to

US way a "selected" for MEJ to consequence size, m. (261.00-07) in ordinary G.P.

-1: 2.00.20: How to get Good MEJ function?

ISN On b. Universality of food for ANN! My complaint about the proof was that for a

function in N space in a unit cube of δ Δ^N points were needed to describe G. form.

Then if δ / Δ^N values are distinct to be simulated were "random", there is no better way! - argument

is with the constant of proportionality - a cofactor Δ^N . How much (if at all) does it

deviate from 1.?

STU on ANN: In fed. fwd ANN: conventionally, one does not repeated runs w. randomized

initial cond. Result is like "all cases" & cancels out pc of initialization params.

[A similar thing is often done in G.A. for same reason]

So it sounds close to a correct ALP method! Expensive, G.vr. - Not so expensive

if one keeps track of mis of data materials in sequencing/analysis & quite

when acceptable precision is obtained.

-2: 00 7.11 This Deficy of G.A. is Adaptive LSrch! - Not being able to watch cond being tested (continued)

So as to be able to suggest (corrections, changes, modifies, improvements) for next generation of search

It would be nice to do this kind of "Lamarckian"? Used example of Lamarckism is:

Tails of all mice are cut off soon after birth - results in their born w/o tails. If the "modifier"

were matching to tails being cut off, he would see Ray would have an utility & not issue them in subsequent

generations, if he put to cutting off would continue indefinitely - However occasionally

is wrong to say
have 297 and do what
conds. is taken
3F

3F

["good reason"]

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

also

4-TM

SM 100

SNESM: A "Non Parametric" way to do Time Series predicn! Consider 2 T's's!

X & Y: Y is used to predict X. We can use PPM, but we can use continuous way: To do k history production! at time t, last k values of Y from

Predn. at time t is $\sum_{j=0}^{k-1} X_j f(\sum_{i=1}^k (y_{t-i} - y_{j-i})^2)$

$f(x) = \sum_{j=0}^k [y_{t-i} - y_{j-i}]^2$ on esp. absolute value smoother (symm?) function

$f_x(\)$ is t. wt. of predn. X_j . It is a \downarrow funct. of error $f(x)$ data is one wants to use.

First depends on α . We optimize α for t. recent data in TSS.

Q: If this method significantly differs from linear predn. w/ k coils? Forb. Sequence of (t), I should think it would be

$f_x(z)$ write as $e^{-\frac{z}{\alpha}}$ i. unoptimized α also normalized α so

$f_x(z) = \sum_{j=0}^k e^{-\frac{z}{\alpha}}$ It would be best better than linear regressor. Hurry, it may not be better than ANN. Both are commonly non-linear methods

i.e. wts. should sum to 1. We may want to \downarrow f_x where as \uparrow \uparrow ... Rec would make normzn. possil. If \downarrow is exponential we have a "x window".

So $f_{\alpha,j} \propto e^{-\frac{z}{\alpha}} \cdot e^{-\lambda \cdot j}$ where λ is \downarrow decay distance.

Other forms of $f_{\alpha,j}$ might be better.

I might try it out on sequence $X_j = Y_j$: $Y_n = S(aX_{n-1} + bX_{n-2} + cX_{n-1}X_{n-2} + dX_{n-1}^2 + eX_{n-2}^2 + noise)$

$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}}$ S is a squashing function like $\tanh x$
 $\frac{e^{2x} - 1}{e^{2x} + 1} = \frac{1}{e^{2x} + 1} + 1 = \frac{1 + e^{2x}}{e^{2x} + 1} = \frac{1}{1 + e^{-2x}} = \tanh x + 1$
 $\tanh x + 1 = \frac{2e^x}{e^x + e^{-x}} = \frac{2}{1 + e^{-2x}}$ So $\tanh x = \left(\frac{2}{1 + e^{-2x}} - 1 \right)$

26140: A "true" Lsach would only consider candidate wave codes for output, and it would list them in order.

So the "true" Lsach has as input, what "real" Lsach has as output! T. did not bother. S- (only used for Lsach v.s. Prefused for GP is that they are arrived at in different ways.

"Restart": One path way to greatly improve GP: $258.06 - 259.05; 260.07 - .20$
Haskell on improving GP: Note Remark of 261.08-11 $-261.00 - .11$

Still, lets compare ordinary Lsach & GP on induction problems: Th. dozen of 261.30 - 35 is reasonable. In 261.30: Any "predictable context" is an (input kind of) "short code". Hvr, context = any kind has 2. $G > -\infty$ - so we want to get some "slope info": which is not true in ordinary Lsach.

The only "slope info" is obtained by "my dirty Lsach". To optimality form for Lsach is based on "All info is in guiding pd" - which implies "my dirty Lsach": - ideally one would use "dynamic programming to get info (by low expected G-fnals) to enable fast \uparrow of G-fnals, not long-run: (which looks like strategies for a reinforcement Machine!).

Hvr, I don't want TM to get into this "Dynamic Pgm" v. n. It should probably much concentrate on "short" problems, w/o. pussy (yet) of ~~the~~ deciding to do "self improvement" doing a problem. This sort of thing should be worked on by TM. Later.

37: First problem is to get to Phase 2 as soon as possible.

38: That said; is GP or Lsach a better way to do Phase 1? A bit of a problem MEJ is beyond.

Phase 1 - it "understands" optimization better than Lsach does in Phase 1.

263.00

4TM

~~is this one~~

LM386

0:260.20 ; T. Discussion of 260.07-26 was for GP using a corpus which index were within bins, or strictly sorted by G. An alternative approach: use tournaments to select top 2 for ~~next~~ winning. Say we get 2 children: we replace the bottom 2 of G tournament w. them (unless they are worse than the bottom 2 — in which case we discard ~~them~~ once or both). The size "in" of tournament determines "elitism" of tournament (see 260.35 for d.f.). We select 2 m roles. $\frac{G_1}{G_{max}-G_1}$ is max. How to get these values for population, must be worked out \rightarrow N: to test phase \rightarrow vs test of G_1 & G_2 are for the children, not population itself. (like 260.07)

07 (The "Gmax" is sort-population) \rightarrow just have track of G & G of cont. w. genes etc. like 260.07

08 In the fugg. models, the MEJ eq. was odd to try to get by G in offspring. It would probably give result? Press was better than current GP. Hvr, I don't know whether it would be very much better.

0 A key idea is to TSO — while it would be nice to improve it's norms (operation of GA), this improvement is not really needed for us to implement TSO's. It might be very much better, since MEJ gives it a much better understanding of OZ than "normal" GP or even "Lsuch in Phase 1"

10 All of ~~the~~ ^{recent} stuff has been to improve GP: What I've forgotten is to specify value of Lsuch for IND probs! \equiv Searching in G ~~set~~ (\equiv PC) order for a cond that will fit. An inverse way: To start w. known codes for corpus & try to find a shorter codes, i.e. compress them. This last is more in the spirit of GP — It is one good way to get Grammars for Corpus (e.g. ASG dir). We start w. Promiscuous & AA Grammar. This sounds more like what "people" do: i.e. look at corpus & try to find reps (compressions) in it. With experience, one acquires a large body of reps by pat & means for recognizing when a rep may be relevant (\equiv "Coding & Recoding"). (24) can be done by Lsuch also! We know ordering of trials (reps) that will be successful compressors. Presumably, as we code more & more of ~~the~~ tsq, our ordering of pc's of our rep trials will "adapt/improve".

12 How does GP normally do induction? One way is to search for productive concepts. A concept could be a rule subset of the corpus & a/o one or more conditions (obs) that corpus satisfies. Circumstances evaluated by empirical logic rules, using not only entire corpus, but "sampling". There seem to be 2 different methodologies: ① GP: A biased MC search for concepts at most poss G. (2) Lsuch: A search for concepts in PC order. MEJ is example of biased optimum biased MC

135 Some confusion in my mind on Press problems: Do some "rambling": For IND probs, Lsuch does a "hard" optimum search — "Lsuch" = soln, we know it's (for "pure" Lsuch) No (long) e. such). GP solves IND probs \rightarrow OZ probs — by MC search over biased S lang. After some subset prob have been solved, Lsuch is also over biased S lang, but it may be about S lang. Lsuch does to search what appears to be PC order (\equiv G) order... but not exactly. The search is for by G cond that fits exactly. (260.20)

So, for $m=10$, $\frac{1}{5}$ off: way down from max, we get $\frac{1}{5}$ decay $\approx \frac{1}{5}$ max down e^{-5} decay.
 The restriction $n \gg m+d$ may not be important. It holds for small d
 m is usually $\ll n$; The d can be almost as large as n , we are mainly interested in $d \ll 2005$
 maximize top; i.e. $d < \frac{n-1}{m-1}$ and maybe d as large as $\left(\frac{n-1}{m-1} \times 2\right)$
 If $m=11$, $\frac{n-1}{5}$ which is perhaps small and so that of $m-1 \approx 10^4$, $d = \frac{n-1}{m-1} \cdot 2 = \frac{10^4}{5} = 2000$
 $d=10$ is $\sim \frac{1}{5} \times n-1$ so t approx is probly still O.K.

At any rate, t decay rate $e^{-\frac{(m-1)d}{n-1}}$ for $m+d \ll n$ is looks very useful

07: 259.05 Put population in G local bins. Take individual from local X random nodes. Take individual
 Local X , random branch - combine: keep track of G locally of comb. structure
 X and Y , we can pit branch of an individual at a later node of some individual
 when we have another state, optimize $\partial t \partial x + \partial y + \partial x^2 + \partial x \partial y + \partial y^2$ (G parameters) so approximate
 \bar{G} Also approximate σ_G^2 as a function of \bar{G} . Then find X, Y comb. structures w. n past
 $\frac{\sigma_G}{\sigma_{max} - \bar{G}}$. (σ_{max} is present σ_{max} - not really in X, Y of a tree)

An easy way to get σ_G^2 : do 2 curves like so for \bar{G}^2 ; then $\sigma_G^2 = \bar{G}^2 - \bar{G}^2$

$$\frac{\sigma_G^2}{(\sigma_{max} - \bar{G})} = \frac{\bar{G}^2}{\sigma_{max}^2 + \bar{G}^2 - 2\bar{G}\sigma_{max}} = \frac{\bar{G}^2 - \bar{G}^2}{\sigma_{max}^2 + \bar{G}^2 - 2\bar{G}\sigma_{max}}$$

} No need simplify, other curve no need to take
 S q. root!

$$\frac{\sigma_G^2}{(\sigma_{max} - \bar{G})} = \frac{\bar{G}^2 - \bar{G}^2}{1 - 2\bar{G}\frac{\sigma_{max}}{\bar{G}} + \bar{G}^2}$$

if one normalizes G w.r.t. σ_{max} so $\sigma_{max} = 1 \rightarrow \frac{\bar{G}^2 - \bar{G}^2}{(1 - \frac{\sigma_{max}}{\bar{G}})^2}$

T . Numerator will be a quadratic in X or Y : t decays, t square of quadratic in X or Y .
 Since we have strategies on what $G \rightarrow G^2$ look like, It probly would be good to use ANN,
 However, since both \bar{G} & \bar{G}^2 are very probly monotonic in X or Y , it might be possible to use "MSE"
 259.22
 259.21
 261.00

21: error ≈ 11 259.22 we from 259.21

$$\frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}} \cdot \frac{\binom{n-1-d}{m-1-d}}{\binom{n-1-d}{m-1-d}} = \frac{\binom{n-1-d}{m-1-d}}{\binom{n-d}{m-d}} \approx \frac{\sqrt{\frac{S}{S-t}}}{\sqrt{\frac{S}{S-t}}} = \frac{S}{S-t}$$

$$\frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}} = \frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}} = \frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}} = \frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}} = \frac{\binom{n-d}{m-d}}{\binom{n-d}{m-d}}$$

For $d, t \ll N$:

$$\left(1 + \frac{t}{N}\right)^{N-d} \cdot \left(1 - \frac{t+t}{N}\right)^t \approx e^{\frac{t}{N}(N-d)} \cdot e^{-\frac{t(t+t)}{N}}$$

$$= \exp\left(t - \frac{t^2}{N} - \frac{t^2}{N}\right) = \exp\left(t - \frac{2t^2}{N}\right) = \exp\left(\frac{t}{N}(N-d) - d\right) = \exp\left(\frac{t}{N}(N-d) - d\right)$$

$$\alpha = \exp\left(\frac{t}{N}(N-d) - d\right) \rightarrow \exp\left(-\frac{t-d}{N}\right)$$

$$\beta = \exp\left(-\frac{t+t}{N}\right) = \exp\left(-\frac{2t}{N}\right)$$

So $\alpha \cdot \beta \rightarrow \exp\left(-\frac{2t+d}{N}\right) \approx \exp\left(-\frac{2(m-1)d}{n-1}\right)$
 I checked f. factor of 2 again $\exp\left(-\frac{2t+d}{N}\right)$ really seems right!
 Seems quite ok
 ≈ 17.5 is variation times
 2 factor indep of d
 for $n \gg d, t, m$

10: : ~~FRANCOIS~~ form of number by ~~some~~

Do 258.36 ff after one has adjusted X_0 via usual MEJ. See if 278.36 increases or

↓ MEJ. If ~~the~~ ~~of~~ m & n , MEJ worse, then ↓ α is scarfed helps

Then ↑ or ↓ X_0 is scarfed adjusting α can't MEJ (so move X_0 & α to try to ↑ MEJ).

25: (SN) Some Q's about G-F | (1) When z / ~~causes~~ ~~are~~ ~~related~~ for ~~branching~~, why ~~exchange~~ ~~branches~~?

If z branch is good, it should be freed on all reasonable places. → 260.67

(2) In tour's mem selection! What kind of DE is obtained for the choice?

Say population is n individuals, all of dist $\alpha = 1$ to n . We select m of them at random.

What's prob that i best of m , is dist d from top? $d \leq n-m$

Approach 1: Prob that i is d place is: $\frac{1}{n} \times \frac{1}{m} \times \dots$

$$= \frac{1}{n} \times \frac{1}{m} \times \left(\frac{n-d-1}{n} \times \frac{n-d-2}{n} \times \dots \times \frac{n-d-(m-1)}{n} \right) = \frac{1}{n^m} \cdot \frac{(n-d-1)!}{(n-d-m)!} = \frac{1}{n^m} \left(\frac{n-d-1}{e} \right)^{n-d-1} \sqrt{\frac{n-d-1}{n-d-m}} \left(\frac{e^{n-d-m}}{(n-d)^{n-d-m}} \right)$$

$$= \frac{1}{n^m} e^{m-1} \left(\frac{n-d-1}{n-d-m} \right)^{n-d-1} \cdot \frac{(n-d-m)^{-(n-d-m)}}{(n-d)^{-(n-d-m)}} \rightarrow \left(\frac{n-d-1}{n-d-m} \right)^{n-d-1} \cdot \frac{(n-d-m)^{m-1}}{(n-d)^{m-1}}$$

More exactly

$$= \frac{1}{n} \times \left(\frac{n-d-1}{n-1} \times \frac{n-d-1}{n-2} \times \dots \times \frac{n-d-(m-1)}{n-(m-1)} \right) = \frac{1}{n} \left(\frac{n-1-d}{n-1} \times \frac{n-2-d}{n-2} \times \dots \times \frac{n-(m-1)-d}{n-(m-1)} \right)$$

$$= \frac{(n-1-d)!}{(n-m-d)!} \cdot \frac{(n-m)!}{n!} = e \cdot \frac{(n-1-d) \cdot (n-m)}{(n-m-d) \cdot n} \cdot \frac{(n-1-d)^{n-1-d}}{(n-m-d)^{n-1-d}} \cdot \frac{(n-m)^{n-m}}{n^{n-m}}$$

regard n or m as constants: so just

$$\sqrt{\frac{n-d}{n-m-d}} \cdot \frac{(n-d)^{n-d}}{(n-m-d)^{n-d}} \cdot \frac{(n-1-d)^{n-1-d}}{(n-m-d)^{n-1-d}} = \sqrt{\frac{s}{s-t}} \cdot \frac{(s-t)^{s-t}}{(s-t)^{s-t}} = \sqrt{\frac{s}{s-t}}$$

$$\left(\frac{s}{s-t} \right)^{s+\frac{1}{2}} \cdot \frac{(s-t)^t}{s^t} = \left(\frac{s}{s-t} \right)^{s+\frac{1}{2}} \cdot \left(\frac{s-t}{s} \right)^t \cdot s^{-t} = e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t}$$

$$= e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t} = e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t}$$

$$= e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t} = e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t}$$

maximize $e^{-\frac{1}{N}T} \cdot e^{\frac{1}{2}(s+\frac{1}{2})} \cdot e^{-\frac{t}{s}} \cdot s^{-t}$ so look at $(N-d)^T = \left(\frac{N-d}{N} \right)^T \cdot N^T$ consider

for $s \gg t$ is $n-t-d \gg m-1$; $n-d \gg m$ - $n \gg m+d$ the $m+d \ll n$ m and $m \gg n$.

Which is not very sharp decay ↓, since $\frac{m-1}{n-1} \ll 1$. Say $n=10^9$ m and m be 10^5

On the other hand $\frac{d}{n-1} \approx \frac{d}{n}$, it's fraction of way down! So: $\frac{d}{n-1} \approx \frac{d}{n}$ where $\frac{1}{n-1}$ is a fraction between 0 & 1

if $\frac{d}{n-1} \approx \frac{1}{n-1}$ where $\frac{1}{n-1}$ is a fraction between 0 & 1

4TM

It may be that Phase I doesn't have to be very smart in order to get to Phase II: It could be slow in Ph 2.

So T. general plan is to use GP to learn a T-SQ (probly Algebra). Some diffies are

1) Getting GP to do T-SQ: (2 recan + Batches of notes on file) -> 252.10 - .40; 253.08 - .11; 253.12 - .21

2) Using PPM and GA for search, (Not so much of an unsolved problem). So 253.08 - .40

A simplifi. that may be - 0.9: 253.27 (No explicit use of R (recapitulation) function).

So! General Plan: Quick once over!

1) Design GP system that can work ALN & similar problems. After each problem solved, we ~~add~~ add 2 (or 3) new

2) After we have population, we implement PPM w. auto population. See if it works.

3) Next do $X_0 = \frac{\sigma_G}{G_{avg} - G}$ See if it's a rate of G ↑. (≡ G')

4) Expand H-SQ: Try harder Alg. problems: Try prob's ≠ Algebra. See if we can get it to invent "R": (??) or perhaps give it "indices" so it can categorize easily, then have it learn to do categorize: (or use A.N. Compression method for categorize?).

in 5) I have to get a PPM working first. Try a regular BZZ, but break it down so I can get pc's directly. Make Modulus in direction of My own version of PPM. See if it's better. Perhaps try variants on Exp text, on Lisp codes, on / Lisp codes from solus of GP problems.

Expand outline of 106 ff, looking for diffies; look up relevant refs in my notes & other papers.

5) Advance T-SQ to Ph 2 type induction problems.

6) Phase 2.

SN: There is a version of B-trees that's designed for large Data Bases using HDD; Could be used in GP?

SN in GP, If we can get G to ↑, try reducing X_0 (2 fair Am's). Arr., this may be done automatically by a system! See 25 to ↑ $\frac{\sigma_G}{G_{avg} - G}$

Another, perhaps better, way to get similar effect: Arranges to have def. over G value of ^{entire} population. Set λ window, w. adjustable width: with wronged ~~parameters~~ max $\frac{\sigma_G}{G_{max} - G}$

Note first: Part of we use "undo" to Gave, G should be "learned". In case of induction, I don't know if pc is linear or maybe ln pc? - or neither. For small λ in G , ln pc is not (so) important.

Since non-linear is "linear" for small jumps. Max MEJ → Max expected jump MEJ

There's some Q out. utility of the farg: After all, all of the s-lays used have universal outputs so any ^{arbitrary} could be pc > 0. Hrs If we could legitimately get a type MEJ by modifying s-lay in any way - Pencil Prow!

Boroughs whether Lisp distrib; involved in Computing pc's w/o generating H. Carlo trials; Look at G of a cont. int. population. & priori reject it w. probly $\alpha \exp(-G/\lambda)$ so large α & small λ fully of rejection. To save time, do several trials for each accepted G value.

A Definit. possl. way to ↑ variance: For s-lays based on "Lisp" rule. In computing numerator, denom. of prob's; Lisp will give a certain constant for $\leq z \leq 0$. By increasing these constants, we flatten the distrib. We can ↑ constant num, denom. by some constant factor. This effect is minimal when $\leq z \leq 0$ is large. A way to do it for even large $\leq z$: ↑ ^{constant} denom by factor α so that it's = to rest of denom. Then a constant

0/5/04

257

FTM

256.90

>: ~~256.90~~: Its not clear as to whether to always use Reis context conc. for predn.

Also, I'm not quite sure that Cover's analysis applies to my analysis of PPM & its variants (like use of ~~sub-trees~~ sub-trees as contexts).

SN There may have been 2 versions of Reis analysis: one in which we use "context" term. PPM:

in the other, we just make defaults of inputs which are supposed to include to token to be predicted.

This last is more like Z141.

→ One big difference is that in PPM, m , the no. of default codes, = no. of contexts, is very large.

Also Reis m is w. comp. length. — How the main difference is that I don't really do all those in codes!!!

For a particular predn, say I consider 10 contexts. I only consider those 10 codes for that predn. Hrs, there are an enormous no. of other contexts that did not occur — they are the \bar{a}_i

parts of various codes: There ~~is~~ is Reis Enormous no. of \bar{a}_i codes that I don't consider.

— If I did, they would have very large w. w.r.t. to 10 "positive" contexts that I do consider.

These \bar{a}_i codes perhaps average to give a p.d. in the base of occurrence of various

tokens.

So, any one of these codes has a w. of 100 w.r.t. "default code", if there were 1000 of type codes, then we end up w. something close to default code for s. entire corpus. (no real compression)

So, if we did it exactly with many codes, the Cover analysis would be correct, and we would end up w. very little compression.

The way I had in mind — i.e. 07-09 is not directly analysed by Cover's analysis. How Good my method is, is an empirical Q.

254.05
also 254
254.06-20

SN DONT FORGET **DNA** analysis techniques. ^{They} Looks like ~~can~~ can be useful for

recognizing repeats of type (255.00-04) in which a ... c f g occurs several times in a row w/ default ^{possibly} names base "a" "c". This is a very simple kind of repeat in RPN terms of functions. Could this fact be relevant to its importance in DNA analysis?

SN Re: Repeats of type (255.00-04): If functions were created so that their first arg. forms a useful name w. function names, then we have a useful rep. Perhaps any random set of GA such would usually end up w. mainly functions of this kind, because they were useful.

On the other hand if we have a binary, non-commutative function, we simply don't have the option of using a default order of arguments. One way to fix this is say w. "Subtract" ^{non-commut.} $a - b$ can be converted to $b - a$ then $sum\ a, neg\ b = sum\ neg\ b, a$. We can do similar thing w. division (non-commutative). Hrs, I'm not entirely sure this will work! The idea seems to be to have all functions to be commutative. So we would have "recip." using functions to implement subtraction & division.

4 TM

10: 255.40 : In Cover and King see pp 255 $\frac{1}{2}$, 255 $\frac{2}{3}$. This is dis. looking, because up to now, I considered

The wts. that I was using to be optimum: They are far from correct & but equality (using of single best PEM. — and a constant mixture ratio may do better than best single PEM!

I don't have a "ready arg." at all to decide what to do. 255 $\frac{2}{3}$. 07 is falling suggests/tells how to do: "constant" wts.

It looks > (if I use ~~the~~ Multilinear prodn. of former Series! Maybe there is a matrix inversion method to get good wts!

We have a seq. of vectors, $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ in components, giving the pcs of E. count that did occur, as usual by the m. means.

We want a simple normalized vector $\Rightarrow \prod_{i=1}^n \vec{v} \cdot \vec{x}_i$ is max

$$\frac{\partial}{\partial z} \prod_i f_i(z) = \left(\sum \frac{\partial f_i(z)}{\partial z} \frac{1}{f_i(z)} \right) \prod f_i(z) \Leftarrow \text{is this true?} \text{ — possibly, since by independence.}$$

$$\frac{d}{dx} f_1 \cdot (f_2 \dots f_n) = f_1' \cdot f_2 \dots f_n + f_1 \cdot (f_2 \dots f_n)'$$

$$\left(\frac{f_1'}{f_1} + \frac{f_2'}{f_2} \right) f_1 f_2 = f_1 f_2 \cdot \frac{f_1' + f_2'}{f_1 f_2}$$

$$\text{eg } \frac{d}{dx} \prod_i f_i = \sum_i \frac{f_i'}{f_i} \left(\prod_{j \neq i} f_j \right); \quad d(\prod_{i=1}^n f_i) = \left(\frac{d}{dx} \prod_{i=1}^n f_i \right) \cdot f_{n+1} + \dots + \frac{d f_{n+1}}{dx} \prod_{i=1}^n f_i$$

So its true.

Back to .09!

$$\frac{2(.09)}{2V_j} = \frac{\sum_{i=1}^m X_i^j}{\sum_{i=1}^m V_i \cdot X_i^j} \Rightarrow \frac{X_i^j}{V_i \cdot X_i^j} (\prod) = 0$$

$$\text{likely } \sum_{i=1}^m \frac{X_i^j}{V_i \cdot X_i^j} = \phi$$

Well, w. Log's m. o.lts, we have a bit more: $\sum V_i^j = 1$

$$\sum \frac{X_i^j}{V_i \cdot X_i^j} (\prod) = \lambda = 0 \text{ and } \sum \frac{X_i^j}{V_i \cdot X_i^j} = \frac{\lambda}{\prod} = \text{constant.}$$

I don't see how to solve this: $(255\frac{2}{3}, .14, -.40)$ looks better!

8: $255\frac{2}{3}, .29 \Rightarrow$ For m params do I need to solve m^2 matrix elements (symmetric)?

In t. usual n-linear optzn each reqn., they may assume all off diag values = 0!

Now this may be good & not so easy perhaps successive approxns

To approximate $\prod_{i=1}^m f_i(z)$ as a function of \vec{v} , we need constant & linear as well as quadratic terms.

I think the method was simply usual constant (more quadratic terms: for m angles = $1 + m + \frac{m(m-1)}{2}$ terms)

Somewhat, the optzn method works if we do not use crossproducts! — surprising! — Both Mult. Reciprocals and

Gourskand's may discuss this.

Anyway, I think the best point is that my old method of combining 11 codes was not so good & — that

it usually ended up w. t. "Best" code having almost all wt.

Say I'm using various production & exp. (E "contexts"): I have a best 5'm using. Each has a distinct p.d. over the next taken possy. Other than give = (or use = wt), I could invent a new sub. series. Co-occurrence of items: use 155 p.d. of the best as a p.d. for the next taken. One problem is SSZ.

2.1.0 : Even if the P.E.M's will have about 3/4 entropy, yet any mix up of them will under noise.

I. (Lower no. of them we mix, the loss rate so → "true" PC.

Hor. in EAC, the individual P.E.M's has rather different entropy - yet mixing did 1/2 entropy to L. Not of "Best Guessor" !

A case where simple mixing certainly ↑ entropy: One of them is ^{character} correct PC.

All other pairs to be mixed are uniform dist. toward alphabet. So in long simply "1/2 chance"

the correct D.P. is moved away from correctness.

Consider Bay method of using: $z_n^{(i)} =$ wt. of i^{th} pair at time n ; $z_n^{(i)}$ is adjusted so that

it is more or less (constant) value in test, the resultant PC of ~~some~~ corpus would be Max.

Hor. finding $\vec{z}_n = (z_n^1, z_n^2 \dots z_n^m)$ seems very difficult. We can make rapid approxs

by sampling ϵ corpus. Small samples, from ϵ only ↑ ~~error~~ ϵ size is as slowly

zero into a good \vec{z}_n function. T. closer we get to max, the larger ϵ we use,

the slower we have to make changes. Usually, our sample will be most recent

data. For m P.E.M's, we can look for a Max. in m dim space.

A good way to find this is to Normal Non-linear open method
poss. by using my ϵ distribution (random) ϵ slow sampling of each pt. of m space ϵ ϵ

For a Best Approx, use ^{mean} / main PC as weight of each pair as their wt.

For King Cover Garry of 12 Guessors: PC ratio varied from ϵ ϵ

$.275$ to $.179$ — 2 factor of about 2. But a factor of $2^{-6} = 1.52$ seems more reasonable for approx. (12)
 2.445 to 3.732 or $\frac{2.445}{.409}$ to $\frac{3.732}{.267}$

(13R)

(2.9)

In ϵ sample of $.04-.06$, the correct PC gains a certain

pc/symbol which ϵ ϵ 1. The others will get less pc/symbol, so they will get ϵ ϵ

So .14 will not be exactly right. Hor. for a larger corpus, it will be found that

giving the correct PEM most of the wt. will be ϵ ϵ max pc of corpus.

Hor. the method of .07-.13 doesn't feel how to deal with a very large set of PEM's

Practical comp. of \vec{z}_n : It's easy to do a run of maybe 100 randomly sampled ϵ points for

perhaps a entire corpus. Then pass an optimum quadratic form thru those points. Use

the parabolic form to get a good guess for \vec{z}_n ; then use distributed/random trials

to see that ϵ ϵ pass a new quad form thru ϵ ϵ data — ϵ ϵ its use, etc.

I've forgotten to detail 1st .24 #! — 256.28

Objection of .22 : It doesn't fit into AIP.

Another from for .07-.13: Instead of using entire corpus to optimize \vec{z}_n , use ϵ ϵ

on corpus: T. length of window ϵ ϵ must be optimized. Using random sampling

of corpus is much harder, hor. — it's all ϵ ϵ !

Comments re King-Cover paper: EET - IT status, July 1978 p 413.

(1) Conclusions: p 439 how they get 64% redundancy is unclear.

log2 27 = 4.755 bits/char. 1.25 bits (char is (1 - 1.25/4.755)) = 1 - .262 = .737 = 73.7% compression

64% compression = .36 x 4.755 = 1.71 bits/char. Only 2 out of 12 subjects got wrong Run Pct!

(2) P 418 II last ff: Letting alpha^(1) = 1/n for committee; means every man of their b estimates for each bet. That it should end up better than "Best subject" is surprising! Perhaps.

(Cover for his co-author) has shown it to be useful/utility always true! -> (3), also (26)

(3) If "Jensen's Inequality" was available online, I could run BZZ on it, using my lower case. - In both test texts, there was only 1 sentence to be predicted, so that the period & double space at end of a sentence, do not occur.

Or, get it from library & scan in a few pp! -> (5)

(4) turn (2): Finding the best alpha^(2) (i.e. (1/n) alpha^(2) indep of n) would be easy (I think this is what I did in my "optimal portfolio") & does a v.g. yield. Was Cover's Univ. Gambling scheme every way over all alpha^(2) values? (again sum over all).? - Any way, if alpha^(1) = 1/n did v.g.,

& guessing it to be optimum would do better! - i.e. better than 1.25 bits/symbol.

(5) try doing ~~it~~ even better than "Jeff's bet" on BZZ directly, see how much compression it gets! I need to make a file w/ cops or GRLP's. Perhaps create one via Basic in some file area, so I can load it into BZZ. Perhaps Delay doing this until I have a version of "BZZ" that I can control better! This may be available online - by Google. Doesn't seem to be - maybe Markov Cover.

For first letter of a word, use freqs of letters of words for context 'space' - i.e. don't use default description. For second letter, use digrams & freqs for starts of words (Plus BZZ may do this normally)

(6) On second part, I shouldn't expect BZZ to do much good w. such a short corpus. I wouldn't know many n-gram freqs. Using 7-gram book would be useful, but. A human wouldn't read such "externals" into, because he has experience/memory of general English text from R.W. Perhaps start out w. large English text from Calgary corpus, or just guess at English text that BZW used.

(7) Hvr., the way I would do it with still be better than "committee". My method is to calculate each b(x_i) = sum_{j=1}^n b_j(x_i) * 2^i / sum_{j=1}^n 2^i

2^i = product of probs assigned to correct symbols up to i by all estimators

I think this is not so simple giving 1/n to each estimator & letting them bet independently & summing, - which is "diverge Capital" method. - Tho, for a long corpus it approaches both the all on best estimator, since "best guesser" gets a trust all w/ for best later on corpus.

As is, total capital for committee was 12.4 / 1.5 = 8.26 times best for Best Guesser

"Average Capital" method was Best Guesser / 12 (EM)

Note p 418 (of King/Cover): Column II (upper middle): they do consider "M4" method. They show that it is -> Best Guesser is n-gram, which is correct. Hvr., he seems to prove that & yield is 1.25 bits/symbol to that - I can't see anything wrong w. f. proof!

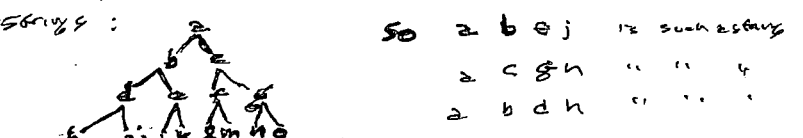
It is disturbing: That one can do better than "Best Guesser" - an example:

Sup the true pc's are P_K^0 & P_K^1 is P_K^0(z) = D z^1 / D D = radix, w. can have in a dist P_K^1 w. P_K^1(i) = P_K^0(i) + noise, then P_K^0(i) = normalized P_K^1. -> 200

: In a PPM analysis of RPN functions: ~~could~~ consider $a(b(d,e), c(f,g))$ (cf 254.08):
 $c \& g$ will be regarded as following $b(d,e)$; rather than "one of 4 args of a ".
 on the other hand $b(d,e)$ is regarded as following "a" ^{being} which is a clearly more relevant
 Can I find a way to (store annotate, ...) i corpus so that both of these kinds of regex
 would be recognized?

A pussy: to include ~~the~~ i. corpus coded w. order of each function's args (reversed/permutated),
 [I did consider this as one of 4 simultaneous problems of the corpus ~~for~~ PPM]

A slightly different tack: Regard all unidirectional sequences ~~that~~ there could be as leaf criteria of



But, if we ~~would~~ want to retain all such strings; for each leaf there is one such string: So
 the no. of nodes in our corpus is $[(no. of leaves) \times d]$ b. depth of corpus. Normally, b.
 no of nodes in our corpus ~~is~~ (for binary branching only) is only $[(no. of leaves) \times 2]$ ^{Not a gross it is} ~~nodes~~ ^{nodes} ~~many~~

In generating a string (top down) we say: in view of the previous virtual context (.07-14)
 what is the left token (rt. token) of this context? After we decide on left token, we may
 use may not visit left token as partial context for diff. or ~~other~~ token. Would be composed of substrings of
 (.07-14) be suitable to do .15-.16? (Not including "lateral" context, perhaps)

Proof of realizing .07-14: each node (other than for node) has auxiliary table of
 "left or right". This enables us to get left's ~~for~~ for each context for left & rt. arguments.

So: A lot of PPM-type regularity: ways to compress:

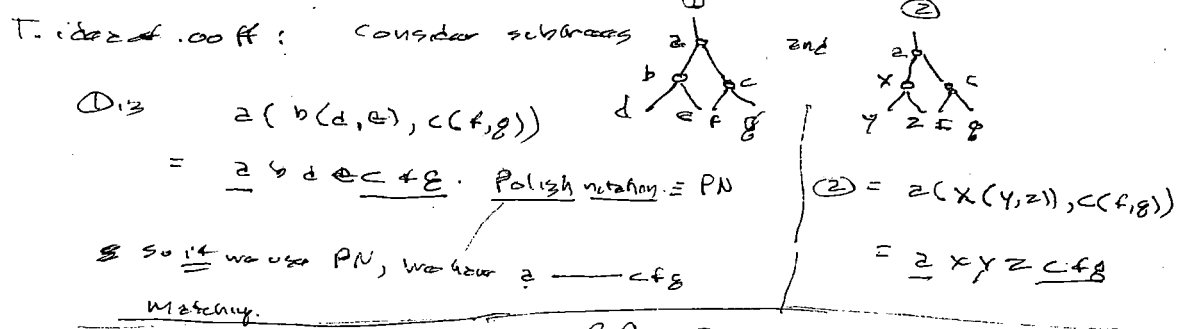
- 1) Normal PPM regex: perhaps pussy my possibly more exact parallel codes: see 227.00-.01 for refs
 of deans. of those methods & 227.00-40 for discussion. of their complexity.
- 2) PPM using substrings as contexts - (see refs of .23 for deans, examples) - also Note (4) (.28)
- 3) For PPM of function trees (Lisp or A-Z): Use RPN then PPM: Also use PN & Normal PPM.
 A cardy + analysis of 254.06-.20; This is a guide to permuting order of args & not changing from RPN to PN.
- 4) As a perhaps cheaper way to do 2): do (.07-.20)⁵
- 5) Use Z14 & make definitions:
 " " & expand it to do arguments & eventually C's & C's.
- 6) Try say, etc for any backward. (This is not to say we vary order of all arguments).

A BIG Q is how much time/money it takes to do these predictions/compressions.

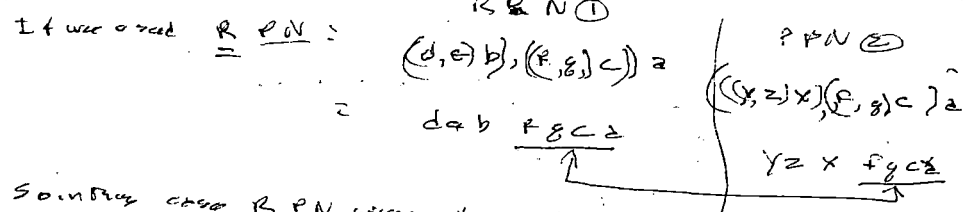
ATM

SN T. kinds of regys detected for DNA analysis (matching of substrings w gaps)
 seems to same as kind of regys found by common substrings. (= common sub-functions
 if focuses Polish RPN)! Also note that Pems for DNA analysis are trying to get
SPEED! - So they may know what I'm looking for.

Also, Pong suggests that DNA may be better understood as a notation
 that decodes a composition of sub-functions. → 257.20



So if we use PN, we have a — c f g
 Matching.



So maybe case RPN gives to similarity by simple suffix matching (PPN)

In any lang function word, if we suitably switch both PN & RPN,
 we could detect any subword as a pure suffix.

→ Is it the same as switching order of arguments in functions? Yes exactly the same.

So N.B.
 Keeping track of
 Corrus in RPN
 & PN does double
 storage cost, but it
 also does a no. of
 regys found.
 T. problems - does
 not affect a note
 of entropy & to
 justify a m.c.c.
 Swaps arguments of
 both PN & RPN? .36
 N.B. For my purposes
 RPN = PN with
 order of arguments
 reversed.
 That's it!

SN In G2D, we generate M. Carlo trials from a corpus of words that
 have their empirical $\sigma > \sigma_0$. Could we generate M. Carlo trials that
 would recognize substrings? .31

SN My Paris Paper "A New method..." is 1960 (or 1959);
 Instead of a "teacher", look for
 candidates regularly in ~~the corpus~~
 entire corpus: If it compresses f. corpus, it is a good rule to try!
 Hrr, rules "discovered" this way must be averaged so one can occasionally
 "back track" (= Undo them) on them, since they will occasionally be
~~wrong~~ non-optimum.

1: .23 It would be instructive to compare various methods of generating M. Carlo trials from a corpus.
 Methods compressors of various kinds would be ~~various~~ subclasses of one method.
 Strictly speaking (this is not necessarily a sample of one), but both are true!
 It would seem to be easier to make such an analysis if Lisp w. RPN is used, than
 if other languages & notations are used.

36 .15R It may be that as we get closer to goal (ordinary RPN, say will begin to "flatten out" a bit
 this point, it becomes useful to start using both notations of corpus. Analysis leading to /124.33
 is relevant (probably).

4 TM

SM

Spec. 10:25:40 decision of the 02 problems and learn to correlate them w. contexts (or other things that help generate cards) in the ~~cards~~ cards that have successfully solved those problems,

23 of rev.

SN Use Top of Page More! If there are any new imp ideas, or imp. continuos of old ones: write note on top of page: Give back refs. if relevant. This helps much in reviewing!

06 **SN SM:** Use External Contexts of stock characteristics (industry, earnings, etc) to help group stocks that likely cross-correlate. Otherwise, there uses very large no. of candidate "mixers" one will "not" it. So this External Context idea enables one to focus a much larger set of cond. in detectors.

08 **SN** Random Notes on Review of 250.11 & ff:
I think recent ideas centered on diffy of getting GA to do TSG's: T. trouble was, in GA, I'd like to include to by G ^{population} of all previous problem so far. ① Asst. TSG gets larger, ~~but because~~ this population gets too large. ② Hrr, another problem is that as TSG grows, we cannot have 22

2 **SN** A poss. way to deal w. this ~~problem~~ For to enter population of cards, there will be an optimum $X_0 \Rightarrow \frac{dG}{dG_{max}} = \max$. As ^{the} population $\uparrow G_{max}$, but $\bar{G} \uparrow$ only a little. G_c may be about constant. So for each problem, its population of cards w. $G > X_0$ will be about constant \approx to total population of cards will be \propto no. of probs in TSG.
To deal w. this: ^{random} sample of ~~cards~~ cards from each problem w. $G > X_0$. One can have, say 100 sets of cards for entire TSG on 1. Disc. Between problems, one ~~loads~~ loads a random ~~pop~~ population from disc into RAM - so between problem one has \approx (usually default) ~~same~~ population for induction.

21 **Pe. 07.13** See 239.07 for details of X_0, G_{max}, G, \bar{G}

22 **11** Eq of time to test each card on all previous problems. ③ A third problem is: Using f. entire past cards w. \approx wts. to solve a present problem ~~same~~ factors to throw away much imp. info we would like to include contextual info, past related "v" problems int past.
This less can be solved in a couple others ④ For QA problems, include external context (distn problem, diff eqs, organic chemistry, given time by Prof Z...ect)

27 also recovery can be, ... ⑤ Don't solve this problem: Let TM discover O_i functions that ~~examine~~ examine each problem & decide what ~~added~~ concepts to ~~include~~ include in it. In solving it. ^{we may try to find ways to encourage TM to find solns. of this type, by including certain primitives in it.} ⑥ Solver's Piz way! TM looks at problem & past & solns to learn, or by other trials

0 Try to correlate ^{features} of problems w. ^{obs} concs (include primitive & deduced sequences \approx contexts) of ~~parameters~~ parameters. This could be done by "Entropy analysis" that Li-vit used to discover "categories" music, plants, literature. Among General Solns: $G_{max} \approx Q_i$

We want $A_i \Rightarrow$ t. entropy of $Q_i A_i$ is min. Hrr, we may want to use \approx copies, t. set. $\{Q_i A_i\}_{i=1}^n = S$; then we want $A_i \Rightarrow$ i. entropy of $S \cup Q_i \cup A_i$ is max or; more exactly, we want e. entropy distribut $S \cup Q_i \cup A_i$ over all poss A_i - then we are using interested ~~in~~ A_i of low entropy.

Hrr, if we are considering $Q_i \approx$ (partial) context for A_i , we must use a more general context that ppm (convex trace & unnormalized context). ~~So~~ All of Q_i is t. context & by abstraction units & depth context.

9/26/04
4TM

20 (SPEC: 250.12, 250.29): Another trouble w. 250.11 ff: Its for INVS Phase I Induction; not OZ. This may be o.k. i.e. we don't need to do OZ until Phase II. The GA-like systems do OZ problems well, they are not able to learn from previous OZ prob. solns.

22

14

250.10 ff was a particular method of dealing w. diffy of Testing all new cond. out. Entire Corpus of problem thus far. Another Method was the R system of IDSIA § 1.2 p.10 ff. 2 ff. 12-18; 24200 ff discusses this: Also soft R functions: statistical testing - variable One output of procedure 250.10
Another approach would be to use user to update rules of IDSIA § 1.2, ~~update rules~~ or as is, or generalize parameters to update rules to be used for "soft" (≅ probabilistic) R functions.

20

Another approach to "improving GA" was modifying the selection cond. down to its Gave (≅ fitness). One may use a S-Grammar to ~~evolve~~ try to approximate Gave or $\delta(Gave)$ for some measure of fitness. The criterion of 239.07 ($\max \frac{G_r}{G_{max} - G_{observed}}$) is a better ~~approximation~~ ^{not use} utility criterion for a grammar. (Superficially, it seems to ~~approximate~~ $\delta(Gave)$ from δ criterion, but it does not. Both G_r & δ properties a good criterion depend on linearization of G . After an Grammar model was formed, Monte Carlo trials were made, & it was found a few new. Good cond., $\frac{G_r}{G_{max} - G}$ will change a/o we know to find a new model in view of new data. (loop to δ) until satisfied. T. Grammar has to be form that enables Monte Carlo trials

20

21

.20 is the kind of OZ solver in the spirit of GA. It doesn't seem to be v.g. at TSQ (true, but). Its not clear how to use it for a QUAL TSQ or QA problem, because we would have to test each cond on entire corpus to get G values that were comparable for diff. cond.

A method that's a lot like .20 ff is 234.19 ff, 234.21-402, 239.00-14. This uses a PPM to make the population and to G values. We end up using the criterion $\frac{G_r}{G_{max} - G}$ (239.07). But its easier to compute PPM (20-30) & to grammar/grammar to find. (T. Grammar of 20-30 is potentially much better: it can ~~be used~~ ^{used} S-CFG or even S-CSE ^{sensitive}).

SN T. main reason GA is poor for OZ problems, is that ~~we~~ it needs a way to look at the problem & decide - form its past experience, what a good S-Gramm. would be for generating cond. This real decision could be made by a function that is created by "standard" GA methods or by any of the methods of 250.11, 252.10, 20, 31. One poss. technique looks at features in s.

4-TM

0:250.40: we can get our randomness affected by specifying a variable in the interval $[0, 1]$. The p.c. of this specification depends on the precision we need, but ordinarily, we will consider all values of x in that interval (as giving rise to all codes for f). A_f being dec'd (i.e. \mathbb{Q}^A , T.M.)

→ T. funcy. may solve a very general Given Probability problem! Say we use start w. x being output on $[0, 1]$. We can use any function we like to map $[0, 1]$ into any interval, & get corresponding probabilities. T. pc of this function (i.e. x form) will depend on f , function ~~itself~~ - i.e. (how freely it was used on the past). In particular, if we need a mapping on \mathbb{N} integers $0, 1, \dots, \infty$, we can use this mapping to get it.

0 — This seems to clarify a lot of my thinking about \mathbb{N} & dependent positive integers or a d.f. on \mathbb{R} (s from $-\infty$ to $+\infty$). It does mean that a d.f. is pretty zesty — depends on the pc of dec'd by f func — which in turn, depends on the freq. of use of that f (or on the freq. of use of "components of f function") in the past.

14 But: What about Discrete P.D.'s? We could have a bunch of random bits (or Bytes?) as part of a function dec'd, that would give a probabilistic mapping. Now, we would like it to be poss., to go from a finite dec'd of deterministic params to include a distribution on an ∞ of poss(-outputs). — presumably, an infinite source of random bits is needed. (This is perhaps \equiv to "R" input to f (i.e. \mathbb{N} of \mathbb{N})).

0 — One way is to have a p.d. over the positive integers, \mathbb{N} (see .07-.08), which gives a pc for every finite binary string. T. mappings between uses are usually not critical: for integer x , we can use "Reservoirs" $A \geq 1 (1/2^x + 1/2^{x+1} + \dots)$ or $\frac{A}{x^2}$ ($x > 1$) or $\frac{1}{x^2}$ for a finite no. of x 's. The binary string x can be one of the "primitive inputs" of f . $A \geq 1$ function (i.e. LISP). would this really be poss.?

We could have > 1 "x" type inputs. Having many such R funcy doesn't cause trouble: we can just use a few (by pc) bits from each of them if need.

25 Another (apparently "Not Bad") method of getting discrete P.D.'s in $A \geq 1$; T. s-functs defined ~~specific by given conditions~~ or by specifying certain parts of its function tree & leaving others probabilistic. Unfairly this gives only a finite no of possys. — not all possys

29 A General ^{complaint} of .29: that the methods dec'd in 250.30 to 251.30 do not take all concrete outputs: i.e. they give pc of ϕ to some ^{some usually} most outputs. — (Clearly a Non-Universal pt of view!)

32 → One poss. direction: S-functs have no termination symbol to start off! They specify Looks promising! possible s-functs T. vector of bits (including a termination symbol): i.e. $\in \mathbb{R}$ inputs.

35 [In the discuss. of .14-.25 the input variables can be of any length, so they could, implement be used to dec'b. any thing.

250.30 & 251.35 give some newish ideas on how to implement universal S-functs. When I actually Need to use s-functs I will probably have ideas as to what they should look like & looking at the funcy. will suggest ways to implement them

9/26/04
4TM

REVIEW

PD'S on S-Funcs 250.30-251.40

Review of Approaches to TM: 11-29

2:248.17: "Working new problems" - depends on what the problem is: INV, OZ or prediction (= Revenue).
After GPD has been updated: INV problems can be solved by LSrch. Predict problems simply use the
current updated pd towards products: OZ problems. well, these are the main problems!
The, TM could start out in QA problems only. - Don't try OZ probs until Phase II.

2:249.12 -> Do this now, w. Final data set would not target: More DOIRS!

1) IDSIA: Phase 1: "Problem Pool" (248.21-249.10): In "steady state" we have a large no. of
solns. to the QA, TSCQ of IDSIA report. The TSCQ idea is not to summarize "Problem Pool" idea.

3 Any way in which a new problem occurs, we generate a "most likely" bunch of conds from: Corpus,
(using PPM to get a P.D. LSrch to get "most likely" conds. T. prod is wtd of these
conds's probab. d.f.}. We save these "most likely conds" to use them for trial updates.

7 To update: we first know the d.f. w.r.t. the entire corpus of successful conds "to
generate conds. We best know only QA (to be updated) only. When we find a
bunch of them, we put these conds in the corpus & generate new conds (w. some extra wh. on

the newly inserted conds), in attempts to get a cond that will work with existing corpus,
when we find one (is probably more), update is done. (→ new problem → .13)

things for
MTM
(= d-induction)

The .13-25 is not done as "Problem Pool" approach. .11 we will use Problem Pool at the
beginning of TM's final education, to get a large enough no. of good solns added to
in our corpus, so we can do .13-21

One poss/ trouble w. 11 ff is that we will be doing "statistical testing" of conds, solns are normally not
very certain that the best cond. really works w. all problems to date. This may screw up "MTM lmg"
- but may be O.K. for most NMTM lmg (= S-functions).
I'm not sure about best Update Algor, but .17 ff all has some good ideas on how to do it: By finding some
conds relevant to new problem (at least) - putting them into "update corpus", then Monte Carlo search - then
adding to update corpus, etc. → 252.00

29 SN PD'S on functions: extended PPM (extended by tree concepts & other extensions
of the "Context" idea) gives a P.D. on functions. If the primitives included some probabilistic
functions, the output could be a d.f. on probabilistic functions. To what extent can (extended) PPM
be used to model a Universal function? - (Both d & s are universal funcs).

Well the func is a spec of the AZ log over which can be made universal.

→ ALSO NOTE! We can probably get universal s functions by including one (or more)
s funcs as primitives: e.g. say one of "primitives" is a uniform d.f. on [0, 1].
We may use more > 1 indiv sources of randomness - but it's likely that more sources
are needed to get universal search funcs. If we have say only 1 randomness source, 251.00 spec



9/25/04

4TM

21 among known cards, there was a good chance of finding a new card that does "small". In general, we should try hard to find cards so that each problem has at least one (or preferably ^{2 or 3} more) cards that can solve it.

23 The situation of 248.37-39 is not so good - try to avoid it by making lots of "quick" fords on "hard" problems - see 249.37-38, using the p.d. generated by ~~the~~ cards that solved other problems.

If we continue to have trouble with several of the problems, try using cards that are not as closely linked to the present. "Most successful set of cards". This might be done by using data from less successful cards. I might be able to implement this by accepting data (in PPT, say) from less successful cards. (This is equivalent to X_0 in 234.19-22, 37)

12 **SN** Make a list of 4 or 5 Approaches to TM that I'm interested in, - Discuss advantages, disadvantages of each. See if I can use them to fix one or two (250.10)

13 **SN** In deciding on the entropy of string b wrt. string a: The usual way is to find entropy of ab & subtract entropy of a.
An alternative way that may be better: frame, no problem viz PPT is a probability of a symbol following any context. We use this to find the PC of each symbol in b; ~~the~~ T. in of total PC is it. It differs from 1.3-1.4 in that in 1.3-1.4 the corpus to determine PC of next symbol, includes the text of "b plus its". In 1.5-1.6 we do not include "b" info in the corpus used to create the "pd for next symbol".

The "problem Pool" approach of 248.21-249.10 is an idea I had for an advanced method of Machine Learning. The machine, given a pool of problems, decides itself which ones to work on, when, & how much time to spend on each. These ~~parameters~~ allocations were to be made by TM2. So, in its "Phase I" I would guide it in its choice of problems to work on.

All the systems that I have this phase I phase II aspect; w. TM2 at work in Phase 2. In "Idia", Phase II may be a little bit out. So I can use any type of Phase I (including conventional GP, say), and use the induction techniques learned to work on Phase II of zug of the model.

I vaguely remember a discussion of a distance based to GA approach and I use Phase II: Whomver this? 241.34-40, 243.00-05: I guess the difficulty was that GA would work for normal GA problems, but not for TSOs: It didn't seem a way to retain info from previous problems. If I simply included by G cards of previous problems in the corpus - the answer is, the problem of finding a card that works on all of the present & past corpus.

What I was considering recently (248.00 ff)

20:24.790: Consider ordinary GP as a u.g. or problem solver: A more difficult way of to do TSO's is that testing would take more time: Each cond would have to be tested on whole past TSO.

A way to fix this for GP: Each problem is solved individually / ^{new} ~~cond~~ is tested on next problem only: The TSO enters this way: For GP, the population includes ^{several} ~~cond~~ of ^{reasonable} ~~cond~~ past successful ~~prob~~ solutions. We can keep pop size ~~small~~ ^{reasonable} by only keeping 2 total of say 500 in population. As ~~cond~~ ^{new} problems are solved, we delete (older / ²⁰ less good) members of the population store S.

A (perhaps) better way is to store / ~~cond~~ ^{store} Grammars of population. If we use PPM-like "Grammars", we have to decide which conds to include in the "population" - a/o what Gores to assign to them.

Another way is to use PPM: store only 1000 ^(or 25 or 50 or 100) ~~conds~~ ^{solutions} to problems (Racoon's heuristic). To solve new problems, parties use LS to generate ^{"solutions"} ~~cond~~ most likely "solving" function. - or ~~use~~ ^{use} a wtd sum of ~~cond~~ ^{most likely} functions. } ?? → 250.00

- Just what is the problem?
 - ① I want to store / have available solutions ~~cond~~ used in solving of previous problems in order to conduct an intelligent search for solutions to new problems.
 - ② I want to be able to narrow down a set of abs most likely to be useful for a new problem ~~cond~~.
 - ③ I want to be able to update both ① & ② after each new prob. is solved a/o after ~~cond~~ a "correct" soln. is obtained.

As part of ② (perhaps) to be able to store info from a ~~cond~~ very large # of previous probs, you be able to narrow down likely relevant things (as in ①).

Another Tech: (I've written on this approach before): We have a large set of sample problems. Our first goal is to get O's first run work as many of them as possible. Usually it's easy to test. ~~conds~~ because they fail ~~cond~~ on a first few problems tried (at random).
Say there are 20 problems in the pool, we have several ~~conds~~ that can solve 20 problems. When we test a new cond on random problems from 1-20, we can ~~cond~~ tell, from 1 to 20, if solved v.s. unsolved prob. Thus far - whether it's at its solving 2 or more of the problems. ... We can use this as a rejection criterion for new conds ("sequential testing").

For small values of "2" (say) we will have to do all 20 or so problems before we reject a cond. ... But mix is 0.4. Each test has a good chance of discovering a cond w. useful "regularity". As I see it now, no matter how many trials of probs we have to make, each of these trials is useful & has possy of yielding import. info.

As the testing proceeds, we will find certain problems tend to be hard to solve, so we will test on these first, to eliminate more conds as soon as possible. (No at certain points in the search, we may want ^{own} more conds that can solve the "easy" problems, some can get a better P.D. to generate trials that will solve the hard problems.)

→ This is a bad situation, see P > 49.03
When we have many conds that can solve 18 of the 20 problems; If its ~~cond~~ some 2 problems that are unsolved by all, then we will use PPM for our P.D. on new conds to try to solve the 2 hard probs. ^(see 249.03) If v. 2 unsolved problems are different ~~cond~~

4TH

0 (5000) (244.18): In working on "S" (stochastic) problems, using soft R funcs: For S problems, testing can be 1 technical. —
 No need to test all cases. We only get a "probabilistic certainty" or a "probabilistic Question!"

There is an optimum amt. of time one should do statistical testing — it would be a "sequenced test" w. no "a priori" numerical cutoff... cutoff pt. depends on how fast it comes. Again, the thresholds of "certainty" would be determined by an analysis like our loading sub.

Equ. of 124.33

06 This "statistical testing" idea may be VERY IMPT... if it works! → (21)

T. necessity to test a new O^j on all past problems is clearly impractical. The "hard R's" is one way, but I think soft, statistical R's seems more reasonable,

more generally applicable. (142)

(5N) in .00ff is almost all of TD's comput., only a certain limited precision is needed. It is very wasteful to use 32 bits when only 4 are needed! IDEALLY we should design our computer so it was able to use only whatever precision it needed and was able to take advantage of fi saved CC in all over CC economy.

4 (27) Oh so picture this! We have a "front-end" part looks at problems & assigns (probabilistic) to a PEM. This is essentially "Phase II" of CDS's Report — starts. It does not count. Same as "soft" R funcs, for problems like ANL — over an S-version of ANL! — where we assign pc's to various eval. funcs of Alg. expressions.

T. way "S" R funcs work: we input problem, T. front end assigns it, probabilistic to a Probabilistic function that assigns a P.D. overall pass A's. P.P. to ff (§ 1.2) of CDS's: on updating: Has ideas on how

to do it w. d-R funcs. They would certainly have to be modified for S-R funcs. P_j Accepting & Modifying S-Func by

21 (106) On "statistical testing": If we make a mistake by ~~accepting~~ not testing & confirm case Q_i, A_i — this is not too bad, because we will eventually test that case on subsequent problems — the only other P_j has been further modified.

Furthermore testing P_j on only some of fi QA, means we will have to decide how good it is by averaging its pc's for a subset of QA's with fi pc's assigned by several different P_j 's. — So fi. Good is not very clear. If fi. pc for all past accepted P_j 's was $> .9$, say, one can easily feel it's ok. new P_j meets this criterion. Also, one might have ^{calculated} mem in pc per A... to use P_j 's criterion for testing new P_j candidates.

0: (20) I guess this "level" O^j function is ok, w. stochastic R. A non-st. approach would try to optimize O^j directly. By elzng O^j into R's P funcs, we probly have to give separate criteria for updating each. Hur, ultimately, the system will probly settle on one P function that does all the work. Since this is ^{basically} Non-st., we wouldn't be really giving long term restrictions out. System by P's R, P divisions. A big Q is how to update R's P funcs. Another Q: when TM "non-st" itself into mainly one R — How do we ~~update~~ test? — how do we avoid testing on all past QA's? This problem would be "Automatically" dealt w. if we use Q & TD's doing fi. updating is improving to update method.

4TM UMC's and Simulation Strips.

10: 242.40: As for t. no. of branches due to duplication (in P2a (idea for Min Sim strips?)), t. corresponding $\frac{M(N)}{N} = 0$, for $N \rightarrow \infty$ $M(N) > 0$ — but I don't know how to show that $\lim_{N \rightarrow \infty} M(N) > 0$ or P2a $M(N) > 0$ for any N . — (i.e. Prob. Prob. area > 1 "Min codes".)

03 of $\textcircled{1}$: Is the thru of 242.38-40 (i.e. $\lim_{N \rightarrow \infty} \frac{M}{N} = 0$) related to the P \neq NP problem? It has to do with rate of growth of duplication programs. $\textcircled{10}$

05: 242.40 NB: Actually t. thru of 242.38-40 says that t. sum of t. no. of branches due to either duplication & "simultaneous duplication" ("simultaneous" means "same length") must be low $\Rightarrow M(N) > 0$ but $\frac{M(N)}{N} = 0$.
 In agreement that t. sum $M(N) > 0$ is for not same length duplications. — "Thus far"

Actually t. $\frac{M}{N} \rightarrow 0$ prim is Non-informative surprising Because as $N \uparrow$, "t. no. of UMC's simulated" so one would expect t. no. of "duplicating a previous UMC" branches would become more likely!

$\textcircled{03}$ $\textcircled{04}$: A possibility related (to P \neq NP) is t. idea of Prim. rec. functs: rather expressed as a finite no. of nested loops of certain kind(s).

Further general ordering relation can be expressed as a Boolean Net. Would it be instructive to look at t. sim strip from this way.

Or t. Branching of t. sim. thru UMC's.

10 (24.40) ^{space}: Also note that L such is mainly for INV problems; one usually doesn't solve INV probs such, but usually convert them first to OZ probs.

Is the proof of optimality of L such ^{applicable} to OZ problems?

SN T. proof of optimality of L such is very \sim to "proof of optimality of Occam's razor: Both are counterfactual "All info is in Φ GPD"

SN The Early sections of the IDRIA report have a lot of reasonable-sounding ideas on heuristics, & ways to design TSP's: Perhaps read this a very once in a while!

An objection to "proof of optimality of L such"! That it doesn't allow for frag (mobility of GPD) during such. Well, if GPD changes during L such, it could still remain all available for it - such: Hvr., to "Plan" it, such

property, something like "Dynamic Programming" is needed. Say "Dyna Prog" into it's m. PD. Could this be

How could this m. "modify" L such "optimality"? Say one of the choices in search is to use "Dynamic Programming" to do calc on optimality ~~could be~~ Is this a meaningful way to ask?

20: **SN** Tree contexts are "cc expensive": ~ 10 to 100 times as much ~~space~~ (probly time) as normal PPM contexts. Hvr., we can balance PC corpus \uparrow v. b. cost of this update system - yielding an ^{equ. of} supply ~~is~~ like that ~~that~~ ~~lead~~ ~~to~~ / 124.93. i.e. decide how much time to spend (contingency) on updating.

A v. way to ~~avoid~~ avoid this cc is by Definition. A definition need not be substituted for its sub-tree (or whatever) every time that sub-tree occurs. T. substitution can be done when one is "fairly sure" it's appropriate. Also, ~~this~~ this substitution is usually reversible w.o. much ditty. - One may want to ~~avoid~~ ~~substitution~~.

"Unsubstitute" certain nodes when responding to corpus.

It will certainly pay for me to spend a fair amount of time trying ~~to~~ appropriate as a shortcut ~~to~~ to implement tree contexts.

Another perhaps equiv. & perhaps better way (if context definitions) is Zick and CRG discovery & Contexts. Grammar discovery.

4TM

See

10: 243.40 Will be few problems, this will not take long. Later, X_0 will be large and we will test on many problems. Hvr, most will fail on ~~code~~ ~~enuf~~ so that their G 's will be $< X_0$ so we can quickly discard them. Only code w. ~~close~~ close to or $> X_0$ will take much testing time.

We may isolate certain problems that have been difficult. Must be a way to use them for testing first or "early". For large X_0 , we want to find failures as fast as possibl, to discard (code w. $G < X_0$) rapidly, keep records of number of failures for each problem. (Presumably it is as $X_0 \uparrow$.)

At certain pts. in TM's training, it may be possibl to test ^{on} certain problems & do some random tests out rest of t problems. For a code that is at or G_0 or $\geq G_0$ we test on all problems.

Hvr, TM should be able to work on a problem of ~~high~~ finding code efficient ~~methods~~ methods of this sort

I could give it an initial corpus of more efficient search methods - have it look for common ~~seqs~~ seqs / subroutines in that corpus.

In what way (if any!) ^{was} my approach to ANL in IDSIA report better able to deal w. this problem, Para 243.27 ff (GA)? In IDSIA report, I had these "Recognition functions" that would look at a problem & do a much more "focused" search for soln. Because of the "Digital" nature of "R" functs, there was no need to test solns on ~~many~~ any but a few of t problems. "No need to test a soln. of a t problem on all chemistry problems."

It would of course, be possibl to devise "soft" R functs, but they would not eliminate ~~the~~ the necessity to test all problems on new codes. (Well - I'm not entirely sure!) - See 247.00! With a "soft" but relatively sharp "F" we could decide not to test an algo on certain probs because it was very unlikely that the algo would be tried on t problems. One could note that if algo hadn't been tested on a certain domain & so one would not expect to use it in that domain

22 W.O. Suitable testing: This is if kind of PSM's operating on problems. Each PSM is, say, usable for only a certain set of problems. If we modify PSM, we only have to test it on t set of problems if it is restricted to.

We would, hvr, like to have a PSM that can work on all problems (deciding, it's self, which sub-progs to use or which external "R" to "CALL")
Hvr, in ANL (1.2) we have a d-problem - we need a d-func to map algo expressing into such "Kules".

Note: In ~~the~~ ~~1982~~ report, ANL was being sold since INV problem, using Lsch. I find this is wrong approach! An INV problem should be given a "Continuous Gorc & sold as an OZ (or GFS) problem. Occasionally INV probs can be solved by Lsch - but this is unusual - only done when we can't assign a reasonable "ray scale" Gorc. 243.27 - .40 is one way to solve ANL problem into a "multilevel Gorc" OZ problem. The it's a somewhat general method, it doesn't find a v.g. Gorc ... It's Gorc doesn't really understand what needs to be done - (or, more exactly, etc "understanding" is not v.g. "P.12 end of § 1.2")

SN ① Noted in our "Analysis" (Breakdown / Unification of R functs) in (IDSIA report): These are fairly big in an essential part of Blumen "Sci Method".

② Lsch is GA (see $\frac{G_0}{G_0 - G}$) seem to be powerful, but essentially definite, a predecessor to prob. solving.
Preceding max is "like" Lsch; but not quite (analogous to res (both population, is not).

Also note "Proof" that Lsch was Best Possibility, we have a number of neighbor 2.

Spec
2:24.40.5

For, ~~GA(B)~~ GA(B) does limit us. (Crazy types macros) of 24.26-37. This may, indeed, be good enough for most all problems, but we do want to system to be able to discover arbitrary reg's.
T. BEST way to do this is to give TM the problem of improving its own Regy discovery Algm.

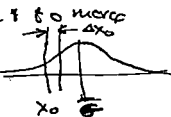
GA(B) is ~~not~~ (w. suitable TSO) almost certainly able to successfully work on such a problem.
→ IN Designing the TSO: I will know what kinds of reg's I expect TM to devr. — So I should not approach this problem w. such triviality!

Another way is PSG Discy + Hoop we devise definitions, subgrammers, recursions. We can have Context dependence to ↑ power of Grammars (perhaps, in theory, Universal).
I could... just introduce Reg's as an improvement of "PPM". — But it would probably be much slower than "PPM".

Re: "CSG's" To detect some contexts is very easy! First obtain CFGnum. Then, when corpus is parsed, keep track of "contexts" for each substitution (Just as in PPM). If a context cond. has an ↑ pc of the corpus, Do it! ... Perhaps, hrs we don't want to context to make its assoc substitution Mandatory, so, like normal CFG's, we want to be able to parse in alternative ways.

Well, the sub rules can be mandatory & still one could have alternative ways.
— Depends on what one means "By Mandatory": In CSG's we will have alternative production rules, just as in CFG's. So: If a production rule is used in CSG, we do have to always respect its context (if relevant).

2:239.18 → BN Ont. "On line" adjustment of X_0 : After one has just chosen an initial or a new X_0 :
Pick a ΔX_0 "region" R $X_0 < X < X_0 + \Delta X_0$. Keep track of $n, \sum G_i, \sum G_i$. Do this for all $G > X_0$: but also keep a separate file for those \geq quantities for region R . From this data, one can tell if $\frac{\sum G_i}{G_0 - G}$ is greater for the recent corpus with a w/o. "region" R . So we can tell whether we want to move X_0 to $X_0 + \Delta X_0$. It may be poss. to estimate ΔX_0 from a Theoretical Model.



27 Consider 24.34 : Start TM w. GA(B)
We begin w. a try ~~method~~ for which we know acceptable answers (not necessarily answers).

Doesn't seem to work w. the problems at the beginning of IDSIA Report. T. diff is: As I give successive problems, I want to see each problem & solve that problem & all preceding problems.

One way to think about it: Say my ANL ~~is~~ TSQ has a bunch of problems (Scheduled) ordered
positions: I start out w. a corpus null corpus: So probles of all units are solved, & no context dependence.
Next I try to solve problems in the TSO, individually — I select as many as I can. Then w. the best corpus as corpus, I use "PPM" to try to solve more ~~new~~ new prob's ≥ 0 to the problems in the corpus — w. k as larger ~~prob~~. (a fitness function like in GA).
So I have this corpus of units: Each has solved to no. G of the problems. I can then reasonably corpus to units that have $G > X_0$ (i.e. non-solved $> X_0$ problems). X_0 is chosen by $239.07 \left[\frac{\sigma}{G_0 - G} = \text{max} \right]$
At the beginning X_0 will have 2 or 3 ~~prob~~ prob's values ... but based on searching out!! X_0 is much larger.

At first we will have to test each card on all corpus TSO problems, but since never space 245.00

4TH

UMC's ordering simulation strings ...

→ On UMC's: Ordering them by simulation insts.

00 (238.40: spec) must be non-bi-universal, so ~~UMC's~~ UMC's can't simulate UMC's, so there can be no construction of S_{12} or S_{21} that simulates U_1, \dots . Yet we know that there exists S_{21} so $S_{12} S_{21}$ simulates U_1 , so S_{12} can't be defined by either $S_{12} \phi$ or $S_{12} 1$ or \dots must be a sim. inst. for some UMC.

03 So it looks like the idea that 237.00ff is correct!

04 (239.29) Consider "random" min sim. instructions. The model 238.15-17 says we have this SRT (238.02). Say the frequency of "1" in graph is some constant, $c > 0$. If this is so, then by moving down to graph and picking the r.t. hand choice every "1" branch, we will get a min sim. code with \dots a fraction c more 1 's than 0 's. This cannot be compensated for by having more 0 's in general, because user of 0 's is worse so, one could go down to have a peak to left (ϕ) choice each time to get a ~~payoff~~ payoff of ϕ .

One possible way out is to have freq. of $1 = \phi$ - they never occur or occur a den. no.

12 of times. If ~~payoff~~ occurs, then the SRT (238.02) there is a single thresholding, and all min sim. codes are prefixes in this string!

They have to occur in ϕ of times! Case 29-36

Assoc. in each UMC is an associated infinite ("random string") that simulates all min sim. instructions of all other UMC's. [A truly surprising result!] But false

16 If min length sim. insts can occur > 1 way, this can't occur "too frequently", if it occurred frequently enough, we could trace paths that would exceed predicted min length codes in a way so as to make no of 1 's significantly $>$ no. of 0 's.

I haven't proved this ~~isn't~~ convincingly, yet!

18 So: Another possible result is that duplicate min. codes occur infrequently (or never)

If 12 is true, then all UMC's are related to all other UMC's in a simple ordered way. Each has a unique subcode that tells how long its code is relative to other's. Also, there are just \dots UMC's w/ den. long (4n). If there are more, then one or more of the 2 assumptions (12 & 20) can't be true. If there is $<$ an UMC per min den length n, then it looks like serious trouble in the proof (mod 2?)

In fact, I've been considering only one UMC's its sim. insts. - However, in the complete tree of all UMC's ~~simulation~~ simulations, we can pick any unique position to know exactly to UMC at that point will have E as its sim. tree, the ~~same~~ same tree that it is to parent of. So ~~the way~~ the way all the sources for any UMC will contain its sim. trees for all other UMC's. These trees contain UMC's only.

29 If we code repetitions only have min codes, the SRT (238.02) will not have SRT's as substrings, this is because every SRT has to have all UMC's & about any substring of it cannot contain all UMC's.

31 T. sorry, ~~suppose~~ that the original UMC tree must have some branches w. bifurcations. If it did not then it would be a single string, and since each UMC has to be able to simulate itself (see 238.36) - say via $S_{12} S_{21}$ then the entire string would be $(S_{12} S_{21})^{(c)}$, with exactly simulates a finite no. of UMC's.

The tree on all UMC's once cannot be able to simulate itself - they would have duplication

33 So, if original UMC tree must have some bifurcations, but it can't have "too many" of them (see 06-10) (Also note (20)) T. max no. of bifurcations per symbol, as one moves down out to root, must $\rightarrow \phi$.

35 More exactly: Say U is a distance > 1 away from root and $M(N)$ is #. no. of bifurcations per sim. inst. Then $M(N) = 0$, but $\lim_{N \rightarrow \infty} \frac{M(N)}{N} = 0$. The argument of 10 can be made rigorous; if $\frac{M(N)}{N} \rightarrow k > 0$ then there must be a way of coding to string in order compatibility. (This is a proof) $\rightarrow 246.00$

$\rightarrow 265.20 \rightarrow$

246.00 ~~246.05~~!!

4TM

Comment on ID's report, P19: $h_{j,l}^{i,r}(G^{j,r})$ is ^{prob} density. Is this density in Γ ?
G or T directions or both? T is fixed, so density must be in G direction.

Eq. (11) (ibid): if we use $G^{m,l}$ should we have $G^{m,l}$ under Γ ?
or perhaps $h_{m,l}^{i,r}$ is Γ -inappropriate.

Eq. (11) $\dots \rightarrow \int_{-\infty}^{\infty} G h(G) dG$. Looks O.K., its 1. first moment of $h(G)$.

There is no natural analog of Γ or σ for Γ or σ of (Γ, σ) , since the origin at $G=0$ is of no particular significance.

In fact the criterion of selection for $G \in \mathbb{D}$, which is 239.12 - Γ is identical to 10.07! So Γ -min difference betw $\Phi \in \mathbb{D}$ is Γ -method of construction of a set of PSM's $\in \mathbb{D}$ (cards) to be selected from! - which avoids very direct.

$G \in \mathbb{D}$ can be regarded as particular PSM: - or rather, to characterize $X_0 \leftarrow (234.21)$ gives a particular PSM. This PSM consists of multiple bits (S and S).

For S-lang of strings defined by Γ -value of X_0 . As we execute this PSM, we simply ~~move~~ $\uparrow X_0$, which changes Γ -PSM, to a more optimum PSM, in ~~the~~ Γ -direction line with.

$$\left(\frac{\Phi}{G_0 - \bar{G}} = Max \right) \leftarrow (239.07, \approx 241.04 - 10)$$

To "sharpen up" $G \in \mathbb{D}$: The original set of cards $\in \mathbb{D}$ defines Γ -strings, are carefully chosen so that they are very Γ -close to present problem - Γ -criterion of "similarity" being based on comparison comparisons: i.e. we look at many "problem deans" / This set of problem deans that are closest to Γ -present problem. Take Γ -set of cards provided in solving these problems

and use method (via PSM) to derive Γ -string using $X_0 \leftarrow (234.21)$

to get a PSM to solve given problem.

16-20 looks V.G. BVT .17-.19 is not so clear! ^{i.e. I'm not sure it can be meaningfully ~~used~~ derived/realized.}

Another tech (padding better), is that ultimately, we don't want to have to match PSM's to problems: we want a single PSM that is able to ~~match~~ solve any problem by

looking at Γ -problem, deciding what kind of problem it was, and decide on Γ -best approaches for that problem. The techniques I've been describing take over part.

Most part, more "el" in that they develop special techniques for looking at Γ -problem, deciding what kind of problem it is, and applying techniques that have in the past

been successful w. that "kind" of problem.

This looks less el method could be done using a suitable TSP. The TSP is developed so that

Γ -methods of .27-.30 ~~are~~ likely to be discovered, but Γ -search method used is for ^{simple} soln. to all problems (oo + with set of solns to all problems).

View in Disney $G \in \mathbb{D}$, using $\frac{\Phi}{G_0 - \bar{G}}$ is a suitable TSP way to proceed to be a "first system".

Note that this is w/o. 6. Attempted refinement of 16-22

Also Note that "PPM" used in $G \in \mathbb{D}$ is only as good as Γ -~~example~~ ^{could} "solns" it has in its corpus plus Γ -techniques for finding very in that corpus,

4TM

20: 239.40 : Phase II is GA (B) try to solve OZ probs in (apparently) quite different ways: GA (B) has a somewhat narrow ~~selection~~ ^{slang} of PSM's to try. This (set is of PSM's) gradually changes as to more Genc or trials. But GA (B) simply tries to get a success in G per trial as poss. (Using Genc of 239.07, it is better than simply trying for a better G than up to now... so it's better than normal GA. It does have a better "understanding" of option than Normal GA does.

In Phase II, I will have an slang of PSM's, ~~if~~ It is obtained by "factoring" a large set of v. G. PSM's ~~is~~ expressing them as some kind of slang. The assignment of PSM's to probs is a diff. problem... looks about to ~~be~~ done in Ph. II of GA (B).

Many PSM's modify themselves to best deal w. a particular problem; for PSM's of this kind, a "feedback" of a function that decides which PSM to use... is unnecessary.

In GA (B) PSM's do not ordinarily try to adapt themselves to a problem. The "pre-adaptation" occurs in P.D. on PSM's defined by the slang. (The, actually, there is no reason why these PSM's couldn't adapt themselves to a problem, just as in Phase II.)

3) DEF Def "Phase II" means Phase II as described in I.D.S. Report: § 3, 11, 18, 19, 20 and applied to OZ probs
 4) DEF Def "GA (B)" means the GA system of 236.21, 24, presumably w. redimensioned 239.07: ($\frac{G}{G_0 - G}$)

It will be well to describe Ph II in a little detail now: My present descn may be somewhat different from I.D.S. report, in view of my recent analysis of PPM = PPM "like" systems.

We have this set of PSM's [PSM's I], which (presumably) were obtained first as discrete set of 10 or 20 PSM's that I inserted into the system. At a more advanced stage, I factor these PSM's, ~~is~~ express f. set of PSM's as a s. language - so (perhaps) any PSM will be in set. ^{That} ^{if it is to be universal}

Assume we solve only OZ probs. (to be out)

Better we do .175 ft, we have a (finite discrete) set of say 10000 PSM's. A problem comes in, P_{ij}. From our experience w. PSM's & probs, we use a table to see PSM's P_{ij} & T₀ (t. < B for a unique problem) we obtain a curve



obtain Genc G in time T₀. The table G_{ij} (costs & time) gives an index as to how good PSM's is in P_{ij} for time T₀. - I would use G_{ij} or time G_{ij} to order the PSM's for our trials to solve P_{ij}.

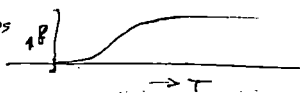
This is all wrong! it applies to INV probs, not OZ probs. Score .31

If we use the "advanced version" (.175 ft), we can create a set of PSM's & select 3. one w. max of f_{ij} G_{ij}.

Same thing wrong here! $\frac{P}{G_0}$ is near or past print! But G. table here from would not apply in this case!

It applies to INV problems in which we have a additive cost function per trial. Does it seem to apply here!

I think it was for inv. probs



P = prob of soln. in time < T.

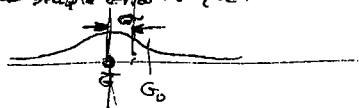
So how did I apply it to OZ problems? In fact, I applied to Gamma House Decm I (I.D.S. P 17 ft) to INV problems not OZ problems. The methods of WOV were used for INV, not OZ problems. I think that on PP 18, 19, 20 of I.D.S. I did OZ probs & used "Expected value" of G as an ordering process for trials. (It's assumed "linearized Genc"), \rightarrow (Score 241.06 ft)

SM → UTILITY of "GAMBLERS RUIN" .20

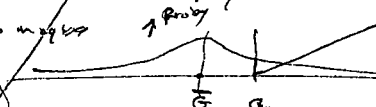
00: (336.40): Another diff betw Ph. II & (b) (236.21, 22): In Ph II I pick PSM carefully, then spend a fair amt. of time applying it, until y. prob's solid on outill that PSM no longer looks "Best".

02: In (b), hvr, I just make successive trials - perhaps rather quickly -
It could well be that they are essentially diff techniques: Tho it may be poss. to apply "Gambler's Ruin" to the Market of (b): To show $\bar{G} + \sigma_G$ order is best poss. order for trials.

05: Hvr. suppose I want trials in order of pc that e. could will be better than Best G. Plus for?
For Gaussian Def's, if I know M_0 of \bar{G} that I want to exceed it's a sample error integral.

07: Or, now ~~sample error~~ t. pc ordering is $\left[\frac{G}{G_0 - G} \right]$ 

There is at least 1 situation in which t. Gov. of .07
may not be so hot: i.e. it may give a very high pc for a very small Δ in G : hvr may prefer a smaller pc of a much larger Δ in G .

11: O.K. Then to "linearization" of G comes into play. We want to maximize expected value of G over \bar{G} . So maybe 

12: P. f. of G is (\bar{G}, σ_G)
we want to max $\int_{G_0}^{\infty} (G - G_0) R(G) dG$. So find \bar{G}, σ, G_0 triplets that do this.

15: T. Gov. of .07 ~~does~~ exactly ~~word~~ word plus Criteria!

18: Also Note: we can use "Steady State" (non-batch) mode: When we ~~test~~ test a new cond: we can tell how \bar{G}, σ compare G_0 over individual - R is will tell us how much (perhaps) to change the X_0 (initials of 234.21 or other param(s) preceding the S-loop Bet process after Guide. 243.20 Cover - Key (1978?)

0: **NB** For the optimal $G_0 - \bar{G}$ S-loop, it would seem that such would not be relevant: Picking trials on Mt. Carlo **NB**

0: **SM** or UTILITY functions; I discovered that probability of "Gamb. Ruin" was an inv. negative Utility in SM strategies: Part 2 (in of Money yield Gov. was lost it was not to alternate Gov.:

One use "helly coils" ≤ 1 to reduce pc of "Gamb. Ruin" in a rather intuitive, Aff. way. How to do it more accurately? Perhaps try linear or ln of Money yield version

PLUS a larger Mag yield for Gamb. Ruin. It may be poss. to design shots in which pc of Gamb. Ruin = 0 (I.E. in response to utility of G . R. = $-\infty$).

So: try shots w. $U(G.R.) = -\infty$ w. infinite Money = Money $\hat{=}$ also w. U of Money = (un)Money).

One Attractive (at first) scheme for $U(G.R.) = -\infty$: Bet $\frac{1}{2}$ of Bank! Invest $\frac{1}{2}$ Bank in fund w. low yield but exp. of $G.R. = 0$. Tooubt is, if one loses all of bet, one ends up with $\frac{B}{2}$ (E Bank) and the process is repeated. again and again. Eventually one can u keep

B so small that one can no longer bet ($\hat{=}$ $G.R.$).

So it looks like $U(G.R.) = -\infty$ is con strat. for 1 bet, it will not be war table for $\hat{=}$ "Long Time Strategy".

2, (02) Phase II vs. (b) (236.21-23). (1) ~~GA~~ GA one uses t. Gov. to derive a TSA for an OZ problem. In GA, we try various Modes of PSM's. There is an Enormous no. of PSM's to try $\hat{=}$ we hit equation on Real Time PD for t. Guide is devised by looking at t. Problem's characteristic properties. it's an answer w. / similar problems of the past:

Phase II has a much smaller set of PSM's, but it is (perhaps) a bit more carefully judge t. applicability of each of them to a particular problem. T. techniques for judging applicability is similar to that of (b) in GA, but loss is perhaps larger (?).

ATH

Simulation Tree of UMC's

Spec: 0:238:0: $P_{i,j}$ is because the length of code is \geq its row. \therefore we always rotate. First row a code occurs in \therefore min length. So only min length ~~code~~ ~~sim~~ ~~ms~~ ~~much~~ ~~occurs~~ ~~in~~ ~~it~~.

1:02 Def(SRT) Super-reduced tree.

(14) For ~~any~~ original reduced tree, is self-similar, in that any subtree rooted at a particular ~~is~~ simulation of U_i , is identical to any other subtree rooted at sim. of U_i .

The super-reduced tree does not have this redundancy, since it consists of min codes only, if may be that this tree is completely random - it has no regularity at all!

(I'm not sure of this randomness, but) - If it is a random tree, then consider the property: As every node \rightarrow prob of \geq (left branch) $\approx 1/2$. This would enable one to determine the density

of UMC's in psu spec!

If .09 is true: does this imply expected length of tree is ∞ ? if not then it's wrong!

T. probability of any particular path is $\frac{1}{2^n}$ (n is length of path) \therefore there are 2^n paths of length n - so we expect ∞ path length.

Unclear: At \geq nodes: \rightarrow paths of \approx prob \dots $1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$

say at the n th length, \therefore prob of \geq states is $P_1(n), P_2(n), P_3(n), P_4(n) \leftarrow$ Heuristics

My guess: That \geq n nodes, \therefore expected no of $\begin{cases} / \\ \backslash \\ \wedge \end{cases}$ is $2^n \cdot \frac{1}{2}$; $2^n \cdot \frac{1}{4}$; $2^n \cdot \frac{1}{8}$

so at n th node expected no of $/$ is $2^{n-1} (\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots)$

I'm still confused about this. T. unreality of ∞ randomness of ∞ previous pages, etc.

A node could have 1 or more children $\approx A$, if one ~~is~~ is not UMC. (Both could be non-UMC but matching is not univ.). One a both children are repetitions of a previous (shorter) code.

from 15-17 T. expected no of nodes at n th level (level is constant, indep of n). - ~~So~~ So for large and n , we ~~should~~ ^{must} have Gambler's Ruin. - so it couldn't be true!

\rightarrow If $\frac{1}{2}$ fraction of $/, \backslash, \wedge$ at each n would \rightarrow some - then there no Gambler's ruin (or \approx prob 0)

\rightarrow But last guess must be true! In fact, the PC of "total end" is $< \frac{1}{2}$, ~~changes~~ ~~to~~ ~~Gambler's~~ ~~ruin~~

PC $\rightarrow 0$. (\therefore PC $\neq 0$ but $\rightarrow 0$ as $n \rightarrow \infty$).

Every Univ. sim code must have at least 1 sim. code child. \dots That child may be (some) ~~is~~ shorter sim code, which case ~~the~~ ~~is~~ SRT(0) that child is aborted. \rightarrow 242.04

Some Q about whether model of 237.00 is correct! T. defn of universality: partial recursive

If $F(x)$ is an arbitrary function, then if U is a universal UMC (or function) \exists a string S_x \rightarrow $U(S_x, X) = F(x)$ for all x for which $F(x)$ is defined. In particular, $F(x)$ may be any other UMC. $U(x)$ has to be defined, hvr. So U is a U/MC machine; so ~~the~~ w. input tapes X .

U eventually produces a finite number string - this is the "output".

Could \rightarrow sim codes be a prefix set? Well, looks impossible. - proof: \rightarrow Say S_{p_2} is a sim code from

U_1 to U_2 , and S_{q_1} is a sim code from U_2 to U_1 , then $S_{12} S_{21}$ is a sim code from U_1 to U_1 .

Both S_{12} & S_{21} are legal sim codes for U_1 , \therefore clearly S_{12} is a prefix of $S_{12} S_{21}$.

Also, if S_{12} is a sim. code for U_1 , ~~then~~ $\rightarrow U_2$, then ~~the~~ $S_{12} 0$ and/or $\rightarrow S_{12} 1$ must be sim. codes from U_1 to some UMC. If not, then $U(S_{12} 0 \dots)$ and $U(S_{12} 1 \dots)$

The Description Tree of (1) Set of all UMC's .off

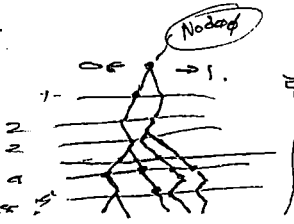
ATM

on UMC's: Say we have a procedure UMC: Consider codes for UMC's:
Prop: For each $U_i(\alpha) = U_0(S_i; \alpha)$, for all α . Consider $\{S_i\}$ such that

U_i are UMC & all UMC have S_j in $\{S_i\}$.

[See 238.32-40, 242.00-03 for proofs of some critical statements below

- 1) every S_i is a prefix of another S_j
2) All prefixes of all S_i are in S .



3) The members of S form a simple tree (list S_i)

But the trees is infinite and therefore no "terminal" pts.

4) One way to enumerate them: Say U is order (row) of nodes. w/ top node row = 0. Start w/ top row, list all nodes. Then list all nodes in 1st row ... 2nd row ... etc row.

Note that $S = \{S_i\}$ is assoc. w/ a particular UMC. Each UMC has its own set, S .

5) For tree $\{S\}$ is (soft similar recursive) in sense that as we move down any branch, we will (usually) eventually come to simulation of U_0 (i.e. doubly). At any such point, we can put Node 0 at that point. (No it is possible to wander down a tree & never come to a point that simulates U_0 .)

6) If ever along a path down a tree we repeat a simulation of any particular UMC, U_j , then we can do the recursive of 5).

7) More generally, in line w/ 5) if the UMC corresponding to a defect pt. out of a branch are identical, then the downward continuations of those pts are identical.

8) At any non-splitting Node, if one takes the choice not in the tree, that alternative will not be a UMC & no continuation of it can be a UMC. So the tree described contains all UMCs and only UMCs.

9) Consider the part of the tree that has no repetitive simulations. We mark each branch where a repeat occurs v. the number of times it repeats, but we don't include that part of branch. If it's not a branch pt. we just stop at that pt. Call this tree w/o repetitions

"The reduced tree". This tree does have terminals. Every UMC occurs once & only once on the "reduced tree" (This last statement may be false! - see 238.31) - But 32 fixes it - see 238.32

10) If S_i is the shortest simulation of U_i (w/ UMC U_0), then must S_i be cut. "reduced tree" will, say have one 2/codes for U_i of some length!

11) (10) supports that "reduced tree" could contain repeats of a UMC: if it repeats over in different branches (i.e. neither is an ancestor of the other).

12) We can take reduced tree of 9) (22 ff) and remove repetitions in a safe way, by taking reduced tree, removing any repetition by an immediate lower nodes. So we take next simulation is "left most" simulation. When we remove a node, we also remove any branches below from it. - we put a suitable marker at each pt. that we remove a node.

13) If we do (12) then only the shortest codes will occur in the super-reduced tree. - see 238.00

4TM

OB-OP. ORTHOGONALIZATION .08

G-PS .05

00: ~~236.40~~ : **SN** ① T. CKT design problems solved by Koza's GP would seem to be as diff. as most problems of trying to find a function that satisfies fixed certain conds. His Gave for ~~the~~ ^{these} ~~problems~~ ~~was~~ a set of conditions! Each of which could usually be expressed in "soft" form, so it had a "gradient" - a "hill slope" to work on. Hvr, instead of ~~looking at~~ ~~f. individual constraints~~ that ~~cond violated~~, & trying to fix them, it just tried a single global optm. - i.e. I think New-Sum. GPS ~~idea~~ would be very helpful in probs of this type. We do have techniques for looking at a sub problem & deciding what actions are likely solve that "sub-problem" - (see 235.17)

08 [It may be best to try to orthogonalize f. set of constraints that f. Global problem solvr must satisfy. TM should have much experience. "Orthogonalization"! A possl. approach to Ortho: we have this "soft" ob-op result Algebra. It is stochastic, based on past observations. We ob sample f. system state: from this we can, for each possl. of action, get a pd other possl. result states. From a matrix of probabilities we can compute ~~an~~ an op that is most likely to advance to Goal (which is some summation of functions on possl. result states.)

② Anytime: I want to understand Koza's ~~the~~ ^{languages} ~~language~~ for Deriving CKTs. If it is likely that such a lang. can be applied to many problems, (like Gabriel Keron's ^{related to km-variables of Keron-Rhodes?} "Tensor Analysis": Koza's Vol III (1999) ^{as 139-140 pp} has much on this. Hvr. He may have incorporated much of the stuff in this book, in more recent papers. Koza's Book IV (2002): 6 coalitions! Grassmann recent results.

19 ~~236.40~~ I also have at least 1 paper (maybe more) extension EvComp by Koza et al on discovering CKTs, opamps, etc.

21 On P.D. for Population Modeling v.s. P.D. for modeling "Populations" as as G. T. GA community is largely to ②. The key claim is a significant improvement over popul using population and cross/mut, it's not clear as to just why it should be better! On the other hand ① looks (in theory) very promising if we can find good ways to do it. 234.19 gives one way to do this (fit it's non-parametric in "G" direction) Another kind of model is 235.25 = .40 (This needs a lot more work!) In both kinds of ① Models (perhaps all kinds), perhaps a good strategy related to evolution for card generation

28 is to choose a card at ⇒ $G(\alpha)$ (given by model) + $G(\beta)$ is max. So of choosing that card has most (chance) of giving a higher G. Choose cards in this order (Max first, of course) How is .21 - .29 related to activity in Phase 2 of "Alpha" (I'd like report)?

29 In Phase 2, we get, for each PSM, problem pair ⇒ T_0 ; a probly first with the PSMs will obtain, iterations applied to Problem for time T_0 , will obtain a Goal of G . In GA language, prob is fixed. (perhaps T_0 is fixed). For each PSM, we have a probly first applied to f. problem, G will be obtained in time T_0 ... for all values of G . So PSMs is a card in v. population. T. problem become part of finding a good PSM. : The prob seems correct, in fact, it is not.

36 There are four PSMs but there are conds in GA. Also Note: The arbitrary Criteria of .28 - .29 is seems similar to the locally optimum ordering scheme I want to use in Phase II - a real solution to a WON problem! So maybe ① (.21) is not so distant from Phase II! (I'd like report) They seem very similar. 237, 238 on UMC Traces (239.00) SRC

10: 233.29 : On Review: T. goal is a simple picture of the project, telling what has to be done & to some extent, what order of execution. Also, possibl/probable ~~differs~~ ~~is~~ ~~my~~ ~~proceedings~~, BoE fixes that I think of now.

How this differs (when ideas) from my previous phases, Phase 2 approach in the IDSA report.

I think the original idea was to simply "improve" conventional G/A/P

T. original idea was to improve G/A/P in 4 ways:

Also use $MSE = 239.07 \times \max \frac{\sigma_G}{\sigma_{max} - \sigma}$

26 (1) Make Model of population: Use ~~the~~ PPM or my version of PPM w. possible for PPM improvement such as sub-tree contexts. This should (at least) speed up GA &/or give it less "parameter conc." problems

28 (2) Allow GA to use info from previous problems by (a) Using modulus (updating) of previous Population P.D. (b) Using a very general Population P.D. & adding sub-trees that seem likely to be useful. [i.e. Don't kill it. Enter after he solves a problem. Use this similar Eng. for new problem]

(3) The system should be able to initialize itself for a new problem (This is really part of 2)

In conventional GA P.D. involves devising initial population (and deciding on size of that population) devising fit/cross rules. Deciding on what Genes to use. $\frac{5}{8}$ Pro Genes is given in the problem data, but it may need finer granularity (e.g. continuous rather than 0. or 1 only!). \rightarrow This extends case of ENV into a 0/1 problem. When we use a P.D. to describe a "Population" we will use the old P.D. & perhaps add definitions as in (2). (0.08)

7 The initialization must be "learned" by the system. It notices relations of problem data to the definitions found most useful in solving that kind of problem. The "kind of problem" can be discovered using the "Compression techniques" as by Li, V. & Yangi... & others. Here, the Compression should be better than just BZZ (GPPM): I'd also want to include sub-tree contexts, & perhaps G/Grammer discovery.

21 (4) Use of Lsearch for proposing new Contexts: Int. system(s) of .06-.21 P.D. merely describe the population. Picking the elements most likely to be in the population does not seem like a particularly good strategy here. A much better approach is to try to model freq G(α_i). (α_i is context) in a way such that it's easy to list α_i 's of likely expected G. One kind of model simply tries to predict G as a function of α_i . An ideal function will have a small error & we may want to minimize something like that. This is a special case of the more general $P(G|\alpha_i)$, in which we want to maximize

Search 2G1.20 to ~255 to estimate & discover Lsearch's G-P

24
$$2p \cdot \prod P(G|\alpha_i) \quad \alpha_i \text{ ranges over known cases. } 2p \text{ is approx of } T \text{ function } P$$

25 If we use the prediction function of 24-28, w. a ms error criterion, G's is equivalent to trying a Gaussian Error df. for $P(G|\alpha_i)$ (of (29)) and using

the maximization of 29 for our "Goodness of fit" criterion"

A possibl. form of $P(G|\alpha_i)$ of $\{G(\alpha_i)\}$ w. perhaps Gaussian error } expresses $G(\alpha_i)$ as equal to the pc output of $\tilde{P}(\alpha_i)$ in some stochastic Grammar (23 fig 5-18) ($\tilde{P}(\alpha_i)$ is the stoch. Grammar).

28 (5) Both GA & Adaptive Lsearch are unable to watch Contexts being tested. They are only told whether

whether a Context passes or fails or whether its error mes, but never "why". (See 263.34)

Newell/Simon's G-P does look at why a new & disc works

(6) ~~When~~ When the Genes have much cc. Try to Approximate it during early phases of GA search.

444

Rev. 00

10:23.40 spec

11 should "live vs pause": A P.D. will be a ways modeling to Population
 (possibly w. some other func of G or $\delta(G)$). On the other hand G or $\delta(G)$ will be trying to
 model some function of t or α , not related to population. However, it is conceivable that t population
 is $\delta(G)$ would come to form a equilibrium eventually, if they were appropriate.

If we added to population viz. Lsry, it would certainly boost population lot!

Another approach: We just use δ functional form of a P.D. to approximate $\delta(G)$ or G .

This functional form is good because it enables us to do Lsry ≈ 0 select cruds of δ easily.

Given a $\delta(G)$ or G one could parse each of several cruds & find its pc w.r.t. P.D.
 Grammar. For continuous phrases one could use NLP to find best values.

However, selecting direct to parse of Grammar is much harder!

One could in theory, find best grammar to fit a set of cruds. To know do such
 a run might be closeness of $\delta(G)$ or G to $P.D.$ as assigned by the Grammar. "closeness"
 could be rms or max value of difference... Essentially, it is a "convexity" problem.

We don't have to use S-grammars to model t . G as we could. We may have a way
 that is slightly expensive to get t models G of a crud, we then rapidly hunt around to
 find a crud w. high expected G , then we can use it as a total it was low
 a from Some that is very expensive.

If we do not try to crud w. models like G , is the true G gives us about the same G .

What do we do to get high G ? We could try to refine around that peak.

One way to get a probabilistic function of G is to find non parametric in P.D.
 "G" direction, however!). Put t cruds in lex order as in BW xpm, but also include
 empirical G at each point. By only using data w/ $G > X_0$, we can get
 a P.D. on future cruds based only on past data $G > X_0$.

So we can be as elitist as we like, by choosing high values for X_0 .

The bad thing about this method is the "Non parametric" in the G direction. This is very wasteful
 of info in limits to goodness of our models. It may have, to be good enough for Phase 1.

To get beyond simple Phase 1, we give Phase II: problem of finding good (better)
 Probabilistic Models for t corpus (then our current (augmented) PPM).

I don't yet have any good ideas for models to try (other than this "Augmented PPM"). However, it
 is a well defined or problem, so it should be able to devise a TSQ of problems leading
 to it.

This stuff from (.00) to .30 is important! (.17-18) is an implicit criticism of one of the ideas I had

That was quite good! i.e. the modeling of G as a func of t cruds α_i . That as well as
 233.00 - .17 which is an using of probabilistic Grammars to model G as a func of α_i , (t crud done).

Needs to be carefully revisited. What the problem was: to propose some of objections to know.

For each value of X_0 we have a different P.D. on α_i 's & a different distribution of empirical G . This G is a
 func of X_0 & G is a func of X_0 ; we can use that to decide on a X_0 to use to get good prob
 of much higher G in new crud, perhaps choose $X_0 \Rightarrow \bar{G} + \frac{\sigma_G}{G}$ is max — $\int \max \left(\frac{\sigma_G}{G} \right)$ is better: see 239.06

Given several Grammars, t Goodness
 of G has hypothesis!
 239.06
 Given several Grammars, t Goodness
 of G has hypothesis!
 239.06
 Given several Grammars, t Goodness
 of G has hypothesis!
 239.06
 Given several Grammars, t Goodness
 of G has hypothesis!
 239.06

This is our way of dealing w. the problem. See 239.06 for a better way to do this!

37

(.23)

20: : **SN** On Grammar Models for GA: I think that recent work on this tries to model the Population as a S-language. My "improvement" (at least) was to model the Score as an S-language: We ~~are~~ simultaneously have to find a monotonic function between the PC ~~are~~ assigned to a card by the S lang, i.e. Score. I.e. say α_i is a card. $G(\alpha_i)$ is the fitness of α_i . $L(\alpha_i)$ is the PC of α_i as given by S lang, L.

26: Then we want to find $\delta(L) \approx \delta(L(\alpha_i)) \approx G(\alpha_i)$. The data set $[\alpha_i, G(\alpha_i)]$ is used to help find good δ and L.

28: ^{positive} Hvr, it is conceivable that for certain G-sets, 06 is impossible. e.g. say 4 cards and 20 integers. $G(\alpha_i) = 0$ if α_i is even; $G(\alpha_i) = 1$ if α_i is odd. } Not sure... 11-17

0: Since there is 20 of even/odd integers, the PC of each card, would have to be ϕ . If we had only a finite no of cards possible, we would say N integers; we assign $\frac{1}{N}$ to even integers & $\frac{1}{N}$ to odd integers. Then $\delta(L) = N \cdot \frac{1}{N} = 1$ if α_i is odd.

11: In an actual GA, this problem would prove no difficulty: the SSS is always finite, so we can always normalize. As SSS \uparrow to normalize constant changes, we do this with large English corpora, as new words enter corpus. So I think we should have no difficulty choosing L $\approx \delta$; δ is normalizing & translatory function, will vary to compensate for \uparrow M SSS.

7: **NB** Any range of $G(-\infty, \infty | 0, \infty | \text{normal} | 7, 29 | \text{can be easily mapped onto } [0, 1])$

8, 23, 40: Normally in GA the form of Score (written $\delta(G)$) can be regarded as "Extreme Cardinal Code".

0: Note on "MASTER Plan": The idea of the Plan is mainly that I should not get lost. I think I did get lost after Sol 99 & P. SAAR work. (?)

The TSC should not ever be a Bottleneck! I should be able to write problems in so usual simple domains in which you could ~~effectively~~ effectively take advantage of a TSC — or something like a TSC.

I do also want to review for (1) Myself... to keep me aware of what I've already done. (2) To have usable review docs in case of fire or other calamity — also so I can easily take with me, usable, usable background for further work, when traveling.

One poss. policy: Actually spend at least 1 hr/day or review: I spend 1 hr (min) on project itself. Rest of time on ~~any thing that comes up~~ (rev, any project or house stuff). - 235:00

29: **Q**: Is the $\delta(G)$ x form ~~any~~ related in any way to "linearizing G"? Can a linearized G go from $-\infty$ to $+\infty$? or 0 to $+\infty$? Consider $G = \text{Money or Time or cc}$. To deal w. neg utilities: balance them w positive utilities. So if B is + value of a neg (Bad) [say] G to value of a good thing. $B_1 \cdot G_1 = \phi = B_2 \cdot G_2$. Zero has property that any amount of it is G-saver. So $B_1 \cdot G_1 = 2B_1 \cdot G_1$. So, super-fittably, utilities may be mapped well onto probabilities $\frac{1}{2}$ — fit into domain(?) $[0, 1]$.

A serious "second look" on GA in general & my Phase I in general. In GA, modeling the population as P.D. is just plain: it models population, not really Score. In my "Phase I" & P.D. models pass successful cards, Hvr, I could use a P.D. to model $\delta(G(\alpha_i))$.

ATM

- 1) IMP sub. projects: 1) Speedup 2-3 Trappem sort works on Cons Computer.
 - 2) (Perhaps) modify it so I can delete selected items.
 - 3) Get 2 fast compression programs: Tixit on Eng text: compare w normal PPM
 - 4) Get 2 GP PPM running. (Notes 231.36) - using Lisp/Perl/fortran/assembly.
- Read stuff on use of Assembly lang in GP.

SN Look at detailed structure of Fortran Assembly pems. What kinds of regys do they have that are common to related pems? How can these ^{kinds of} regys be detected? Can PPM be ~~used~~ adapted to do it? Would ~~Common~~ ~~Sub-tree~~ discovery be useful? Would use of deltas (say in Fort) be useful?

SN Werners Monte Carlo pen scores by randomly choosing previous tokens; was able to generate desirable biased Mt-Carlo pens by randomly choosing sub-trees from past. Had to do do Mt-Carlo Lsuck, which is normally expensive. Would have to keep pens in a lexical order or use some kind of Hashing to strengthen. But, ~~generally~~ my Mt-Carlo Lsuck did require a lot of RAM (or slower disc Access). Perhaps I could put a normally Mt-Carlo pen Generator into a Lsuck order!

[I don't see how! ...]

A trick! If a cond doesn't have a reported "call" within a certain # trials, we drop the cond from ~~Memory~~.

This involves a daily garbage pc cond.

SN Two ways to speedup (i.e. effects of) sub-tree search in a PPM. ① Make Definitions ② Do $\approx .14 \rightarrow .21$

GA Tournaments: Pick n^2 candy at random from population. Take a Best of these n^2 ; remove $n(n-1)$ of them. Make $n(n-1)$ children by picking random nodes in each of n "top" and exchanging branches. Replace $n(n-1)$ discarded candy by these children.

The degree of "elitism" \uparrow as $n \uparrow$. It is most extreme, when $n^2 \approx$ population size! Say one ~~cond~~ $n^2 =$ pop size! We continue this until it appears that the "top k " are not changing.



One problem is that in "top k ", some of branches from different candy will be the same. So we must have to check this out - if they are the same & choose another.

When it becomes that we have "converged", do an exhaustive generation of children of "top k ".
 How, effect of n has k branches, that's $n^k(k(n-1))$ poss. children (if not branches are the same).
 n^2 is pretty a doable number.

Or, perhaps each time one chooses n^2 for tournaments do all $n^k(n-1)$ children and keep $n(n-1)$ best of them.

NB Tournaments use, whereas GA much less scores remain for any δ (Gen) ~~known~~!

Can be used to "standardize" score so it is more continuous, but nothing else!
 E.g. f. 2 valued score of a CNV problem can be adm "interpreted" for partial credit...
 In this case $f(\delta)$ is critical. If does with an CNV into a 0/2 problem. \rightarrow (233.17)

10

ADMIN: Perhaps Make outline of TDM project with time goals/Deadlines for each section. As I move thru the project, I may be able intelligently modify various deadlines in view of my experiences.

My Guess is that it would be best to first, make detailed outline w. refs. to how to do each parts & w/ alternative ways to improve &/o overcome expected BUGS.

Next, Good to get a "more or less working model" of final system using minimally functional techniques for each section.

Next, Fix various parts - doing the most critical first.

This way I get F.B. (Both Intellectual & Political) as soon as possible.

Deficys of GA: ① slow, & (occasional/local extrema) ② Doesn't remember previous problems. ③ Can't do ts Q's ordinarily ④ Population as Model of P.D. ... Other Models can be better. ④ L such rules much better than M. Code Such.

By just doing ordinary GA problems, using PPM (perhaps Heuristic) & L such, I should be able to do much better than GA & perhaps be able to do ts Q's by starting in the B-N X form (to BZZ)

SN

The system I have in mind for induction is oriented toward finding good functions for QA problems. Whether such a system would be as good for the kinds of problems that have been addressed by GP is unclear.

I did an analysis of L such, vis. GP for $n=$ Multiplexer. T. system was able to learn $(1, 2, 3, \dots)$ & eq. w. not much extra cc as in P. But I had n't tried for it.

recursive soly. - that TDM should have been able to induce from the successive non-recursive soly. Q: Does the present TDM model Automatically normally look at the seq. of previous QA solus & try to find a (possibly recursive) way to solve "the next QA"?

Each of $\frac{1}{2}(n_1-2)^!$ is rather expensive: I will try various simplifns, approxs. $\frac{1}{2}(\frac{1}{2}(n_1-2)^!)$

I have to sum over many of these, But in fact, if I start with the longest context of $5 \leq z \leq 1$, & work toward shorter contexts... perhaps I will quickly come to a peak, & not don't have to go much further! A big Q is "How much cc is spent, on the average, on this particular" (29)

Function Evaln?

SN HWSW Note: If I want to use Lisp: I can realize Lisp in Fort w. perhaps much less speed: Further if I use any partly compiled code (in which the no. of "New" statements is reduced by expressing very short sections of short codes "in line". A number poss'ly to express APL in Fort. [See Advances in Gen. Progr. 1994 Kim now ed. pp 295 - 324 (2 papers) for Code Editors I have this book. Other Good Papers in this book]

29
0

36

4:TM

Spec

00: 22.37 → While the Z141 method is ok for normal sequential predn, I don't see how it can be applied to tree contexts. In Z141 I could do define / sequence / corpus using PPM ngm. In t. case of tree contexts I can see how to define T2m, but I don't see how to code corpus using PPM → well: try P2m: If t. corpus is a set of trees, we can use a particular sub tree P2m as part of "many of these trees, & as a function w. (often) several log elements, i. deriv. i. corpus under w. that function as one of t. "primitives" — I think that would do it → (37)

08. **SN** On using ideas from old pps, in new pps. (TSAQ's): One way to do P2m for a general problem solver (like GA, say) is to transfer a set of definitions used in old pps, to a new program. Here we look at our new problem but deriv; we decide which problems in the past it is "most ~ to". We then transfer the defus used in those problems to the new

System that is to solve the new problem.

- ~~Other poss. ways to transfer into from one solved probl to new unsolved problem B.~~
- 1) Defus. used in A (.08) &
 - 2) rules & subrules used in A (as in PPM is "tree PPM")
 - 3) General freq. of use of inst (as in OOPS "boosty" construction) or pairs, triplets, etc. elements.

So: A (perhaps) better way to solve TSAQ's, & analog. I've been thinking of this. Instead of just adding a new problem into the corpus: Analyse it & try to find similar old solved problems: P2m use definitions ("tokens") w. assoc frequencies, that were used in solving these problems. If definitions are allowed, transferring freq. use of insts, includes transferring / wtd. defus.

The "UNDERSTANDING" can be of any size, depth, degree, detail, fuzziness, etc. — and of ways to "understand" ...

A (sort of) objection to .16 ff: that we have to worry while having to write always definitions, materially <JS> of solns of new problems, it's not quite the same as "understanding" why the old defn (function) works, under just what cond. it should be applied. (31)

24. **SN** Designing a cor w. certain specs in as short & low as poss. is a Hutter-type OZ problem. Yet it looks different from Hutter's account of Multiplying Matrices. Actually, .24 looks more like a ENV problem, similarly, perhaps Mult Matrices is best regarded as a kind of ENV problem. Hutter's concern w. a special kind of problem, much more general variety, would be longer as said, generally. Just as Juergen took a long time to finally understand what Lsrah was, Marcus (perhaps Juergen) don't yet know what an OZ problem is & how to solve them by Lsrah. The Jew do know how to "sort of" solve OZ problems, it's not clear that he understood why he was by no means an optimum soln.

07. One thing I'd like TM to be able to do: It sees a problem solved by "invertible substitution". eg an integration problem. It is able to sense the technique used to, say solve of quadratic, & cubic equs. (possibly quartic). More generally, it could try to get it to learn (or "improve") to have: "try to sense the soln. technique of any problem you solve or see anyone else solve."

37: (07) Also relevant: It may be good to use both derivatives of functions and PPM (a. P2m via a w. subtree contexts) Doing an just regular PPM would be very fast & perhaps V.G.! Defus would involve new symbols — which PPM is able to deal with!

4.TM

spca

x: 2.28.40: Li-vitzki's use of compression along (just BZZ) would seem to work \approx 2 bit in TM problem,

0) \approx its Very cheap. Cleverer Compressors, using subtrees, should do even better. (20)

2.28.41 ff Gives a slightly different view of t.TM project! At the beginning, I can simply try to improve GK, in all of the 6 aspects. Later I can concentrate on (2) (TSC's).

05 (5N) Note that approximating trying to model t / original G_{ore} by a P.D. does of necessity involve model of t . G_{ore} by a monotonic function. This is in addition to extra penalty to G_{ore} or otherwise mobility \rightarrow so (a) if by hyper "resala" (\rightarrow 2 levels) (b) is more amenable to t . GK idea.

In recent work on \rightarrow improvement on GK, one approach was representing the population by a P.D. This seems different from approximating t . G_{ore} by a P.D. (as (5S)) — but it may be that it is not. Perhaps the difference is: In the "recent work" \rightarrow use of non-fuzzy population, and the attempt to make it fuzzy (probabilistic). In my formulation, I try

to model t / entire population with its assigned core values, by a P.D. whose value = $\delta(G_{ore})$ where δ is some monotonic function. This uses t data better \approx is likely to give better results.

0.01: Another thing was from Prosa Compression Algorithms: They tell us how to make new PSM's also initializing. The problem deals \approx t. RSM's / initializing from a given member of β , we find a set of α

members for which $\text{Pr}(\beta)$ would work well. We choose β as a subset of α as a corpus to compress. Any new problem α can be added to this subset & we can see how much extra info is needed to compress a augmented subset. Maximum Pr of this subset tells how close our new problem is to $\text{Pr}(\alpha)$ class (in α) of old problems for which Pr element worked.

By looking at other members of β & sets of prob & cons for which they have worked, we can well find members of β that has the best list members of β in order of how likely they are to work α . Also, we can use t. compression Algorithms to generate new β members that have high likelihood of solving α . 2.0 ff \rightarrow needs to be worked out more carefully, and integrated w. my idea of factoring a good set of PSM's / initializing subset of t.TM. how to correlate PSM's to problems, initializing as well as generate new (PSM's / initializing) for new problems.

statistical HT for long Q to A mapping, mapping.

I want to write "sort of paper" for a sequence of projects concerning in a V.G. General problem Solver. Possibly get sponsorship from IDISA or (R. Holway?) or ... get job in staff. Another posy would be to start co. & get Van der Capital" (David S. J. Kotiz, Seth Prudhuan) maybe not good to get "stuck" in USA!

4 T01

So! a Bit of JOY!

20:227.40 : ~~227.00 - .40~~ ; (228.10 - .24 in particular) seems to clear up a bit all of my uncertainty of how to assign pc's in PPM's Tree contexts; Also it suggests that Z(4) is probably basically better.

A poss. advantage of PPM is the Burrows Wheeler method (BZZ) - that it could be much faster than Z(4). If this is true, I may want to trade speed for ~~pc's~~ pc's - possibly for tree contexts: - but look into that later!

Audier form is 226.40 - .23 which ~~doesn't use~~ sticks closer to Z(4) & doesn't exactly use "Context".

But do write a detailed ~~in~~ disc of 227.00 - .40 : As is, it's kind of vague to a person reading it (yr from now (or maybe Mar 04)), 2.20.

Much of PPM is involved w. deciding what to do when a "new symbol" occurs. In Z(4), I didn't address this problem at all! In match (or ≥ 1) of what I've ~~thought~~ thought that about this does not occur (except, perhaps in Laplace's rule). T. PPM analysis

Suggests how to use LZPS to deal w. continually occurring "new symbols" - i.e. At all times, whenever has a corpus of n symbols plus for, there is a pc of $\frac{1}{n}$ of a new symbol. T. total pc of this feature is $\sim (1-\frac{1}{2})(1-\frac{1}{3})(1-\frac{1}{4}) \dots (1-\frac{1}{n}) = (\frac{1}{2} \cdot \frac{2}{3} \cdot \frac{3}{4} \dots \frac{n-1}{n}) = \frac{1}{n}$.

But it occurs for alternative codes, so it sort of cancels out... The codes w. fewer ~~pc's~~ pc's

→ Symbols for the whole corpus would get more wt.

10

Understanding the relationship of "classical" PPM & Z(4) is important because I want to be able to map the "improvements" in classical PPM into Z(4).

Many of the features of PPM are approximations/averages over Z(4) values. These approximations can save lots of cc & ~~even~~ even tho the pc values of PPM may be worse than those of Z(4), the cc saved may make PPM superior! (Hr. for Expensive Fitness functions, the better pc's of Z(4) & of "subtree finding" will prevail.)

SN

One of the difficulties in adapting GA to a problem is modification of Fitness function: In the hard of the problem some fitness function with a large problem. In INU problems, it is usually a MPSE function in which GA can be used, but often almost always, a better Gove (R-function) must be devised. For OZ probs, the Gove must also be modified. If worse

≥ 2 ~~is~~ is ≤ -1 ng. to model the Gove, we have to modify the Gove by some monotonic function.

If the original Gove is very expensive, we want to replace it by a much less expensive approx., at the beginning of the "hill climb", & possibly other modifications during the climb.

When we get to "phase 2" presumably, modification of Gove is unnecessary - But I'll probably want to look into this in some detail.

which is really part of (3) initialization

So: GA, GP has simple difficulties: We add a new one in 26 ft: 5 difficulties:

- 1) GA is slow
- 2) GA does not remember & use ideas/solns of previous problems (No ISQ ability)
- 3) T. PPM initialization of the RUN is Big Task. Getting initial population, Lang. to decide - rules, mut/cross rules.
- 4) To initialization we add Modif. of the Gove (Fitness function) 26.
- 5) It should use LZPS rather than MPSE.

In all of these initialized problems, we "look at the problem down" & propose initialization. → 51cc 229.00

34

ATM

20: 225.40 : T. Q. 4: f. Legitimacy of the codes for PPM of 223 32-22470 (+ 224.21-225.70)
21: codes for Trees: of 222.16-19, 20-24, 25-29 + 223.24.

Any alg. that gives pd for each "next token" in corpus, is a legit code & a legit PPM.
At present time: I think I understand why both PPM code & tree codes are legit:

T. PPM code is legit because: At each ~~state~~ pt. of Token prodn, one has several contexts to base on. We can give each marked wt. A reasonable wt. is to effectiveness in coding f. corpus of that context (is context). Since we only have k different contexts to choose from, we can legitimize wt. PPM w. some function (is a prior) on f. nos. I think k

— an optimized form of this function on the past of the given sequence.
The form can be anything we find empirically "good": Hvr. since we have @ limited we can use

is $A \beta^{-k}$ in which β is f. const. of alphabet & A is a normal constant

12 If $A \beta^{-k}$ is not so bad, this can correspond to f. pc of defining each context.
13 More exactly, this process should be $\propto \prod_{i=1}^k f_i$ where f_i is f. frequency in corpus (is now)
14 of f. symbol in f. context. We get this by simple Z 141 type coding.

15 On the other hand, every time we get a result by considering that a context
16 must be (chosen) by the symbol "s", then will have an "escape cost" (in pc) of f. frequency of "s" in
17 f. corpus plus for.

18 Hvr. the argt. of is seems backwards: Because in normal PPM, they start out in f.
19 longest poss. context, then escape down to one context that can predict the next symbol.
20 (i.e. its pc for next symbol is > 0). On the other hand in normal PPM, they are burdened

21 w. f. additional pc of "escape": This is always an additional symbol in alphabet. Whoops! Note .38
22 Also, PPM doesn't multiply by f. prob of α & $\bar{\alpha}$ — so perhaps normal PPM should be prob of α
23 & rather with escape of which context to use. (That f. prob of α & $\bar{\alpha}$ should be very small)

24 So .12 - .14 which is Z 141 may be kb is perhaps slightly better than PPM.

This appears to be a MAJOR advantage of Z 141 over PPM

Note 226. 16-23 on using Z 141.
So: Z 141 would seem to give a good rigorous model — a way of understanding how to improve PPM.

It does seem to give a very reasonable way to do my modif. of PPM.
It gets rid of f. (some) confusing idea of coding by means of positions — since in that case, the
25 entire (ll) codes of f. corpus are not shared. My most recent (10-24) analysis makes it poss.
26 to get is is probably much better than PPM using the normal ll codes & prob of α 's of
deductions, of Z 141.

33 Also, it clarifies f. assignment of pc's to tree contexts. We effectively have a
pc of deriving each tree used for prodn. T. Now will give no. of symbols, between and
34 β^{-k} (β = no. of symbols) is an approx. $\prod_{i=1}^k f_i$ (in .13) which we use as

35 a suitable normal factor to get a group of f. various ~~tree~~ sub-trees used as "context". → Hefty in trees, given
36 In normal Z 141 we also have to add a new symbol to f. alphabet — usually of
is a normal assoc pc for it. This does correspond to f. pc of "escape"! (When it does so quantitatively is unclear.)

Solns of PPM; 01 || Tree Coding; 02 || AZ41 (10-22)

9/15/09
4PM

- 0: 225.40 : In view of 225.35 - 40: Perhaps drop Ritz problems re turn to Tree Compression coding.
- 01 T. Soln "pro" form of PPM is 22332-224020: (with objection on 224.21-225.40).
- 02 T. Soln to Tree Coding w. PPM: I think ~~or~~ ^{or} ~~Bill~~ ^{Bill} outlining is on 222.25-.36: ^{pass!} ~~is~~.
i.e. ~~(220.16-19)~~ : 220.30-40) ~~222.28-29~~ ^{Man discuss} ~~221.28-29~~ ^{How legit April of Nov.} (Diff of 221.30-35 is resolved on 223.24)

SN IMPT!

I have been using Contexts to do predn. This is in spirit of MDL, but is not much good for OSL (≡ One Shot Long). In PPM (is problem in Tree version of PPM), we can easily modify our techniques by simply considering alternative contexts of f. corpus, & each computing f.p.c.'s of these alternative contexts of f. corpus. Here, we do not consider contexts as having prob. distribns on f. following tokens. The calcns may be slightly simpler than in the usual PPM way.
Hvr, the Mechanics of 0133 ~~is~~ — how we do predns — is not yet clear!
One way: ~~the~~ ^{the} ~~advantage~~ ^{advantage} of longer k includes last k-1 symbols of known corpus plus 1 symbol to follow.
We expect f. corpus 222 0000 No!
I think I would have to use the matrix of AZ41 in my ~~is~~ ^{is} ~~bluff~~ ^{bluff} letter.

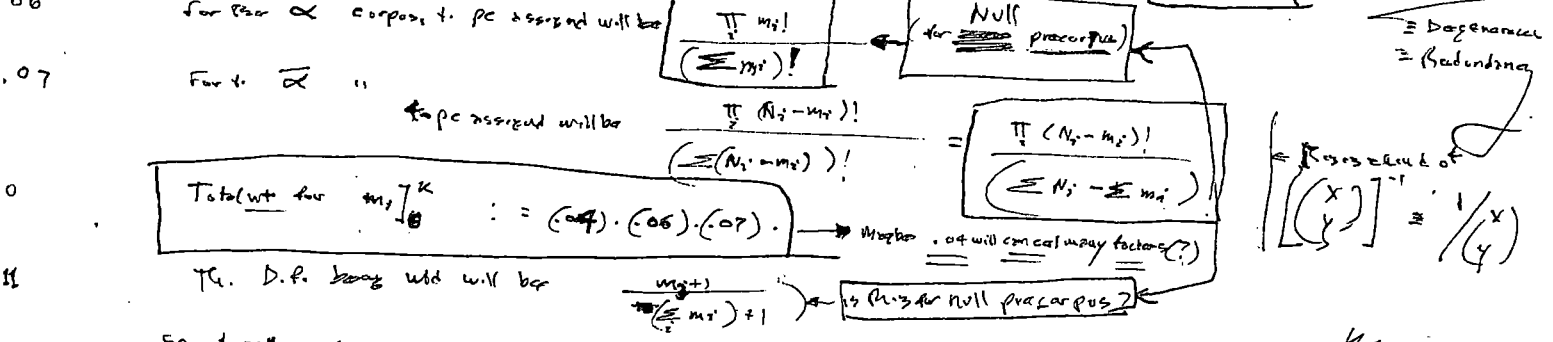
500
227.00 -
227.40
Estimate for 2
continued discussion
of 2141 v.s. PPM.
is their close relationship

Perhaps use it ^{Both} and ^{Big} to previous version PPM for comparison.
~~the~~ ^{the} ~~advantage~~ ^{advantage} of AZ40 method is that it enables (eventually) many deductions & repairs of corpus & definitions of updates → PSG discy.
One poss. advantage over other methods is that we just use BW km to find good news as in usual PPM, — then common substrs are then discovered (probably) by PSG discy techniques — that don't ordinarily involve direct search for substrs. This may be much less expensive.
SM Consider ~~the~~ ^{the} ~~pc of corpus~~ ^{pc of corpus} if context of length k was used every time (as a first of k). This should give to P(k) — not apr. wr. for contexts of length k. We have to be sure that each context gives p > 0 for every symbol (via Laplace or Dirich).

4TM

Perhaps express it as a sum! say we have a specific set of m_i $\sum_{i=1}^k m_i \in N_i, i=1/k$.
 get the redundancy (\equiv "degree of freedom") a'wt.: Permutation. Maybe simplifiable by factoring, by
 various identities, approxs, etc.

for a given m_i, h_i , there are $\frac{n_i!}{m_i!(n_i-m_i)!}$ ways to choose m_i subjects from a set of n_i . $\equiv \binom{n_i}{m_i}$
 Say N_i are the code counts for the entire corpus. So for entire set of symbols: $\prod_{i=1}^k \binom{n_i}{m_i}$ (\equiv no. of possible ways)



So, the red pc of symbol j will be $\frac{N_j}{\sum_{i=1}^k N_i}$

$\frac{N_1}{m_1+1} \cdot \frac{N_2}{m_2+1} \cdot \dots \cdot \frac{N_k}{m_k+1} \left(\frac{\prod m_i!}{(\sum m_i)!} \cdot \frac{\prod (N_i - m_i)!}{(\sum (N_i - m_i))!} \cdot \frac{m_i+1}{(\sum m_i) + 1} \right)$

$\frac{N_j}{m_j+1} \cdot \frac{\prod m_i!}{(\sum m_i)!} \cdot \frac{\prod (N_i - m_i)!}{(\sum (N_i - m_i))!} \cdot \frac{m_i+1}{(\sum m_i) + 1}$

$\frac{N_j}{m_j+1} \cdot \frac{1}{\binom{\sum m_i}{m_j}} \cdot \frac{1}{\binom{\sum (N_i - m_i)}{N_j - m_j}} \cdot \frac{1}{\binom{\sum m_i + 1}{\sum m_i}}$

The $\bar{\alpha}$ factors will be simplified directly by simply considering the no. of edges it has to add
 \equiv using mean entropy for the entire corpus: mult. by fractions of counts $\bar{\alpha}$ has to add.

Other factors will be approximated by Gaussian, Binomial, or Gamma distribus,
 (.04) \equiv wt. j distributed $\bar{\alpha}$ could be a 1 dim Gaussian
 .06 & .07 could be multidim Gauss or Binomial if $\bar{\alpha}$

Other idea - it may be best to neglect Process (1 codes! (sur, I really don't know if Play
 can be neglected w. impunity!

It may be best to drop that "term" for a while! Just see if works well w.o. it in
 the PRTY. The real improvements pre-22 probably more worth spending time on it than
 to search for common subtrees. Have found some useful contact subtrees, use
 and also try them as cands for bothers (to be predicted).

The total wt part of .14 is $\sum_{i=1}^k \frac{n_i!}{m_i!(n_i-m_i)!}$ and this should sum to $2^{\sum n_i}$ we
some sum since because we divide (14 by 2 (some sum) because that's the pc of each of k terms.
 "same sum" $\equiv \sum_{j=1}^k n_j$

$\sum_{j=1}^k \frac{n_j!}{m_j!(n_j-m_j)!} = (1+1)^{n_j}$ so I guess .25 is the product of factors $\left(\frac{1}{2} \right)^k = \frac{1}{2^k} = \frac{1}{2^k} \sum_{j=1}^k n_j$

effective wt. of .14 is one code (It's the wt. of $(1+1)^{\sum n_i}$ codes), $\bar{\alpha}$.
 Q15, How does it compare in other codes - (Having an exp wt. of n^2)?
 We can use the approx method 18-22 to get estimates.
 Another way to try to answer .25: T. codes of .14 are all codes of subsets of α . If we use Play
always get << pc MCM single code for all

47M

10: [SN] Just how does the "precorpus" enter into the coding of a sub corpus? Is it "precorpus" shared by α & $\bar{\alpha}$? if so, how is this done?

View ①: That α & $\bar{\alpha}$ are regarded as separate corpi: each has its own precorpus.

View ② (Unlikely) that the whole $\bar{\alpha}$ & α corpus has this precorpus — so α & $\bar{\alpha}$ interact somehow in using it.

Since α can be a rather small subcorpus, its precorpus can be a big deal. We may want it to only have $\frac{1}{2}$ a symbol for each of the alphabet (which is a common Dirichlet distribn.)

07: 223.40 So: we code the "entire corpus to read" by assigning its symbols to either α or $\bar{\alpha}$, and we obtain a pc for it: $P_{\alpha} \cdot P_{\bar{\alpha}}$. Say this context is of length l . Then we have $P_{\alpha}^l \cdot P_{\bar{\alpha}}^l$ as the pc of α corpus wrt. $\alpha, \bar{\alpha}$: The α tokens to follow, will be those of α^l .

Each of the contexts of $(l=1 \dots n)$ will have a wtd. $P_{\alpha}^l \cdot P_{\bar{\alpha}}^l$ assoc. w. its D.P. on the next token. The final probn. will be a wtd. mean of all of the D.P.s assoc. w. all of the contexts. Then wtd. will be $P(\alpha) \cdot P_{\alpha}^l \cdot P_{\bar{\alpha}}^l$. $P(\alpha)$ is the prior wtd. assigned to contexts of length l .

13 $P(l)$ = ratio of various forms used to assign prior to integers (see 220.30-39) — but ultimately the format $P(l)$ will be empirically adjusted for Mark up past, the pc of the corpus had been

15 maximized (see 221.28-29) — i.e. for f.n. of $(221.24-27)R$ — my bias is to start with α at $l=0 \dots$ rather than pick the largest l w. α or $\bar{\alpha}$ to start with — but I think this has to be empirically determined — which assignment method gives the highest pc to past corpi. (Note 21 on $P(l)$ has also started to do in, t. corpus to date) } see 227.18-24 for an insightful discussion of this Q.!

18 So: 223.31 — 224.18 is a complete down of how to do PPM — except for rescalings of 224.00-06 and 224.15-18: also the dtd of (21) } Of course this is to Q of how best to realize or approximate the technique described.

21 [SN] Note $P(l)$ of (13) will have a function of n (the length of corpus to date): If n is small, then larger values of l can't have larger $P(l)$. Perhaps assume functional form: $P(l, n) = P_1(l) \cdot P_2(n)$?

[SN] While 223.31-224.18 may be a good Bayesian approach, the recent "improvements" in PPM may be expressible in Bayesian form & be integrated into my approach — so I can optimally incorporate the recent improvements into my model.

[SN] A theoretical criticism of my l & large PPM approach is we can give lot more

11 codes by treating some of the α cases as $\bar{\alpha}$ cases. This was done in my "Letter to Gerry Wolff, on 1/14/04." See below Φ PTM 2141 [2000] for hardcopy letters + correction.

One Q is whether we have to say which pairs of the corpus we're using. If α & $\bar{\alpha}$ are parallel passages then it's easy, which one we're using has $pc = \frac{1}{2}$. This seems correct: So, given n case counts & \bar{n} precorpus case counts, one should be able to get a distribution for the next token and so on.

Considering any subset of α as a $\bar{\alpha}$ subcorpus: code. If α has k symbols with each being one just 2^k different sub corpi (including null sub corpi) each sub-corpus β has its assoc $\bar{\beta}$, so code of entire corpus.

One may considerations of precorpus: if $n_2 = \sum_{i=1}^k n_i$ are the case counts of k symbols in α . Then the subcorpi will have maximum 0 to n_2 case counts of symbols (see 2).

00: 222.00 : T , pointers, t. Temp. & desc. is 2 left (PEM)
 Now ($\equiv P_{Nov}$) to assign pc's to sets of PEMs. (Each P_{Nov} has a different set of PEMs). ~~Each Assoc. w.~~
 02 each P_{Nov} is a "single current PEM", that assigns pc's to "next token" on basis of previous known corpus.

03 The "PEM" of .02 consists of a bunch of PEMs in \cup . They are to standard α , $\bar{\alpha}$ pairs assoc. w. α & $\bar{\alpha}$ rules
 04 (or Dirichlet's rule). T sub. PEMs of .03 can be anything I like, I could assign them to every poss. context, including
 05 - contexts that do not occur at that point) - Hvr, I chose to not include those contexts. If I did
 include them, it would seem to (be close) to not consider context... but this is not exactly true; only those
 27 ~~each context has a α , $\bar{\alpha}$ and a pc for corpus~~ and ~~next~~ ~~occurs~~.

But I think we want to only consider pc's assoc. w. contexts that did just occur. - If we did not, we would
 be summing over ~~almost~~ ~~the~~ set of PC's - indep. of t . context context.

10 Now say we do .03 & include just PEMs assoc. w. contexts that just occurred & we use recursive "uts (.07)
~~this~~ we then ~~sum~~ sum over all P_{Nov} (sum since total P_{Nov} sum \cup). This is a legit. P. & M, since
 12 it assigns an (α already normalized) pc to each pos. ("next token"). - Th. Q is $\frac{1}{2}$ it is by good?

13 To better legitimize us to ditch $(.00 - .12)$ express it as \cup codes of t corpus.
 - well, any PEM has an assoc. code, so ~~the~~ ~~code~~ ~~is~~ ~~legit~~ ~~to~~ ~~be~~ ~~the~~ ~~extent~~.
 5 $\rightarrow T$. Q: α is next to codes, dec since? α Also there is (transposable / good) approx. to Dirichlet diff.?

Well according to .13 - .15 we would then want to include all of \cup contexts (most of which did not just occur).
 - Hvr, all those "read out" pairs are (identical) ~~to~~ ~~the~~ ~~same~~... it's not clear that they have different codes.
 15 is it pc's ~~are~~ ~~different~~ they have to have different codes. If pc's were ~~the~~ ~~same~~ they may or may not
 different codes.

19 For each context (any derivable context), there exists a α , $\bar{\alpha}$ code for t entire corpus: Since
 all those contexts derive sub-derivs, they correspond to different codes! Hvr, if we define contexts that
 20 are not occurring now, they are ~~very~~ ~~expensive~~ to derive. To derive a current context, you have to
 first where t symbols are (2. dimensional): For non-current contexts you have to tell position ~~times~~ ~~identical~~ ~~to~~ ~~the~~ ~~current~~
 of α & $\bar{\alpha}$ symbols - which has to be record of (k) ^{Nov} ($k = \text{number of distinct symbols}$).

24 \rightarrow So, for non-current contexts: If there are N of them, they have each a pc of $\frac{1}{N}$ so total wt is $\frac{1}{N}$ (only).
 So: (24) way solved! All of the non-current contexts have 1/N wt of (so we can probably
 simply disregard them, since each current code has a wt of

To get pc's of current contexts for each P_{Nov} , ~~and~~ ~~number~~ ~~of~~ ~~symbols~~ (indir of No. of symbols),
 (depending on functions used & their "activity"), there is a number of poss. contexts. (see 220.15 for 2
 form values for binary functions only) - we can use ~~the~~ ~~reception~~ (of this no. to get pc of deriv of each
 20 current context.

31 So: "FINAL SOLN" to context weighting for sequences (or PPM) ~~is~~ ~~the~~ ~~same~~ ~~as~~ ~~the~~ ~~ideas~~ ~~on~~ ~~the~~ ~~relation~~ ~~to~~ ~~2141~~.
 for derivs in which contexts are sub-derivs \rightarrow 223.32 - 224.20 ~~see~~ ~~possibility~~ ~~of~~ ~~deriv(s)~~ ~~at~~ ~~224.21 - 225.40~~

32 ① Sequential compos. (\approx PPM): ~~formal~~ At each point in corpus, we have a nesting (sequence) of
 sequences of contexts (\equiv ~~derivs~~ ~~derivs~~) that extend ultimately over first symbol of m 's corpus.
 for a corpus of length n , there are n contexts, Assoc. w. each context, there are 2 sub-corps:

② Tokens are preceded by context α & $\bar{\alpha}$ next of corpus in which tokens were not preceded by context.
 α & $\bar{\alpha}$ have different case counts for each of t symbols; using these case counts, we can use Lap or Dirichlet to
 assign pc's to each of those sub-corps $\rightarrow P_{\alpha} \approx P_{\bar{\alpha}}$. $\left\{ \begin{array}{l} \text{Note 224.00 - .06 on} \\ \text{+ Lap - Dirichlet diff.} \end{array} \right. \rightarrow$ 224.07

9/1/04

4PM

20

So describe problem exactly (2) Give poss. solns. (3) Tell what's wrong w/ poss. solns

bed: $\sqrt{min x}$

List alternatives in pro-con Args. (4) Is there a way to avoid the problem?

22.30 summary & bit. counterargt.

Perhaps list some good features of proposed Soln.

Consider case of nested sub-caps. (sequential predn., standard PPM)

Since each "deeper nesting" is a smaller, more restricted sub-corpus, each such nesting costs more (bcost) to decb. \therefore lower app. Tho, in the present case, this decb. cost seems much smaller than usual!

To decb. a 10 bit prob. requires at most

$\log_{10} \times (\log_{10} 10 + \log_{10} 10 + \dots \text{bits, rather than } 10 \text{ bits. (i.e., } k(10) \text{)})$

i.e. $k(10)$ v.s. $k(2^{10})$. A similar Economy occurs in 'sub-tree decb. (in the "tree-compression" problem). We just have to use the positions of the symbols in the context, not their values

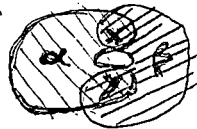
$$\begin{matrix} 1111 \\ 1010 = 10 \\ 4 + 3 = 7 \end{matrix} = 10?$$

10

in the nested case, this is particularly cheap. Tho, contexts themselves, even not "nested"

12:22.00

Perhaps consider more general case of several Coups that have some overlap, but not all elements in common.



So Coups α & β overlap at x & y . Ideally, to do induction $x \neq y$, just consider total Coups.

15

This approach doesn't seem to help much in sequential PPM. Since in the narrow case, α over β completely!

First note that as SSZ \downarrow , the distribution need not get narrower, but it could!

- Case 1: narrow corpus has symbols A, B, C . Broader corpus has $10A, 10B, 10C$. Broader corpus has 10 cases of A predicted; no others.
- Case 2: narrow corpus has $1, 2, 1, 1, 1, C$. Broader corpus has $100A, 1B, 1C$. Broader corpus has 100A, 1B, 1C, 100A predicted.

Poss. predic. rules pick context w/ smallest var? or combine contexts w/ w like $\frac{1}{var} = \frac{1}{\sigma^2}$

Tho, SSZ would seem to be very imp!

25

Another possy ist. $\alpha, \bar{\alpha}$ approach at 220.03, but if N is length of longest context, (so there are only N contexts), one uses Base N contexts in Π , w. w. α & $\bar{\alpha}$ of each corpus w. each $\alpha, \bar{\alpha}$.

26

While 25-26 seems not too bad for sequential predn. for the "subtree" contexts we have such a large set of N codes. But for reasonably fixed corpus & reasonable compression,

The codes that don't use suffix contexts may have small var. (?) No actually, i. no. of contexts that have never occurred, is really enormous!

30

So var. for long being, just consider contexts that have occurred. But have a string of symbols (linking from to taken to be predicted).

So this is essentially the same as 220.06-120 & 220.30-40 & (221.23-2.30)

See exactly why I discarded it & first time I think it was 221.30-33

The idea there was that the null context was by far the commonest context - so if we counted all the PPM's of all contexts in Π , the null context would win because it had so much redundancy.

30

A poss. soln. is that the null context only occurs once. See 223.24 for soln. to this difficulty

See if this soln. applies well to tree induction. If we only consider contexts w. one or

more cases, it may work! Also look at about (12-15)

But it looks like (12-15) is a nice approach to the problem discussed there!

ATM

Huv. in to seq. produ., The overlap properties over "special" a may bear a different

Scratching!

SN Can I adapt Long Methods to Dimensional Images?

SN How does folg. (w/ Dictionary) method compare to other Compression methods?

For direct char of sentence we use char tree, from large English corpus.

"sacra" " " " " " " " " " " " " " " (for first 2 chars)

"sacr" " " " " " " " " " " " " " "

for fourth & subseq., we need dictionary to narrow possib., then 2pm frequency or 3pm frequency.

SSN For comparing speed of 2800 part vs. 500 part, use 32-bit compiler!

can give some insts in "Consul Compiler" or a regular Power Basic Compiler.

1. (02) Consider sequential Cases, P_{now} is what? ! It wa make actual probs for only 1 context each time

the P_{now} is to long term past pc of R. choice of next particular context.

Consider 2 solns. soln. 1 is + symbols to predicted that will occur 4 or never occurred, we give zero pc. & use a different produ. method.

soln. 2. In all cases, all poss. tokens are implicitly considered. If N tokens have occurred

up to now, it is always reserved for new tokens, (There is usually "wasted"). Using a process rule

we have a ~~list~~ of $\frac{1}{N}$ case counts of symbols that have not yet occurred.

Soln 2 will always give a true prob. of $P > 1$ for any future symbol.

Soln 1 will give $pc = 0$ for a prediction, but this will occur rarely in large corpus.

T. discussion of .11-15 is about "escape" problem: EC will not occur in certain kinds of corpus

in which wraparound known to alphabet (say for certain English texts), a known amount set of primary marks & no delimiters are used. Let us discuss this case first (i.e. no escape problem).

So we have many contexts; each can be used to make a code for C. corpus.

We want to use these coding methods on 1: how best to assign bits to P_{now}?

In all cases of produ. methods, we have to assign a code to each method, & this code will be somehow related to part of C. corpus.

In the (worst case) (Bad in this sense at available info), we use to (in part of discourse to assign a priori. Using data from C. past, we assign a priori so that PC of corpus is max.

(Note footnote of (24-27)R) - this can influence a priori much - we want to try both ways.

We can start by listing contexts from N=0 to N=number of words. Maximum available work down to N=0. See discussion of 224-15. 227.18-24 is more inside!

At, perhaps use all contexts for all passes. One uses either or each time. A trouble, perhaps! See 223.24 for resolu. of this diffy. An imp't feature of

"Context" is based on words well in a particular situation - not only of whole corpus.

If we have 2 produ. systems: each w. mainly its own corpus, but some overlap: How do we out of 2 systems when they both are usable? The normal sequential PCM system is a special case in which C. corpus always includes one another! 2.22 12.

10 : 219.20ff & 31-32 do point out how PPM (for sequential) is for tree produ. output (Dmchlett d.f.) to be ideally done.
I'm not so sure that the alternative parsings of 219.31-32 are relevant to 219.21 ff.

13 Say one declares a context α (for "corpus up to now"). This also means defining $\bar{\alpha}$,
1. Complementary: whenever α does not occur, $\bar{\alpha}$ occurs.
So defining α gives us a set of case counts, & each case count can be assigned a PC via Laplace or Dirichlet. (See 219.31-32 for exact eq.). T. PC of corpus is product of ≥ 3 PPs!
① PC of default (②) PC assoc w/ case counts (③) This product of 3 PC's gives wt. of P.
predn of context α ($\bar{\alpha}$). (13)
So a big Q is how to get PC of α (defining) for A good soln. of this problem! (23, 24)

14 An assoc. Q: Given a corpus, what does it mean to ask for "prob. of next token"? What is Criterion of "Correctness"?
One answer is Part 4: present known corpus is a basis for the prep of its context.

15 16 3: OE For 103 ff. Think about sequential (non-tree) coding to start.
Separately 21 each of in f. corpus parts, for each context-size (corresponding to N_{Nov}) (know in any case possible contexts: for all other contexts, no use $\bar{\alpha}$ at this point, & $\bar{\alpha}$ is to "no context" case, in which f. PC's soft. follows naturally their raw case counts into corpus - predictor. This is "default d.f." - pure Bernoulli (Laplace/Dirichlet).
for each "Nov" there is an assoc. no. of poss. contexts. — for seq. coding it is $k^{N_{\text{Nov}}}$ ($k = \text{no. tokens}$). (for tree coding No. of contexts $k^1 = 2, k^2 = 5k^2 = 12k^2$)

15 16 9 Then for all contexts α of this N_{Nov} , one has a code for corpus (via $\alpha \geq \bar{\alpha}$). We know "PC of corpus" with some of all those contexts \geq this is assigned (arbitrarily) P. — (This result is indep of how many contexts of $N_{\text{Nov}} = k$) one sums — i. more, & better, N_{Nov} — Two usually only few have much wt. (30)

15 16 9 Some new PPs in α invitation for functions to facilitate such!
Say f. used for A was α where α means $\bar{\alpha}$.
A B C
D E
F F

15 16 9 So this subtree occupies much more RAM than strictly needed, (A tree of depth d , needs 2^d bytes!)
So f. string: A|B|C|D|E|...|FG|...|000000| — points.
The value of worked many is that it makes it easier to get N_{Nov} of 2-dim. symbols into corpus.

19 The idea is: for each N_{Nov} value, we have a prediction α w. a priori P_{Nov} . It looks like a good soln. to the problem. I'd like to put my mind in a state so that the soln. seems obvious — so I would automatically & suitably generate it.
d.f. function
 α (PC of corpus w.r.t. $\alpha \geq \bar{\alpha}$)

P_{Nov} (P as a func. of k) should be obtained w. a priori computationally: some possible forms are:
 α prior of f. integer (in some possible): (R as 1st approx. a), my approx. for integers only
Normalized $\frac{1}{k}$ choose α ; ($\frac{1}{k}$ normalized: finite no. of values of α used)

So 116-119, 20ff seems to solve PPM for seq. & trees theoretically: Also clarifies exact form of f. kernel, for both cases:
Actually, f. exact form α should be based on our experience in the past. T. function that relates Prob. to N_{Nov} is selected so as to maximize PC of α of corpus. (See 221, 28)

4 TM

$N_{ov} \equiv OVN$ with better to use number.

Nov

Def: $N_{ov} \equiv OVN$. $N_{ov} \equiv OVN \equiv$ overlap number.

10:218.20: On production Basis of "subtree overlap no.". I had that first I would pool "all data from a given overlap no." But by 12 not so good. It is likely that certain substrings would be used for prodn in others "not much". OVN ("is overlap no.") could be an index of PC of subjct. subtree. Also an index of how much bc. occurs zero by defining it. So say we assign PC to overlap trees of $OVN = k$.

Def: Hvy, for coding, we have to specify OVN which is a f. many possible substrings of that OVN, we want to use for prodn. \Rightarrow many OVN variants of PPN of seq. context, copy. There is an automatic way to choose: context. In a case, it's the longest context that a comp. predicts to occur following token. Instead of picking one context, one can use all contexts \geq w. l. from so no "best choice" need be expanded. This is the method I was planning to use for seq. context copy using a kernel to rapidly implement this prodn routine. (My mind is a bit vague; I better make sure soon that I write up a/o how to do better versions.)

I: why f. sequential prodn work, we do prediction in the \bar{L} domain (Lazily ordered contexts). To prodn was a wtd sum of probs of contexts on both sides of the "present" context. Wts could be a function of (either or both) ① long R of common match ② How far away along \bar{L} to assoc. context occurred. To wts would be a func of what the predicted symbol was actually (actual symbol that occurred). $\left[\text{I'm not sure this last dependency is neg} \right]$.

For each tokens there is a 0/1 sequence along \bar{L} , telling whether that token occurred at that pt. " " " we have a smoothed prob of 0 or 1 being at that pt. we have a "smoothed kernel" to get those probs. For each context we have an assoc prob vector for it poss. predicted tokens.

Consider sub-tree contexts: Also our N_{ov} , i. no. of diff. contexts with first \uparrow then flatness $\rightarrow 0$. Hvy, f. SSE's for f. contexts w. small N_{ov} will be very large a func of N_{ov} .

21 I did often want to wt: prodn accuracy to how much PC of corpus, i. assoc. concepts, context. I'd have to use a main concept per symbol that existed under PC's special code to this special concept. O.k. some If a token occurred in times out of m (i.e. β tokens after); Prod I had this factor: β rather formula for α . β (Dirichlet's eq) Prod gave to PC of coding. We have to be heavy in f. PC of defining f. context. Also divider by f. PC of coding to corpus w/o f. depths.

This isn't done exactly would involve Prod on f. ① case counts for all tokens in corpus up to now. ② case counts of tokens in corpus following context of interest.

If we can't use one context to code to corpus, Prod (1) codes can only up f. PC of corpus. (never!) I was afraid that (1) codes would "swamp out" any v.g. code... but this can never occur!

31 T. details of f. prodn. (including alternative parses of f. corpus w/ f. actual data) is done w. 17 my "Letter to G. Wolff" is f. TM notes at about that time.

32 The consider 2. ans of 2. f. f. should apply (more simply) to a sequence (non-free) corpus.

Well I'm not so sure about f. Direct reference case p. 31-32 to 121 ff! In 31-32 we use f. defn of grammar etc. to f. PC of corpus. In 21 ff we use the definition of context to modify pc's of following tokens so as to f. PC of corpus. Hvy, much of the MA of 31-32 is relevant to 21 ff, since we should consider all par. symbols and the effect Prod defining a context modifies f. case counts of tokens not (produced by) f. context.

4TM

$\sigma^2 = \text{var}$

Algebra: $AIC = \frac{1}{2} \ln(\text{var error}) + \frac{n}{T}$ ← no params
 minimize AIC. = $\ln(\sqrt{\text{var error}})$ ← no. obspts.

10:21:40 : The idea of "Value" of some Numeric (or non-numeric object) being assoc. w. many objects, is a very fuzzy heuristic idea. One way to put it into TM: There are many objects will have a number assoc w. them. Or, they will have fuzzy nos assoc w. them. Also Real objects correlate w. one another — the correlation being via proximity (i.e. context) or functional relations. So a certain function will tend to have certain argument types. If a funct has zeroed inputs, they will tend to correlate, so knowing one zeroed gives "in a" a level of other

AIC = $\ln(S_n^2) + 2M/T$
 $M = \text{no. params}$
 $S_n^2 = \text{sum of error}^2$
 or AIC = $T \ln(RSS) + 2k$
 $k = \text{no. params.}$
 $RSS = \text{"Residual Sum of Squares"}$
 $T \ln(G) + n = \text{min}$
 $\sigma^T \sigma^T = \text{min}$
 $T \ln(n) + n = AIC$
 seems
 Dimensionally incorrect.
 We can say of by your calculus factor: But this means adds a constant
 $\sigma^T \sigma^T = \text{min}$
 sets one!

[SN] In thinking about Ngrams, branches, substrings etc: In a sense of PPM, these substrings would give a PD order to them and an adjacent part.

Here, one big use of Ngrams is that they are often predictable following others ngrams. This seems like a decent (i.e. important/different) way to look at things! If we define objects as make tokens of Ngrams from these generalized tokens, can be both predictors (context) as well as things predicted.

We can achieve this by making deltas, parsing corpus several ways w/ diff data. i.e. periodically re parsing & perhaps periodically removing old deltas as well as creating new. Reparsing is not such a big deal if we can reverse O's each time we modify it. We could also parse previous O's

we can also have several || parses for a given O's. → 213.00

[SN] Res Ops: One big problem in Phase 1 is that there is a O^j that has no work for all Q's: So if we make a window, we have to test for on all past Q's. If we are able to work on each problem anew (as in Ops) would not have that problem... but fr. cont. of transfering is from very limited.

[SN] Re: GP Complaint was that "GP sometimes broke up good branches". This was perhaps because the selection rules were poor. A better way would be to keep top k percentage of population as in no cross/mutate from that only. This is close to what's done when best population is selected by a p.d. we may want to keep all or almost all of population & use a selection of parents based on G (= fitness)

[SN] Use default methods for function trees (e.g. the one described in: How use PPM on corpora: Compare redundancy and various n-grams. There should be nice way to combine the 2 sides here (Backward & forward) on the codes of Look at being and how far good way to combine them. Maybe if G will be used to decide.

[SN] A big Q is how much to use each O's in corpus (for updating): One (seemingly "logical") way is to use G. no. of times it was successfully used w.o. Modify. So Design of TSC will end up deciding on Wts.

47M

0:216.40 : update : $\left(\frac{\text{overlap no. limits}}{\text{time/cc}} \right)$ can be adjusted so that we get best cc for corp's for

02: 216.37 : Actually ~~the~~ ^{the} trace node notation (e.g. 2.3. traces) is not so expensive. W. z. corpus of ~~2.10⁶~~ ^{to 4.5} I have 500 x 10⁶ ^{bytes of} - So 500 by problem. Certainly adequate: My 2.3. trace record "only" ~ 36 by per problem & corpus.

This trace notation does make comparison cheap.

A poss. relatively "rudimentary" way to get a lot of "points" across w. large overlaps: write a particular product p_i . Say α is a pt. of comparison: Get all p_j in corpus w. "big overlap on (left) using (RPN) and PPM". ~~using~~ ^{using} $\{X_i\}$ \rightarrow complete overlap comparison to get the more complete sets. (Note 216.38 - 217.01 on cc of "trace matching")
It might be usually to RPN, PN matches will give overlap more of > 8 , say, so we need not search further. Only if they do < 8 will we fall out to "Big Guns" of complete comparison of $\{O_i\}$.

0: 217.09 : TSQ's in English : For ANL (i.e. SHARBBG W): Need to do just what the problem is ...

What needs to be done. So, to some extent, we must admit, choose a language/notation
So list = set of poss. notations, for Algebra: e.g. RPN, PN, infix \leftarrow Continuous kinds: APL uses a mixture
What notation does APL use for algebra? Maybe infix makes it. Computability of $x \in \Phi$, easier to deal w/ false advantage of.

Then there is to ~~be~~ ^{be} considered mechanical. It would seem that all of these choices must be made before the ANL problem is definable.

Are there some hours that are useful ^(adequate?) independently of some ^{which} choices ~~is~~ ^{is} made?

Some "indep" hours $\text{Incl. } \text{array} + - * /$ problems; The idea that $+ - * /$ are simple "controls" may be "ob" output that tells us what to do. [how does this differ from $\exists, \forall, ? \dots$ that also "tell us what to do" ?]

Abstractly Stated: we have $\geq T \rightarrow Q$ of $[Q_i, A_i]^n$ pairs: To derive a sequence of machine symbols (\equiv instructions) O^j such that $\prod_{j=1}^n O^j (A_i, Q_i) \equiv \text{Max}$. $\equiv f(O^j, [A_i, Q_i]^n)$
TM is just given. T set $\{Q_i, A_i\}^n$ is for any O^j and, \forall value of

I guess that operators like $+ - \dots$ are numbers and "different" in ways we can deal. A small no. of operators, a larger no. of numbers. Numbers are all similar in certain ways.

The initial lang: TM sees $\exists \geq 7 \rightarrow 10$ figures that $\exists \geq 7$ are "yes" and are acceptable signs of binary functions. "+" is an operator name. ~~so~~ ^{we} ~~we~~ ^{learn to} ~~we~~ ^{associate} it w. +.
operator $\text{sum}(x, y)$. [Polish]

We could initialize TM w. f. knowledge that $\text{sum } x, y$, $x \neq y$ must be "numbers", and "sum" is the name of an operation $x \neq y$ that gives a number as output. We show "tell" or "teach"

TM what "numbers" are: how to recognize them.

A Big Q is how much we want TM to know a priori.

It may be that choice of lang for defining probs, & choice of instr set for the machine is mainly important for ANL: Most for problems past that point, discussed in "English" can be made that do not depend on these choices. E.g. we should be able to solve e.g. eqns. quite

independ of these choices ... debr. our hours independ of these choices
Hvr, a value of having TM know what numbers & functions are, is that later, these techniques could be useful in lang. new, novel concepts - quite unrelated to res. & functs.

2 : 2:15-30 SM2 On Cramer's early formulation (1965) of GP. ^{← somewhere in PS! 2004 I think} Seems he mentioned a lot of good heuristic tricks that were not taken up by the GP/EA community -- look at them!

I have a recent paper written as a MS Project in Eastern Europe (Nov 2002 meeting at Stanford). In searching for continuous cases (regular or user ODE's) to solve a problem, he restricted search space by limiting forms of functions considered: This is one form of HINT that I may have to use in trying TM!

Some restrictions he used: (1) equis of ~~max~~ max degree 2 (2) He gave it to function $f(x,y) = \frac{x}{x+y}$ or $f(x, const) = \frac{x}{x+const}$

This fit to PC form $\sim 10^{-15}$ to $\sim 10^{-8}$. Since T. fitness function was very expensive (he had to step out the ODE's w/o optimizer param of functions & code), this was very imp.

? $pc = 10^{-15}$ would be very expensive, if pure Lench was used, but perhaps $\}??$ quite accessible if GP were used.

I found .02 while looking for papers on common subtrees in Google. Much more on largest common subtree. Parikhony on other such things

0: On searching for Subtrees: A possi (expensive) way to represent trees: Way I did it in the 2-3. Tree rep: Each node has a list of nodes that are its kids. They are arranged in an array: The nth row of the array was next: content of Node n: followed by node nos of children of n. OR nth row has contents of children of node n. each content followed by node nos of next child.

24 This may make it faster to find similarities of nodes. To "compare" 2 nodes (to find how large their overlap is) first we note that they have same "content". If they are the same, we examine children. (Have we an order of children? Commutativity is a Real symmetry.) We do this by putting one around (as per) children out of stack. "tracing out of identical nodes"

Actually, we may be able to do this w.o. the expensive representation of 20 to 24.

0: E.g. look at the graph of (20, 24, L): we trace out "forward": A B C ... P: Pure process of 2 nodes Graph. Suppose we want to find "overlap no." of 2 diffrnt graphs at point α . To 2 graphs have same child at α . We trace along as "A B C ..." two nodes at which to 2 graphs are identical. If they don't match we drop that node & go on to next. We keep track of number of nodes that common nodes that 2 graphs had as we traversed them, we only ~~traverse~~ "traverse" to common nodes that are "connected to α ".

The mechanics of how we cut out a non-identical node/branch, is unclear: How we do it depends on no. of nodes: Straight RPN (or PN) is cheap, but cutting out non-identical nodes is expensive.

37 T notation of 20 & 24 is expensive, but omitting non-identical nodes is cheap. In 4. loops, we may not want to trace out beyond "overlap no. $\leq \alpha$ " - In PPM, they did not get much better by considering much contexts... (This may be because they didn't do it correctly!) This "limit" can be adjusted to - depending on how much extra we want to spend on 217.00

for Radix 3

0: 214.19!

In f. sequence of integers representing $n: \log_3 n, \log_3 \log_3 n, \dots$. The leading char. can't be 0.

So we start out when f. leading char. is zero, then we know the next integer coded will be "0" for us. we want to represent. Say $f(n)$ is the no. of ternary bits in n 's expansion in 3.

$f(1) = 1$ $f(3, 4, 5, 6, 7, 8) = 2$ $f(9, 10, \dots, 26) = 3 \dots$

So to code 25 we need $f(25) = f(23)$ and to code 3, we need $f(3) = 2$. To code 2 we need $f(2) = 1$

so for 25 we need $f(25) + f(3) + f(2) = 6$ symbols, + 1 extra zero

Note 1: leading symbol number 1/2 being 0 we get into infinite loop. So only number 3 and 4 can be coded properly, so 0 is unnecessary; but 0 is still useful; so for code 25 we need $3 + 2 + 1 = 6$ symbols.

Since we can't code 0, 1, 2 if we want 24 we have to add 3 to the integer before we use radix code.

It's not so clear of 07 ff. how we choose to code $2 = f(3)$ if we wanted us 3.

The first 2 symbols always represent 3, 4, ..., 8. If the first symbol is zero, then to have 2 symbols to represent $n = 3, \dots, 8$.

If the first symbol is not 0, then the first and second symbols represent 3, ..., 8; but then we have the no. of symbols in next number (which will be 1 if it starts w. 0).

So, if the next symbol is 0, the next (3, 4, ..., 8) symbols will represent n .

So the Q is, will the prefix of any ternary sequence be 2 (ways represent an integer).

12 = 5	22 = 8
11 = 4	21 = 7
10 = 3	20 = 6

There may be an awful lot of infinite nos!

certain points in the generation a zero will start a termination. But, chances of this sort occur less and less frequently, so it may be that almost all infinite nos. to infinity are of almost all infinite ternary sequences are of length ∞ ! Hur, lots of long sequences

would make $\sum 3^{-n}$ converge! - it may be that $\sum 3^{-n} < 1$!

0: 214.05

On choosing between LISP & Forth: In Forth we generate f. output sequentially, so when the cond. is only "1 desc'd", we have implemented 1/2 of its "test" (say for INU probs).

In many cases, this is impossible - that's. Test requires full entire cond. Probably if we used

(as we did/does) to splice from the test cond., we usually need the whole cond. before

we can begin testing. If we can do this, then LISP does not have that disadvantage.

Second Q about LISP vs Forth is whether PPM or other regularity

could more easily find reg's in LISP v.s. Forth code (Also consider Machine Code)

2 APL and 2 BASIC, Fortran, C // what about Object Oriented langs?

Third Q in choice of Lang. is speed - The it may be poss. to modify slow

langs like LISP - by compiling them 2/0 by making primitives close to Machine lang.

Another speed point is to average lang. definition so parallelism of modern chips can

be used. There is a JAPS chip in development at \$15/Buytop - very parallel, but

One must properly design search tasks for it.

4TM

Also ideas from ACP.

10 : Prospectus for TM: A path of work for future. Give various alternatives & choices. List various problems &, if possible, histories of where I've done work on these problems. Put this on Computer, so can modify & rearrange independently.

15 → Perhaps 2. MAIN idea: Just I want to get started on T SQ as soon as possible, because many other problems (like choice of R 's, "Machine", \approx "syntax" to use; Forté v.s. LISP (i.e. sequential v.s. ~~complete input~~), perhaps form of a proof on "integers" (2.12.20) ~~write also~~ solid by f. T SQ's in Rear Hours. → (2.10.11)

Describe in some detail "Interrelated Approaches to TM" of 10.5.18! Have up your Phase 1, Phase 2 approach and to TM, = The approach described in Abstract of the I D S/A Report.

I think the main idea is to start w. T SQ writing, in English. Get all off hours I want to use, for each problem. And an idea of which corpus I'd need to discover the hours. → 2.17.10

10: 2.12.32 → SN This system gives a certain strip of integers - (I can adjust the amt. of the strip for short v.s. long R 's, by using different α or by other means). There is, however, Q of all other inputs to the 3 input one - that is, to have all the different pc 's... How can R deal w. this? Well, perhaps "No problem". If we are writing & taking using R , we use up pc in the same way as we would for any other one input - we just multiply the pc 's together. - But, my mind is not \approx together & Q

→ about 1/4 of question! Perhaps considerations of Prefix set properties & Normalization and order - that 1.13 is adequate for the problem!

Well, the Q input: We don't know what the pc is. - But the O & R inputs are "self deal".

An aside: Could we treat $\log x + \log \log x \dots$ seq. divergence base e : (is possible for base e .) Can't devise a system (like 2.11.14 for base 2) but use Base 3. - it has to converge!

→ But it may not be the expression → 4

20 SN Re: T. PSG-dicy problem. It may not be so difficult! T. problem is dirty may have been that people really didn't believe in the proper Game (2) If Ray had, Ray didn't try many pos's. (2) May may not have considered parallel parsing as part of the Game! If Ray tried to solve the identification problem (other than the prediction problem).

Also note: I'm not very in involved in PSG. I mainly want to find common sub functions (Sub trees). There may be work on this: try Google. Also after search engines:

SN Look in Google: Phrase Structure Grammar Discovery
Context free Grammar Discovery
 Also try Yahoo, Excite, Altavista, MSN.

I tried Common Subtrees in Google: lots of pages, usually "Maximum Common Subtree" ... See 2.16.02 for Best SN

Try other Search Engines.

Tree Compression Got a review article! (See PS/8.2.8.04)

Also Look in 2 citation indexes: forward to work by Shen (?) or other PS Also, Review of 1976 at Harvard will have Bibli that I can use for forward refs.

30

2:212.40 : I could use R to "Modulate" the effects of the best corpus, on the present problem.

In COPE: It is possible to have a "problem fixed" O_j program $\rightarrow Q_i$ objectives } R is a universal
many, or on a stack, a have a R program generating P codes, sequentially, } def on A_i or
} subset of O_i & Q_i .

DEF Going back to 204.39-40; 206.00ff model (which I'll call "Model M"), we can regard each A_i as being a function of Q_i : But by having the "property list", we help to search for good codes for A_i . Say we had many properties & a tuples of properties & we correlated them with the freq. of use of various Tokens & Macro-tokens in the correct A_i 's. Even with best properties in N-tuples of properties - all we get is a v.p. of ~~tokens~~ expected use of tokens in the target A_i . Even if it were sharp enough to give us exact freq. of use of the best tokens in the best code of A_i 's it still not $\pm PC=1$, (R could not be rather good ... it would still be extremely "fuzzy").

This expected freq. of use of tokens is simply a general improvement of COPE's boosty. It's a good heuristic, but it is certainly not entirely adequate.

A possibly adequate ~~Model M~~ "Model M" (03) would have lots of Macro-tokens - so the final soln. would, ideally, use only a few of them.

Another ^(input) ~~Big~~ augmentation of Model M: Model M finds the best p.d.'s for tokens - Best in the sense that they give the best PC for the corpus. Hvr, the p.d. of a token at a particular pt. in the generation of a code, should depend not only on properties of Q_i , but on previous tokens that it has used: This is in direction of PPM...

How can we "mix" PPM in. ~~is~~ effect of properties of Q_i (\equiv Prob (M))

Or the p.d. of tokens should depend on "order of the token" (how far ~~they~~ along in generation of the program for target A_i) appears? (This list is "similar" to previous local "previous history" as opposed to corpus "previous history".

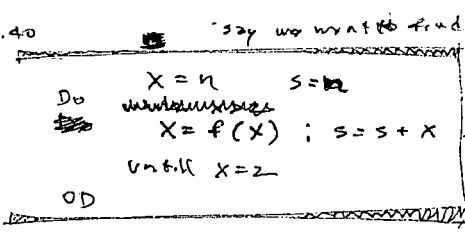
The Model M "token probs" are like "unconditional" probs: ~~the~~ conditional probs would be based on local recent history. So the conditions should be local ~~corpus~~ previous corpus is local previous history.

8/24/04 Note that it's O_j that has maximum about $Q_i \rightarrow A_i$: R ~~is~~ used by O_j to implement the p.d. Would it be ok to put the continuous into about a p.d. into "R"?

: 3 p.d.'s? ① Original ~~program~~ before TSQ; ② ~~program~~ of Sumac after TSQ has been run. ③ Guiding p.d. Actually, a GPD is a p.d. on {Token/Macro-tokens} - R ~~is~~ is to a p.d. on the soln. of latest problem.

4 TM

20: 211.40



say we want to find how many bits we need to represent n
 S will be required output if $n \geq 2$
 $S = 1$
 Print I computed this to f however const.
 $A(\log_2 x + \log_2 \log_2 x \dots)$... but I may have computed incorrectly.

$$f(x) = \lfloor \log_2 x \rfloor + 1$$

$$= \lceil \log_2 x \rceil$$

05

06: 201.15; Hungary returning to the shop; In a 3 input case:

One input is Q_1 , another is Q_2 , the third is R . We find $Q_1, Q_2 \in R$ into Per Machine & R outputs f , standard function. For Per R input, we can use binary integers of fixed length. We can assign R bits to P 's via $X = \lfloor \log_2 n \rfloor + 1$, $P = \frac{1}{3} 3^{-X}$ ($X=5$)

or via $2 \log_{10} .00 - .05$ 2^{-5} (computer for $n=17$: Alg. $.00 - .05$ gives $2^{-11} = 2.048^{-1}$)
 $X=5$
 say $X=17$ $s = 28 (= 17+11)$ $\frac{1}{3} 3^{-5} = 3^{-6} = \frac{1}{27} = 729^{-1}$
 2^{-28} v.s. 3^{-18} $3^{18} = 3.87 \times 10^8$

So for small n , $\frac{1}{3} 3^{-X-1}$ is larger, but for $X \geq \log_2 n + 1$, 2^{-X} is larger!

But $X=17$ means $n \approx 2^{17} = 128k = 1.28 \times 10^5$.

Well, P 's pt. is not critical.

Print P 's for large n , Per 2^{-5} & 2^{-n} but for small n , it's much smaller

hvr $\frac{1}{3} 2^{-X} \approx \frac{1}{9} n^{-1.585}$ ($\approx \frac{1}{9} n^{-1.585}$)

$$3^{-X-1} \approx \frac{1}{3} \lfloor \log_2 n \rfloor + 2$$

$$X = \lfloor \log_2 n \rfloor + 1$$

$$\log_2 n = \frac{\log_{10} n}{\log_{10} 2}$$

$$\frac{1}{9} 3^{\log_2 n} \cdot \log_2 3 = \frac{1}{9} \left(\frac{n}{2} \right)^{\log_2 3}$$

$$\frac{1}{9} \frac{1}{n^{1.585}}$$

10: Use power binary integers for R , 0, 1, 2, 3, 4, 5, ...
 = 0, 1, 10, 11, 100, 101, 110, 111, 1000, ...

2nd we see if f outputs A 's

we want to simulate R for Nat outputs. The PC assigned to that R can be zero or fixed methods, 2.28

Or other methods; Normalized $\frac{1}{n^{1.585}}$ into book. We can use a normal $\frac{1}{n^\alpha}$ for any $\alpha > 1$.
 (The normalized output for $\frac{1}{n^\alpha}$ is related to Riemann zeta of α).

SN for normalized $\frac{1}{n^\alpha}$ v.s. $(.00 - .05)$, $A > B$ for smaller n . — but there is always a value of n for which $A = B$: for $n < n_0$, $B > A$, for $n > n_0$ that $A > B$.

So for very large n , B always wins; for small n , A is always better.

28: While this works fine in theory it is quite wasteful of CC: Each value of R has to be tested individually — for 2 different R 's, there is no sharing of common tests, & it would be true, if f random machine worked along R seq generally (as in COYS or for P).

[Think: distance amounts to a factor of d , where d is f , dist of f from \mathbb{Z} + length of f could decrn.] R is may be indep of Radix... but I'm not sure } $2.14.10$

Hvr dropping that for a moment, suppose I had to Q^j that summarized previous corpus using Z ideas of $204-34-10, 206.00, 40$. This gives a default pd for A as a function of Q^j . Should this be pd for null R ? If so, how does R modify f pd on A ?

well to default D.P. Z is pd of all tokens in machine to learn on previous corpus.

A more general pd could be based on P 's out-copied R 's form. Also there is a sequence litym f -ideas & $204^{st}, 207.00$. That's quite different from the Lisp/AZ approach. R 's is a non Lisp/AZ approach!

N.B.
 One of the Rad. Machines
 I must consider is \approx APL

ATH

Conclusion (proton) .11 ff

20:2081 to 10 remain old pems: we just want to vary wts so as to max (i products of P's of solns of all of the pems for t. A₂. ~~In~~

Also in "update" we can try to find new Obs (features) of Q's 2 ways to relate them to t. pems for t. Assoc "A's". These "Relations" need not be simple wts. They can be any functions

T. ~~func~~ stuff: we have a Q_i's obs on it is a func of these obs, PCs or tokens 2 "Macrofuncs"; that give a PD on Assoc A₂. The only way we can vary this PC on Q_i's is to change ~~the~~ any new Obs, change relation of obs to tokens & Macrofuncs. ~~the~~ removes a subset Macrofuncs.

We want more flexible way to vary the relation between Q_i's (its obs) & t. final PD for t. associated A. The "present" way is by varying PCs of tokens; Macrofuncs - this is pretty

Not to much General way: "Real Macrofuncs" is a big step towards Universality

Conjecture: When I have hours for QA problems ~~is~~ - a TS Q of them. - I ~~do~~ expect that any attempts to implement these hours will surpass ~~the~~ to forge ~~the~~ Questions

T. hours will ~~be~~ suggest solns & perhaps ways to implement these ~~the~~ hours using ~~the~~ forge models & Modifys of these Models.

ATH: Discussion of t. R funcs (= Obs = properties) in t. early discussion of QA probs. in EDIA reports is very closely related to forge. - ~~Actually~~ The R funcs are meant to deal w. d Funcs (cf. NMTM). But Actually, R funcs act. Best to Q's "Best to Q's Gen":

R funcs are pure Obs: they map into Yes/No. They then control P_i(A|Q_i) functions, which are s funcs of t type that I'm trying to Define/ Create/ Engineer. If t. R² funcs were narrower, they would arrange P_i(A): i.e. t. Q_i need not look at Q_i. R² has already looked at Q_i. ... We have 3 ~~cases~~: ① No R at all in which Q² look at Q_i & give PD for A (O²(A|Q_i)) and ② P_i(A|Q_i) in which R² recognizes a set of ~~the~~ Q_i's such that a certain P_i(A|Q_i) will work for this set of Q_i's ③ R² ~~can~~ will look at Q_i ~~and~~ decide Yes/No where its associated P² will generate P_i(A) for that Q_i.

So I ~~still~~ don't have an adequate set of ways to realize implement parametrize P(A|Q_i). 207.20 A is still looks like beginning of a good way to give More general P(A|Q_i)'s

So far I have been using only a few ways to "Modeling" P(A) (= PC of A: a s func) - Obs on Q_i correlate to PCs of tokens & Macrofuncs parametrize A.

In 207.00 - 19 I had some "Special Situation" ways to characterize s funcs:

- 1) Fixed form, but set of real param. - funcs of Q.
- 2) fixed set of func func & param func of Q
- 3) Universal set of funcs; funcs of Q; also param of f₁ & func of Q.

8/22/04
ATM

Limits of Human Reasoning/Induction 1.00

o: [N] On HUMAN (very, induction: T. way / ^{human} languages are created, lang: Humans are born w. certain built-in deduction, induction facilities: Human langs are distinct so they are ~~unique~~ (unique) languages. These facilities are probably not universal, so certain ~~languages~~ Synthetic langs would not be (ruled) by a new ^{born.} ~~born.~~

Also, ~~the~~ humans may use ~~these~~ non-universal reasoning facilities for ~~all~~ ^{all}

Domains - so we would not be able to ever discover certain ~~new~~ ^{new} ~~things~~ ^{things} by ~~our~~ ^{our} ~~own~~ ^{own} ~~power~~ ^{power} ... ~~certain~~

kinds of physics (cosmic natural) laws.

On the other hand, we may be able to discover: machines/people that can do reasoning

way our ~~our~~ ^{our} ~~own~~ ^{own} ~~power~~ ^{power} available to humans.

~~What~~ Turing/Partial Rec func's are able to express all Human ideas of functions.

Could there be ~~the~~ functions not accessible to humans - perhaps accessible, but perhaps not accessible to Turing.

As it is, Human induction does seem limited to ~~the~~ "Phrase structure" langs.

~~For~~ = CFL's: Those we can include context sensitive langs... which perhaps makes them universal (TAG lang is ^{very} "simple" context sensitive lang ~~problem~~ investigated

by Post, is is universal).

It may be that humans can only deal w. very limited context sensitivity.

T. ability of humans to understand ~~the~~ (non CFL) concept "recursively"

probably does not generalize to all context-sensitive langs,

→ Another possy: That humans can get beyond their inborn reasoning facilities, but they do this ^{awkwardly} wordily as ~~by~~ ~~CC~~ ^{by CC} - so machines could be way ahead, here.

Humans may be able to transcend ~~to~~ ^{to} ~~newborn~~ ^{newborn} reasoning by writing things down, by drawing pictures or their logical equivalents.

ATM

0:207.90: The D.F. on A_i depends on t previous solns, but in a fixed, standard way. So it's a fixed function of past. A ~~(perhaps)~~ goal: To make it search "around" that past
 → One search can be over properties (Obs) on t Q_i 's.

A possibly related problem: we have a coin which in random input, gives a certain (continuous) d.f. on its output. How do we modify that coin to generate an arbitrary pdf on its output?

Since it is universal, ~~we can't get every dist. with a coin~~

∃ a finite ~~input~~ initial input string, that will convert this coin to any other ~~coin~~ i.e. any other output d.f. This is "non-destructive"

Since t pdf's on Tokens induced by 207.20 ff are all non-destructive (i.e. they give a ~~pdf~~ $p_i \neq 0$ for all poss. choices) I think t resultant d.f. has to be Universal. ?

I'm not so sure of (07-18) $p_i \neq 0$ for all choices means that no coin has zero p_i

A machine w. ~~pd's~~ pd's at each instruction point is ^{"invertible"} equivalent to a deterministic machine w. a pd on each input symbol. If this pd on input is invertible (i.e. ∃ a pd on strings that converts it to t uniform d.f.), then the code for that machine will be perfect, ~~and thus?~~ original machine into a Universal machine. — so that ^{invertible} machine must also be univl.

If a machine of 10 simply words has a constant t some modification is, unimp'd $\frac{1}{2}$
 $\frac{1}{2} \rightarrow p_i$ for all definitions; we can use Normal $\frac{1}{p_i}$ $\left\| \begin{matrix} p_{i, new} \\ p_{i, old} \end{matrix} = \frac{k}{p_{i, old}} / \frac{k}{\sum p_{i, old}}$

If t pc of a token, F_2 is a function of the set of all previous tokens, we can use t method of 15 to invert t d.f.

If a pd. assigns p_i to n as well as tokens, we can also express it as a pc on t next token only by summing p_i 's ~~over all~~ of all inputs that began w. a particular token. (? — what about t pc's of t second (say) token in one of t inputs?)

In 16-19 Powers is a parsing problem: If t token $\alpha\beta\gamma$ occurs we don't know if t first char. was obtained from th. n or from token, α .

Perhaps split: Say we have a grammar that assigns p_i 's to tokens and to n 's.

At first consider grammar:

$$\begin{aligned} S &\rightarrow \text{S N S} \\ N &\rightarrow \text{a b c} \mid \text{a b b a c} \\ N &\rightarrow \text{P}_1 \mid \text{P}_2 \mid \text{P}_3 \end{aligned} \quad \left\{ \begin{array}{l} \text{a, b, c are terminals, S, N nonterminals,} \\ \text{P}_1, \text{P}_2, \text{P}_3 \end{array} \right.$$

This gives a pd over all strings of a, b, c: Is there a grammar that can, w. these strings as inputs

put out a pd. w. ~~a~~ pc of a, b, c = $\frac{1}{2}$? Well, that's not what I want.

to get t uniform a b c is not hard: for any input symbol, output $p_i = \frac{1}{3}$ for a, b, c: !!

What is it I'm really looking for? A deterministic grammar?

I could get that by listing all poss. outputs as a pd. assigning $\frac{1}{3}$ of them to a, $\frac{1}{3}$ to b, and $\frac{1}{3}$ to c.

More precision by considering seqs of 2, or 3 or ...

Go Back to 207.20: We have this set of ops on t Q_i 's. We have t set of associated A_i 's. We have the parms. But have created the A_i 's. The pc's of tokens & subfunctions, a c's of A_i c's of A_i parms have been obtained as "correlated" (or ~~other~~ other functions) of t obs on associated Q_i 's. When we get a new Q_i , A_i i.e. t parm assoc. w. that A_i , we try to adjust t relations betw. t obs & their correlated tokens, n 's a c's, to be "better" i.e. so that t total pc of t corpus is max'd.

Updating will consist (of each) of optimizing continuous parms (als) involved of t new data pt. Since t parms for all of t solns are invariant, we don't have

210.00
210.00

SFUNCTION realization Optimizer for induction with

Memo on Sfunction Generation, representation, induction:

I can use many different forms of Sfunctions, for different kinds of problems:

I can give TM the form & have it find continuous/discrete params for the problem.

E.g. for $h(x)$ problems, I can give the Gamma D.F. or the Gaussian D.F.

is 250 for 4, 3 params.

In other cases, I can ask for a pretty, messy not D.F. & one or more other moments of the D.F.

For discrete D.F.'s • one (or more) cluster centers, and a "distance" rule w. error continuous & / or discrete Params.

For some problems I may want to use the 3 input time model, but I really don't know how to do this effectively. — The 1. OOPS "four" model would work.

I have worked out some details on how to do this, here. How practical it could be.

OOPS Forth, is unclear. There are 3 input params: O_i (the general model)

Q_j : No 1st question is R's "random" no. We could arrange so TM knows it has to read

O_i first, then Q_j . Then try to make a pop to get A_j that is as short as possible. ... Then we will

be interested in several other pops as well. — (See .20 - .23)!

The idea of an "Epic" TM that knows it can read for easy to get info; perhaps this would enable (.12-15) !?

if I had some how is for realizing the soln to the stochastic QAP problem — then, I should try to fit the 3 input time of .09 ff into that / those heads

A v.g. way to get Sfunctions: 204, 24, 40; 206, 00 - ff: — This considered QAP problem!

A bunch of Properties (≡ Obs) on the Q_i 's are (found/given/true). These Obs is pairs triplets of them, induce d.f.s on the individual tokens & on functions, branches, etc. The output does correspond to that of a 3 input time. (Huh, this language doesn't look at the details of Q env)

I think it was a very fuzzy answer, but perhaps I can modify it to "do it into Phys"!

Also: I want it to give nice answers in deterministic (MTM) case. The methods of 204, 206 don't seem to do that, but perhaps I can get them to!

Say Q_i has scored better a TM has submit into answer. If Q_i has some answer, ^{Rest time,} no problem — but if the answer is ident; All properties of Q_i are the same set will try same instructions & functions that worked best time, then "move on" to less likely conds: & long, of course, will still be universal.

Say Q_i is "big" Q_i : lots of properties bits: But a few differ from the last (correctly answered)

Q_j : Depending on properties, Q_j is Q_i may or not look "identical", but for Q_i , TM will usually try the tokens & functions used as Q_j .

What this system does not do well: in .29 - .30 say we have very long function w. only 1 numerical argt. Different from rest of properties Q_j . (Like $Q_j = \sin(\pi(3 + \cos T))$)

TM's "way" is "4" in Q_j . We'd like TM to be able to localize the difference between Q_i, j & make a "narrower" trial set.

Problem. .24 as an SEARCH implementation of 3 cv-form decoupled Sfunctions: The property list of Q_i contains a complete set of the problem's & new one to Q_j . There is an inclusionism to search or to try new kinds of Sfunctions (via D.F. or A_i). It is not a D.F. on D.F. is on A_i 's but a single D.F. on A_i 's.

I think I need a Universal lang to desc. the kinds of PD's on trials that .24 affords

Modifying Property list & / or relation of parts of tokens & sub-trees, would not seem to be "universal". 20600

See if S can use
t. "EM" (≡ expanding
Maken) method
for any of my
param opten
problems.

FTM

00: 204.00:

Also, for ~~the~~ English, parse trees ... of properties, we can get probes for Sub-functions of various pairs,

~~the~~ rules of token - This would be a v.g. way to find subfunctions.

T. Subfunctions (or ^{or} ~~rules of tokens~~) could also correlate w. previously generated parts of ϵ cond. So, ϵ function/leaf chosen as a p.t. in ϵ cond. generation will depend on ^{properties} of the problem also on the previous choices in cond. generation.

Both of these choices will suffer from small size so we will want to ~~select~~ ϵ of ϵ token selected, depend on whatever "property of problem" we have data on - ϵ size of the data: The recent context in ϵ cond, ^{its} size (since there will be many contexts - each w. its own size)

The "properties of the problem" can include External Context (e.g. Chemistry problem, Dynamic Programming, Linear Algebra ...) as well as context based on studying problem down itself.

We can use "Encyc" to occasionally decide what kind of problem it is:

In all of these, TM does not have much "understanding" of what ϵ problem is: The properties ϵ show correlates ~~are~~ ^{are} ~~indeed~~ a kind of "understanding" - but not as good as a person understanding "the text of a problem. Hopefully, after TM has worked lots of Algebra (say) problems, it will have found lots of useful properties of the "problem down". Would it be able to look at a partial Diff Eq & decide what kind it is so it would know what solve methods were possible? ... well, this sort of thing can be found like TSCQ.

One kind of Understanding of Q: It will have to understand "solve" for various Alg.

~~Maximize~~ Maximize (for Discrete, Continuous domains)

At first, we only give it Induction, INV, OZ probs. Later, we only INV, OZ probs, later, it finds out how to convert INV to OZ. ~~In~~ Each problem type has a standard form, so TM really knows what is required.

How, even after a lot of training? Consider QA problems: TM could have lots of examples "understanding" + Q's - But part of the QA problem TSCQ is learning to understand what the Q means, e.g. in 18-19 w. "solve", "Maximize"

Also, I expect to teach it to understand English about ~~the~~ domains of its Competence

APL (KJ)
put A Q-19
1900 + 1901
Kings

4TM

Ranz. - Husband - Bladder cancer. 23 wks ago.

00: 203.22 : Also, continuation ~ 199.32: in this next situation: GPD is able to learn:

.01 At first, it adaptively modifies pc's of fishery. say ~~the~~ Lapsor Dirichlet's rule

.02 Next PPM is used, Next turns on PPM, next P5Gding - using Stochastic w/o my improvements Analysis

03 (3) on next level, say I look at TSO, & surplus of sales, & I notice ~~copy~~ that could be used for prodn, but Program not accessible to TM Rev PPM ... P5Gding act.

04 -> I could modify root (say so next this (these) kinds of copy would be accessible via

06 PPM, extended PPM (201.00 - 202.60 - 203.00 - .20) P5Gding, act.

If .05 - .06 fails, I want to look at finding ways in a couple of PPMs as a problem

It is solved by ~~the~~ (search) over a random set of resolutions defined by a universal law.

- So essentially, need like Phase 1 to be able to work on this problem. (TM1 = TM2)

10 TM1 looks for solutions in the primary ~~the~~ TSO, using a GPD that is at first, updated by methods like .03 - .06. When Ross fail (or even if they still improve (as in work), we try to view

12 TM2's problem as being "improving GPD".

13 .03 - .12 looks like a return to the original IDSA reports: The system that was described in Abstract of this report. It does not mention: Phase 1, Phase 2 system. What's more that abstract,

I did not know about the ideal soln to the OZ problem using h() function evaluated

and the ~~the~~ slope param. (Ideal solution to Gamb. House problem ... better than search ... that

optimum

18 It is ~~best~~ desirable / good idea, to integrate the old "Abstract" approach, w. the Phase 1 & Phase 2

19 Phase 2 approach. Int. "Integrated" method; we start w. Phase 1 & use the methods : Discrete, continuous

20 of .01 - .02 for updating the UPD. We eventually begin giving TM problems in induction : using various

like Ross needed for GPD updating. When it gets mature and in the area, we have it : standard forms?

work on its own UPD updating. In this last problem, there are many ~~many~~

methods TM is already familiar w. (like the methods of .01 - .02) If gives certain

pts. to Ross methods, but ~~the~~ ~~develops~~ ~~and~~ ~~has~~ it has already had experience

using the h() methods (.13 - .17) on OZ probs & it can be used for the END prob. of

GPD updating. [Note that in general, it will look at any problem, & use Obs (properties)

of the prob. even to help solve it (e.g. 204.34 - .40, 206.00 - .11)]

How to get work
P.D.'s:
Discrete, continuous
using various
standard forms?

FTM

00:20:25.90

Re computation of top of 203 1/2: using "add on" only:

(2) $21 + 4 = 25 = \text{no mts}$ $\frac{(25)^6}{2 \cdot 2 \cdot 3 \cdot 2} = 1.6 \times 10^9$

(3) $\frac{(29)^5}{3 \cdot 2 \cdot 2 \cdot 3 \cdot 2} = 2.8 \times 10^5$

(4) $15 + 21 \rightarrow \frac{(36)^8}{2 \cdot 2 \cdot (3 \cdot 4) \cdot 2 \cdot 3 \cdot 4 \cdot 5} = \text{unusable}$ 4.9×10^8

(5) $\frac{(43)^7}{5 \cdot (4 \cdot 5) \cdot (5 \cdot 6) \cdot (3) \cdot (2)} = 7.55 \times 10^6$

These are PC values. They should be ~ no objects freed.

Dist 5:2
P 30: 50
31 12
32 6
≤ 68
P 33: 2 (# 69, 70)
P 34: 3 (# 71, 72, 73)
≤ 73

Discuss of Inst. times, ref: P 30; TOP

P 20 1 to 2 M steps/sec on 1.5 GHz CPU: somewhat so with clocks/step.

so 0.5 to 1 us / time step.

For (6), its not clear how to calculate! It may have been done as 2 completely separate ppm:

5 instructions: but c1 & c2 had been "boosted".

$(21 + 2.5)^5 = 23.5^5 = 7.17 \times 10^6$ which is rather small.

ppm (10) Tower of Hanoi, P 21: 10 mts $(21 mts)^{10} = 1.41 \times 10^{14}$

Comparison: 5×10^6 for OOPS! The disparity is due to the very success of boosted Boost

for this problem. Boost (2) uses 5 mts & returns the problem to (6) (how this was called "2" is unclear)

(Also note, in my 21 mts, several were never used better: - Day should be lower PC!) \rightarrow Also Not in (10) (TOT)! Several "New" mts were used, so they probably would be at very low PC (in my formulation)

GENERAL Discussion of OOPS!

OOPS was able to solve Tower of H. Because

- 1) It had set of instructions that gave a smart (to cost) soln to it.
- 2) By using boosting, it was able to get hyper PC's for 4 of 4 mts: for 4 of 4 mts, c2, c3, c4, c5
- 3) As I pointed out to ~~the~~ amount of mileage one gets (in CJS ratio) depends critically on the no. of mts used. w. 23 mts, he gets factor of 1000

Only 21 mts used for both problem types, ~~etc.~~ The ~~the~~ improvement ratio is ~ 1 for maybe 74 mts. It may be a 4th power effect, since 4 mts are involved.

so $\frac{23}{21} = (\frac{3}{2})^4 = \frac{81}{16} = 5$ so he'd just get 2x. CJS improvement is $\approx \div 5$

if he only used 21 mts.

4) His version of FORTH & his method of getting recursive of interpret: (The Using & DO loop is faster.) Its not clear that a pair of ~~the~~ greater PC. One Q is how it fits in with other functions, problems, hours. e.g. How PPM works there

5) Boosting could be a very useful idea: We look at each problem & we know & properly list for it.

The properties are correlated w. N. prob. solns. When a new prob occurs, we get its properties. Each property & was a use for all previous solns. We then use a wtd \leq of 2/1

& inst. used in all correlated problems. - Perspectives app: wts for PC of Tokens in new problem.

For more precision, we have a 552, make ~~corrections~~ corrections below

pairs of properties, no mts triplets. Actually somewhat work on how to do this \rightarrow "Opt. Dependancy-Trees" Scott Danes 1997

for each pair, triplet... of properties, we get a probability for each Token. \rightarrow 206.00

510 Do experiment in 26 Assembly mts. See how recursive can be implemented. Criticism of OOPS Problem solved: OOPS has no way to merge ideas benefits how warning on TOT under Perm 1424. Except: Re setting of pairs of tokens; Also important shortcut ppm each time an add-on best TOT for each token. So OOPS can't hold the variety of properties to describe problems.

Have copy I have in copy

only one of 3 factors listed in 202.29-32

20: 202.40 : In view of 202.32-39, it would be easy to drop 202.26-32 & use the simple mean of f. & PPM predictions of 202.38-39

Cover and King had a ~~idea~~ (I'm not sure) method of combining prediction systems - Perhaps look into it.

The PPM has a way of ~~letting~~ why. ~~with~~ context of dist. length. I had a method of doing it using a kernel in Σ space. The idea (perhaps incorrect) was that f.w.t. should depend on f. length of context, as on how much it's length differed from (h. long dist. ~~kernel~~) but ~~not~~ usual PPM may be w. this

stands on how many ~~distances~~ distant it was from: closest ~~to~~ context.

This may be close to the "To The Brunt" method of Burroughs & Wheeler.

It would be interesting to try 202.26-32 on simple sequential PPM as a method of combining context products. It ~~might~~ be v.g. & it might take less CC than the "kernel" method I had in mind.

While makes 202.26-32 not so by CC is that data for all contexts of same "overlap no." are pooled. In v. of methods of 202.38-39, it would be well (to try them individually, to see that each had about same product capacity. (Note: 2 reversals of 202.38 would seem to make any change in power capacity unlikely! 202.38, however, could give significant ~~different~~ ^(complex) ~~than~~ ^{than} ~~others~~.)

Another ~~4~~ ² types of context (in addition to 202.38-39) can be obtained if context is toward Root v.s. toward leaves: so this gives ~~5~~ ^{6?} context types that are easy to find & keep track of

Consequently, only 2 more contexts, since each node has ~~only~~ ^{only} 1 parent. (The I think we can still invert order of binary functions, non trivially, so maybe 4 more context types ~~∈~~ [∈] ~~8~~ ²⁷)

SN This ~~is~~ ^{may be} ~~an~~ ^{not} conventional "Phase 1": Go Back to 199.20 ~~to~~ ^{to} 26 I really need to work out exact details of it. → O.K. - Say we get to 199.32: we have 2 TSC → 205.00

FN Continuation of 203.20: Improved methods of PPM & then PSD using viz. . . . & my own tricks.

SN Gaily ~~is~~ ^{is} version of Conversion ~~from~~ ^{since} 1 = 0, or U will always be an output,

The Norman factors are always Exactly 1: {The output set is a Complete prefix set.}

SN In the "4 methods" of .17-.20: "Double counting" is not a problem! It is a feature! i.e. when the "double counting" is (I think) a legit example of alternate coding. e.g. when "sum" has its args counted. If it is really an alternate code, it should be counted.

A Counter Arg: That any self-limiting code can be written backwards & still be a legit. self-limiting code. Since this is true for All codes, considering C. backward versions leaves resultant PPs invariant.

SN In (decoding or designing) a code long: Look at ~~the~~ ^{successful} ~~code~~ ^{codes}, & prove in each lang. ~~the~~ ^{how} PPM ~~is~~ ^{is} able to find regys? Would PSD help? Longs to Consider: LISF, AZ, OPS (other forms of FORTH e.g. the lang I used in SMARTS for ANL) Assembly lang (26 mts). Fortran/BASIC/C.

8/18/04

205 1/2

TM

Usage of insts in OOPS

00

Problem 2: $2+3=2+4$
 " 3: $2+3=3+2$

$$\frac{1^3 \cdot 2^2}{2 \cdot 4^5} = \frac{2 \cdot 4^5}{4 \cdot 3 \cdot 10} = 2 \times 10^6$$

50

← pe's or solns for problems 2,3,4

prob 4: $4 \cdot 5 = 47$

$$\frac{4 \cdot 18}{4 \cdot (4 \cdot 5)^3 \cdot 2 \cdot (3 \cdot 4) \cdot 2 \cdot (2 \cdot 3) \cdot (2 \cdot 3 \cdot 4) \cdot 2} = 4.2 \times 10^{22}$$

Its clear those "Add on" solns would have not been obtained w/o it.

Special reference to previous soln.

Problem	No. insts used	Σ
1	2	2
2	5	7
3	10	17
4	18	35
5	25	60
6	5	65
8	1	66
9	7	73
10	10	83
	7	

8 add	30
9 mul	30
0 power	30
1 sub	30
12 div	30
13 inc	30
14 dec	30
15 by 2	30
16 by 2	30

55	c0	130
	c1	
60	c2	
	c3	
	c4	
	c5	

ONLY 20 Insts used (of 21 (opt insts))

at 73 2 insts are simple examples of how to write programs FORTH

Inst	Usage	First use	Count	Other	Notes
def up	1 2 3 4 5 6	10			
2 to D	2 3 3 4 4 5 5				
prt	2 3 4 5 6				
c2	2 3 4 5 6				
end up	2 3 4 5 6			10	
boost	3 4 5 6			10	
Do it	3 3 4 4 5 5				
bst	3 4 4 5 5 5				
from D	4 4 5 5 5 5				
del b	4 4 4 5 5 5				
by 2	4 5 5 9 9				
cpnb	5				
c1	6				
call p	6 10 10				No Problem #7 exists.
mv disk	8				8 Row 10
c4	9 10				Tower of Hanoi
c3	9, 9 10 10				
cpn	9				
EXEC	9				
dec	10				
c5	10				

boosted for prog 6 (thru 6)

boosted for prog 10

4TM

①

Scir vs. com / Joseph

Continuous and Convergence Rows

Graph

corpus

00:20:18: HA! This Partial ordering Net can be used to summarize data: just as linear ordering did in PPM! We can devise a pgm (like "2.3 tree") to put how nodes in the Graph!

Not so fast! ☹: Each node has a region around it! ~~There~~ for a given fixed node,

→ common regions w.r.t. each other nodes, but partially ordered.

04

Given a node N_i , every node N_j is partially ordered w.r.t. i (including $i \dots i$ has its entire graph Nodes in common w. itself). So we have relation $N_j \geq N_i$ (\geq means ordering w.r.t. Node i 's "influence")

for all $i, N_i \geq N_k$ if $k \neq i$. So this partial ordering is a binary relation.

But, if $N_j \geq N_k$, this tells us (little nothing about $N_j \geq N_l$ for $l \neq i$!

e.g. say $N_j \geq N_k$ has nothing in common w. N_i (i.e. they have 2 different tokens at 1. node from N_i) Yet $N_j \geq N_k$ could have a common token at their nodes.

0

So, to represent all of i 's info for all i 's, we have to place it on each graph of every node in the corpus: i.e. each Node has to have its own graph.

Actually, we only have to describe into graphs of nodes that have same token - other info being trivially uninteresting.

Also, we are only interested in the "down" (toward leaves) direction from 1. Node.

14

while 04-14 refer to \geq types of p. ordering & assoc graphs. There is no p. order graph: i.e. to graphs corresponding to tokens in the corpus.

A poss. way to use the info about contexts. Say we have a $(\frac{N_i}{\text{node}})$ in its assoc.

(toward leaves) context. (or, we can generate token toward leaves, so context is toward Root)

6

Any way, we have this Node's we want to know P.D. of nodes just below it.

Each context from 1. node of the corpus has a certain amb. of overlap to N_i (it may be zero)

22

The overlap can be the no. of Nodes that are ① sources that are N_i context ② are connected to N_i via "same" nodes.

Each context has a sub of tokens that it has produced at its Node point.

Each context has an overlap no. w.r.t. N_i . Assoc. w. each "overlap no." is a

26

To get a prediction at Node N_i : we take a wtd sum of p.d.'s induced by each of the contexts. The wt. of a context will be the product of 3 factors:

29

① Some constant w.r.t. function of "overlap no." (22) ← "Overlap Number"

30

② A Laplace or Dirichlet's rule involving the case count of each token w.r.t. that context (includes "the corpus")

32

③ α^{+N} : No. of tokens (including, perhaps, the corpus) and α is the (geometric)

mean pc per symbol of the corpus up to now.

The product of ② & ③ is the sum of pc saved by using that context (saved w.r.t. ③) ... I'm not

quite sure of this ... we may simply know to find α by "correcting out fast".

N.B. we can pool all the bottom contexts of Ξ (Overlap No.)

Sum of bottom.

At present, I don't know a practical way to get the overlap nos. of all of the contexts

relevant to a particular prediction at a Node.

I can think of 4 easily obtained context types, here; They are:

38

1) PPM with normal RPN. 2) PPM, but RPN corpus is von Neumanns.

9

3) ① but order of binary functions in counts is reversed 4) ② but order of binary functions in counts is reversed → 203.17 for 4 more context types → 203.40

4TH

0: 200.10! By doing PPM for RPN functions — both forward & backward, one would get somewhat

different branchings

So forward its m) α (phonetic) order. α A B C D E

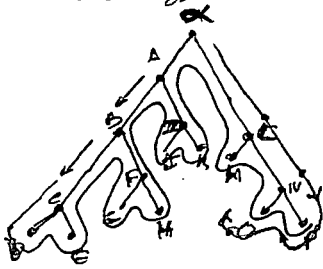
Backward its α RPN L N P O M A I K J B F H G C E D

So its quite different from backward of α forward RPN

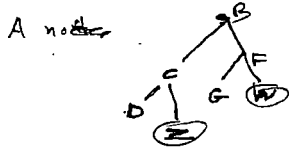
So it looks like using both forward & backward decs will no. of produs four
It will no. of a apparently similar nodes

A different approach:

Consider node B, we can "match" it w. all



in corpus to some extent. Each B node in corpus will match



will match up to \forall p's α α W. So for every B node in forest we can say how far out we can go in any branch & it'll have matches.

So assoc w. any α B nodes there will be a common match region, consistency & Branch, but

usually parts left out. These common Branch will be a function of α parts not

expanded. In α above case we have a function at α (E), (H) positions

Suppose we had a way to group, relate, each subtree w. every other subtree that had same α for

How could we use this structure/relations to assign p's to nodes?

Start at top node. This constitutes a "PARTIAL ORDERING" of the contents of any nodes into corpus. AS such; it can be represented by a Boolean Graph. (Also Note equivalence by commutativity (or possibly symmetry) of functions. \rightarrow 202.00)

Groups! α discussion of α α to key to wrong! The modulus I have desc'd is not really α .

discussing backwards, but inverting α order of (binary) functions! \leftarrow which is certainly an interesting α on which to use PPM!

But tree backwards would seem to have less good contexts (i.e. good for produs).

P O N M L K J I H G F E D C B A α 130

SN On trying existing PPM forms on various corp: I find it standard pems take α α any kind of file, expressed as a sequence of Bytes, — So I should be able to fast various ideas w. \pm 4 bits precision. E.g. I could try it on ~~the~~ Backward English, Backward Lisp, etc. I probably want to remove formatting instructions — like line feeds, etc. Perhaps represent each Lisp instruction as α a different 1 byte string.

In looking for matches, a search could take advantage of Commutativity of certain functions. (Also associativity, distrib. etc... perhaps).

N.B. The Magnitude of speedup of such is Enormous! Going from raw 8 bits/token to 4 bits/token (of L.Z.) is square root of CJS. α 8 bits to 2 bits/token (i.e. "PPM") is 4th root. For search of 10^{10} words; $\rightarrow 10^4 \rightarrow 10^{2.5}$ ~~reduction~~

Even a factor of 10^{10} (Jacobian COPS) does involve some compression... just how much? Also 2 bit compression involves long corpus (which α we don't have). α raw BZip2 I could start w. pure roots search for COPS $(\frac{1}{23})^{10}$? factor of 3 is $(\frac{1}{23})^{3.3} = (1.6 \times 10^6)^{-1}$ $(23)^5$ (factor of 3) = (2.07×10^9) which is α what COPS does. so we get $75^{1.66} = 1.24 \times 10^3$

4TM

00:199.40 That we can go very far. The system will learn to do induction in this particular way.

by concatenating useful heuristics. It will discover ways to express many (perhaps all) regys in this form --- so PPM would be e.g. update system,

To work hard problems, it would have to have an adequate set of concs. of this type. TM could develop a set of MS sets but it would perhaps (probably?) involve a set of heuristics different from those used by f. Thamer (?)

For heuristics that Thamer is conscious of: Could they always be put into form of concat. (Cases/Means)? Perhaps best way to understand this Q is to write various TSC's & see what heuristics they have & see if I can always put them in form of concat of concs.

My guess is that it can always be done: That a "most general" form is to be found of compositions of functions (primitive Rec funcs perhaps).

Can one always express any func of k vars., so that it is in r. form: $f(x_1, x_2, \dots, x_k)$: i.e. k variable args all follow + function defn.

Well we can have a function expressing that does just that: It looks at a function tree w. args. It has a way of referring to any no. of nodes of this tree. Say it refers to k different nodes:

If T_n is a tree of k original function trees & n refers to k of the nodes of T_n ,

then $h(n)$ T_n defines a function of k vars. --- i.e. that tree, w. those k vars replacing the k nodes of T_n . Finding h given n would also mean many funcs/identifiers: $h(n)$ would also mean many funcs/identifiers: $h(n)$

So looking at a larger corpus of function trees & finding common sub-trees is about the same as finding $h(n)$ for T_n (15, 16)

Is it problem of 15-19 + same as finding CFG Grammars for a corpus of function trees?

So there are 2 interesting kinds of "functions": Omega's $h(n)$ (15-16) & Rec others + do loop

For h takes an n (integer) & a function & a function as inputs & crosses a function $F(n) = \alpha$ for $n=1$ to n $F(n) = g(F(n-1), n)$: output is $F(n)$. (See 196 on h clause)

There is a function of 2 vars (i.e. $g(x, y)$) --- g could be regarded as h clause.

So, given a 2 input/output function, a real α & an integer, we get $F(n)$.

T. foggy suggests that it may be poss. to get a universal induction using PPM & a "common sub-tree search" & some kind of PEG-like Alg. (See 196 on h clause) $h(n)$ would also mean many funcs/identifiers: $h(n)$

On the other hand, it would be well to continue on discussion Goal of 199.20 & 2.6

An alternative more intuitive (?) way to get universality (& recursive).

ATT! Use GA to find Algs for fast detection of common sub-trees! First do it by hand to get good set of "primitive" operations. --- Perhaps it is reform can ROOTSOP itself in this problem! Perhaps use a (WEB) computer to constantly work on this problem.

This problem may have had others work on it. I think it known to be "NP complete", but still short approxims may be available.

If there were many people working on a phase 1 or even phase 2, "common sub-trees" is a major sub-problem & it will be poss. to get some fraction of the community to work on it.