

10 : In a Phase 2 ATM, we would be modeling $h \approx h'$ (194.30 : PP 16-19 of ^{my}IDS IA Report)
 01 ~~often very close~~ found by PPM, it's modulus, ≈ 198.30 (166.11-90)
 Actually, in phase 1 we would be doing $100 \approx 01$ as well.

0 : 194.17 : T. Monte Carlo ~~so even~~ does not at all work well for inv probs w. only 1 soln.
 Hrr. how it works for normal OZ probs is unclear! It is "sort of" used in GA, ^{where} it clearly
 is much better than pure Monte Carlo!
 How GA works & use M.C. to get a subset of cards of $>$ average G. Using P/M/Cross, derive
 2 - new P/D that fits best subset of cards & loop to α .
 Actually, the only difference (w.o. Lsich) of my method v.s. conventional GA, is that I use $\epsilon \leq \epsilon \geq \geq 2$ as
 priority to use ^{much} better models of 1. GPD.

T. ~~algorithm~~ Magnitude (i.e. other effects) of Lsich relative to McCork, is not clear!
 M Carlo ~~also~~ has more variance, so it would ~~be~~ more able to find "crazy far out solns" - ^{degree of goodness}
 but using better GPD models ~~may~~ ^{may} be better. In any case, the "elicited" effect of ^{badness of}
 Lsich is not clear

0 : : On Hyper level / rmg: ≈ 2 TM2! Say we are in Phase 1 we start with TSCQ:
 26 The ~~regg~~ needed is not findable by PPM. What we do is Go to hyper level: We use TM1 to
 use 2 complete (unw.) lmg. to find reggs needed to inv. TSCQ. P.D. hyper levels can have
 2 choice betw. several update models. One is PPM, another imposed ~~the~~ PPM...
 2 report!
~~use of~~ use of ~~various~~ ^{various} PSGD. or other: Go over this More carefully: more exact statement
 of what to do - A picture (in "English") of what I expect Phase 1 ... ~~(TM)~~ to do.

1 Simple case: TSCQ of QA problems: TSCQ is "slow" in the sense of not many new concs.
 introduced. Also ~~not~~ solns. have short solns in ~~ref.~~ (avg. GPD = ref. p.d. Ordering of
 probs in TSCQ is irrelevant to speed of soln. Essentially, we have this list of rules & their
 20 pe's & we do an Lsich to solve all probs in TSCQ. There is no modification GPD!
 no "hyper level lmg.". Static GPD: Perhaps based on a prop of Tokens only.

2 As above but GPD is updated, say using PPM. PPM may need large ϵ 's - so one way $\rightarrow 20500$
 of something reg is to report problems. In each QA problem we can give TM ~~some~~ a number
 that tells how many times it was reported. This gives Signif. Mod. for PPM updates.
 We are still limited in pass). TSCQ's fast system could be to understand/predict.

Theo it may be that very many reggs are expressible as a simple concs of ~~tokens~~.
 That once ≈ 2 v.g. set of users is discovered, & one has probs for ~~various~~ various concs of ~~tokens~~ (20000

4TH

10:196.40 : On status of Recursive funcs: If f. lang has Do loops (I don't know if while or 'until' are used) it can do all perm. rec. funcs: - which is ^{probably} adequate.

How, recursive funcs as such (via recursion delay) are (probably) necessary to get any PC for certain funcs.

- So I'd like a way to define Recur. - Presumably ~~the~~ ^{the} compiler creates a "Do loop" ~~to implement~~ ^{to implement}.

Kleene's $\forall a (p(a)) \equiv$ my Page 196 $\frac{1}{2}$: $\varphi(0) = q$; $\varphi(y+1) = \chi(y, \varphi(y))$

is a simple way to define most prim. rec. funcs.

{ I had ~~known~~ ^{known} $\varphi(0) = q$ $\varphi(f(y)) = \chi(y, \varphi(y))$ ~~known~~ ^{known} $f(y) \equiv y+1$.

109 One (common) practical way to get recursive solns: say we notice that a is b and both acceptable objects; Rec. to verify simple relation r, to a (i.e. $b = r(a)$). We then hypothesize that ~~the~~ ^{set} $r^n(a)$ will be "acceptable".

Unfortunately, recognizing that $b = r(a)$ may be unlikely, since b is a may not be in proper form - for this reason, it is best to keep ~~some~~ ^{several} solns. to each problem in memory. for both $0 \leq i$ Inv. probs, one should "over search" for solns, if at all possibl. This makes it more likely that a soln to $b = r(a)$ could be found.

115 Whatever (Formalisms) ~~are~~ used to define recursive functions, I suspect that they will usually be of very low PC a/o by CJS. (Unless found by .09-15... which is a different kind of heuristic solution) I'm thinking now of simply generative ~~various~~ ^{various} kinds in "Guiding Pd" order - Recursively defined ends will usually have much longer delays than non-recursively defined ~~various~~ ^{various} kinds. This may also be true of "Do loop" defined functions. So it may well be that recursive "Do loop" defined functions will not be found until we get around to very CJS problems. - But even for humans, small solutions of Rec's sort ~~are~~ ^{are} of very large CJS.

20 Probably necessary of defining w. recursion is via "English" § 194, 24-27; 19610-13; 20th § If I think that recursion is a good solution and, then I can b. y. recursion in English. Whatever way we decide to implement. recursion, the English statement of it should guide us in evaluating it PC.

Re: Recursion in COPS! COPS took maximum out of time code Tower of Han w. $n = 10$ or 20 or 30 - In Lsrb, it must have tried other ~~much~~ ^{more} improbable plans for shorter times... See also Rec's "Scus" - Did I really do "Lsrb"?

30. Yes, Even if we stick to English, we do have to have at least 2 ^{non-human} languages: One for ~~the~~ ^{the} lang is part of the description of problem (which can be close to English), One for description of solution. - "problem desc" The choice of Rec's second can be imp. we can try out various langs to see which ~~seems~~ ^{seems} to work best. Take a look at the COPS language. Just how Recursion was implemented. Was it of all efficient? It may have been a bad problem! Even w. correct solns, it too enormous and of time to solve. (for large n ... for small n, not so large).

amine this in more detail, let us write the pair of equations

$$\begin{cases} \varphi(0) = q, \\ \varphi(y') = \chi(y, \varphi(y)), \end{cases}$$

ss the definition of a function $\varphi(y)$ by induction on y , where q a natural number, and $\chi(y, z)$ is a given number-theoretic function variables.

for example, the value $\varphi(4)$ is determined thus. To generate 4, rate successively 0, 1, 2, 3, 4. By the first equation, the value 1 be the given number q ; then by the second equation, the value 1 be $\chi(0, \varphi(0))$, i.e. (using the value $\varphi(0)$ already found) $\chi(0, q)$, ince $\chi(y, z)$ is a given function) is a given number; again the value 1 be $\chi(1, \varphi(1))$, i.e. $\chi(1, \chi(0, q))$; the value $\varphi(3)$ shall be $\chi(2, \varphi(2))$, $\chi(1, \chi(0, q))$; and finally the value $\varphi(4)$ shall be $\chi(3, \varphi(3))$, i.e. $\chi(1, \chi(0, q))$.

ve have a process by which, to each natural number y , on the the generation of y in the natural number sequence, a corre- number $\varphi(y)$ is determined. Since a number $\varphi(y)$ is thus associat-, for each y , a particular number-theoretic function φ is defined se numbers $\varphi(y)$ as its respective values.

unction φ satisfies the equations (1), when (1) are considered as ul equations in an unknown function φ , since every particular comprised in (1) (namely, $\varphi(0)=q, \varphi(0')=\chi(0, \varphi(0)), \varphi(1')=, \dots$) is satisfied in the course of selecting the successive num- b, $\varphi(1), \varphi(2), \dots$. Also this φ is the only function satisfying unctional equations, since the process by which we determined the e numbers $\varphi(0), \varphi(1), \varphi(2), \dots$ from the equations (1) can be ed as showing that any function φ satisfying the equations must values selected.

er definitions by induction, the function φ defined depends on 1 variables x_2, \dots, x_n , called *parameters*, which have fixed values ut the induction on y .

LE 1. Consider intuitively the equations

$$\begin{cases} a+0 = a, \\ a+b' = (a+b)', \end{cases}$$

encountered in the formal symbolism as Axioms 18 and 19. ine the function $a+b$ by induction on b , with a as parameter, a previously known function. Then the equations

$$\begin{cases} a \cdot 0 = 0, \\ a \cdot b' = (a \cdot b) + a \end{cases}$$

KLEENE: Meta-mathematics

define $a \cdot b$ by induction on b , with $a+b$ as a known function; and

$$\begin{cases} a^0 = 1, \\ a^{b'} = a^b \cdot a \end{cases}$$

define a^b by induction on b , with $a \cdot b$ as a known function.

An example of a definition of a predicate by induction will be given later (Example 2 § 45).

What number-theoretic functions are definable by induction? To make this question precise, we must specify what functions are to be taken as known initially, and what operations, including what forms of definition by induction, are to be allowed in defining further functions.

We shall now select the specifications with a view to obtaining functions definable by induction in an elementary manner. These functions will be called 'primitive recursive'.

Each of the following equations and systems of equations (I)—(V) defines a number-theoretic function φ , when n and m are positive integers, i is an integer such that $1 \leq i \leq n$, q is a natural number, and $\psi, \chi_1, \dots, \chi_m, \chi$ are given number-theoretic functions of the indicated numbers of variables.

(I) $\varphi(x) = x'$.

(II) $\varphi(x_1, \dots, x_n) = q$.

(III) $\varphi(x_1, \dots, x_n) = x_i$.

(IV) $\varphi(x_1, \dots, x_n) = \psi(\chi_1(x_1, \dots, x_n), \dots, \chi_m(x_1, \dots, x_n))$.

(Va) $\begin{cases} \varphi(0) = q, \\ \varphi(y') = \chi(y, \varphi(y)). \end{cases}$

(Vb) $\begin{cases} \varphi(0, x_2, \dots, x_n) = \psi(x_2, \dots, x_n), \\ \varphi(y', x_2, \dots, x_n) = \chi(y, \varphi(y, x_2, \dots, x_n), x_2, \dots, x_n). \end{cases}$

((Va) constitutes the case of (V) for $n = 1$, and (Vb) for $n > 1$.)

A function is primitive recursive, if it is definable by a series of applications of these five operations of definition.

This definition can be given in more detail, analogously to the definition of provable formula for the formal system (§ 19), say using the second version, as follows.

We refer to the above equations and equation pairs (I)—(V) as schemata. They are analogous to the postulates, with (I)—(III) in the role of axiom schemata (or more strictly, (I) to a particular axiom), and (IV) and (V) in the role of rules of inference.

A function φ is called an *initial function*, if φ satisfies Equation (I),

"BASIS A"

induction

Handwritten notes and diagrams:

- Diagram showing $x' = x + 1$ and $q = 0$ with arrows pointing to (I) and (II).
- Text: "isn't III a special case of IV? maybe not!"
- Text: "I don't see why '0' is necessary for counter-ex. Sp..."
- Text: "shouldn't $q=1$ by (I): No!"
- Text: "It defines φ for $0 \leq y < \infty$ if χ is defined for anal. argts."
- Text: "223 is more relevant, but still, have $q=0$ in case 50."
- Text: "composition"

8/11/04

1965

4TM

It may be poss. to continuouse (vs to prim. rec. functs) & avoid that recursive defn problem!

— See 193.10 on Paths (in Klamro)

I vaguely remember McCarthy defining
$$\begin{cases} \text{Fact}(X) = 1 & ; X=1 \text{ else} \\ \text{Fact}(X) = X \text{Fact}(X-1) \end{cases}$$

Use 2 IF statements. Def Fact(X): IF X < 1 Then 1 else X * Fact(X-1).

This would do some simple rec. functs — more complex functions like Ackermann — maybe more diff. In particular ... would one be able to "unroll them" (i.e. not use f. stack?).

0:194.27 : This seems like a very imp. idea! I should keep down lang. in "English" as long as

possible because it will be better to get (user sample of kinds of abs. D.2) it will need to in order to get. Redundance (lang/unc). Also use several English translational solutions in ||, since I want to be able to get as much "Mileage" from the Corpus as poss. (Understand how to get what I want from the corpus). → 20

SN On NMTM: For real output, a distribution can be characterized by μ and σ^2 . For more exactness give more "Moments" or other params of t.D.P. Or, for Multimodal dist. — have several μ, σ².

For discrete dists — One common approach: ~~Abstract~~ At "center of cluster" w. and/or finding how far other choices are from the center. For more complex dists — have > 1 cluster, possibly different ways to express distance from center" —

One way, when ^{one} has a "not so happy" soln. — one retains it because it might be useful for "backtracking" later.

On Prim. Rec. Funct.: P119 of Klamro pretty well characterizes Prim. rec. func. — They are definable by simple functional operations. The most complex involves a Do loop in which the upper limits defined by p.r. funct.

I don't know if "until" loops can ~~be~~ be used to define ~~the~~ p.r. funct. I suspect that they can but I haven't yet proved it.

We can define rational nos. as integer pairs.

P119 of Klamro suggests us to define a func for all integers ≥ 0 ("natural nos").

Then Def $\rightarrow \alpha \text{ BR } \beta$ just defines completed branches, not "functions" (ordinarily).

In branch generation ~~BR~~ "BR" ^{is} a NT in P&S generated P&MS. In generation

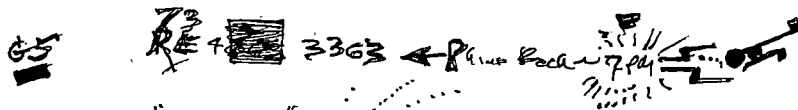
of 1. branch from $\alpha \text{ BR } \beta$, if any Branch occurs, we can write "X_i" for it.

The compiler regards P&S as an expr. & keeps track of t. "strip" of t. P&S being generated. If $z = 3$, say, the next expr will have $w = p, t$ subscripts (P&S)

We may want to modify P&S d.f., from $z = 4, w, p, t = 5$ $z = 1, 2, 3$ w.p.c's $\frac{1}{6}$.

This looks pretty much like the AZ language described in ID&M Report.

4TM



On Minsky's Chapter 6:

Progress in Statistical Learning: Theoretical, practical.

1) PPM is its improvements:

~ better than 2 bits/spatial approaches 1.3 bits/spatial of Hutter's.
 ~ 2.25 bits for LZ; ~ 2.25 bits for BZIP2 - < 2 for improvements
 → 1.3 for human
 (also some techniques seem promising for getting much closer.)

2) Fogel's checker player: feedback from win or lose only using ANN and hill climbing. Feedback from win/lose only: A neat way to deal w. "credit assignment problem"

3) KOZ's GP using Lisp prims: Designing of Opamps, Invention of Patterned

filters Spectrum filters: New common types of filters, - Also new patents

4) Mech Translation - Statistical } Long understanding as a kind of MT problem... Common Human lang + in formal lang.

5) Recent improvements in GA sorts more like regular statistics - No population } Busson Modem GA!

Passy of LZsch - ... not yet used. "Cracking of PEG-discovery problem?"

7) Ref good, statistical lang is a adequate soln to General A.I. problem:

It encodes all the OZ's: INV probs. Perhaps foggy methods are adequate for TM to begin lang. early into.

100 possible moves for Checkers. fairly good game playable over "lookahead" - i.e. a best move for given position

Main traps that the statistical prog haven't been able to go to hyperorder abstractions...

Also, RealTime has been much news about what they can do: Much less about what they cannot do.

SN How. A vital point: Prob's ^{like PPM} update algm need not be Turing complete... it can be relatively rudimentary... ^{like PPM} (over Lemp. Zip!). That it is important that the Ref lang be universal. Perhaps this is equiv. to TM₁ = TM₂ - i.e. we get to update problem to be worked as a "regular" TM problem... having as "co-args, the rest of TSC: scores"

SN In LISP, you can do partial evalns as soon as partially completed "S expressions".

eg. $sum \equiv \#mul\ 0\ 3$ doesn't have to record args for "sum", but we can evaluate $mul\ 0\ 3$ & perhaps put it on stack. We then have to miss feature of Lisp asking for more inputs" so we can use it as "3 1 0", to assign pcs!

can recursive function Lisp have short-cuts, but take long time to evaluate if stack is used. We will always try to complete such expressions, so they take much less time (eg. expressions that as far as next loops. As far to mechanism in it in working the "constant" in the definition [Ifn $factorial(x) = factorial(x-1) * x$; $factorial(1) = 1$ - y. "factorial(1) = 1" is a "constant part" the "boundary cond" ← a better way of looking at it.] Perhaps we can arrange to put after giving to important recursive data, the boundary cond. depends very heavily. eg after writing $fact(x) = fact(x-1) * x$, we submachine ~~write~~ write $fact(x)$ & we have to complete the expression as to the pc of "fact" set to zero.

4 TM

TSQ's : Searching for Roots : 14

Mr, in this last case, the A.H. hypothesis has to be very A.H. : very visibly so! — As a result, it is ~~is~~ much less dependent on a principle . So this is a good feature of Cor. val.

[SN] In SAARB paper for ANL! Not solved! I could do facilities for Definition!

Def $\rightarrow \alpha BR \beta$ \leftarrow CRITICISM (196.30)
 \uparrow \downarrow
 At \leftarrow α is bare bones Nat def. in a defn. BR is a Branch \rightarrow Non-terminal
 That represents a complete branch. So defns can be for constants or Forms. After a defn. has been made, it is given a name automatically & its PC is obtained by the Dirichlet rule (Ganzel & Lutz rules, maybe \exists connective rule ("continuum of inductive Methods"). A (perhaps) way to do recursion. A defn. can only refer to past defns. So we first define $Fact(1) = 1$ then $Fact(n) = n \cdot Fact(n-1) \leftarrow$ No! #
 See how Juergen did defns, & recursion.

I could do recursive defns then "in role" from n is computer to $\downarrow \infty$.
 Look up Peters' work on Prim rec. functs. I think she had a simple, useful () \rightarrow 196.25

Way to characterize Nam. — (Not so easy to find, but. — (not in K(ance).)
 Would it be ^(possibly) to have BR in .02 be able to refer to a branch that has already been ~~generated~~ / generated?

4 Seems like empty ideas!
 In devising a good "closeness Criterion" ~~is~~ better conds:
 (1) In writing TSQ's, I will have various heuristics in mind by which previous experience is brought to bear on present problem. Use M's for "closeness Criterion"

16 (2) How can I (file/store/access/recall) relevant "close" instances in past? —
 Would work on "Case Based reasoning" Be useful here? They may have derived some useful policy systems. Perhaps Ask Alice for rules — preferably on web.

Re: (2) (1.16) An idea: First use PPM type look order to get cases w. some "near context"
 then take that set & somehow look at next ~~closest~~ closest context, loop to (2.20)
 We may or may not have to re-order (re-sort) after each occurrence — but to stress eff. population to be restarted & rapidly

(3) I'd like to make it possible for TM to discover (better indexing methods, better "closeness Criteria") — these may have to be done simultly.

2: A disturbing possibility: that for rather simple tasks (in Algebra or even in non-A.H.), logical reasoning is needed. — That w.o. it, SSZ has to do very large to get any induction at all.
 — Otherwise, all ^(possibilities) choices are of ~~the~~ \rightarrow PC. Well, so then in certain areas of induction, we will do poorly w.o. deductive reasoning. In General, there will probably always be areas in which a given TM will do poorly. — whatever a priori or search P.D. we have, there are probably regions in function space that are \downarrow will remain, very distant. — balances suitable TSQ's are provided to suitable bias applied to search.

4-TM

SN Admin. Note: If I'm in good creative mood. (usually detected by good handwriting (like now)) then I work on Serious TM problems. If not, do Reviews. The logic is: Reviews need to be done, but should require less cleverness. Also, if I'm in a creative mood, I couldn't get very far into a review before being side-tracked!

Also form meta (Jung's version) makes it easy to implement PD's via a 3rd func. = "3.10"

SN An impt decision is whether what ref lang. to use. Present choices: Lisp v.s. Forth. Lisp seems to have nice tree structures for functions, but Forth is much faster, & enables one to do partial pyms before switching to another branch of it. Lisp tree (This last gives a factor of "tree depth" in speed over "Lisp" ... I assume a Lisp ^{PSM} has to be completely specified before one can begin to run it ... (But this may not be necessarily true!)

Jung's version of Forth seems v.g. (Pro has not tried to get it to be used at all! ~ 1000 machines insts. for Forth inst!)

T. Frazer I'm interested in speed, is Ref it enables me to write T&Q's w. larger CJS's ... which is much easier than writing them w. small CJS's. There may be a dubious "Economy". The main problem seems to be (1) writing T&Q's (2) designing langs & assoc update Algms for ~~the~~ T&Q's

Another (perhaps) big problem is update Algms for Phases. Different updates give various (speed of update) at various cc's for updates. eg (24.33 calls how frequently to use a particular update Algms. That eg. assumes Ref & update Algms can be "Mixed" (i.e. use one then use another later. Update Hvr, update Algms can do various things. (1) They can find new solns to problem ^{old} T&Q. (2) They can recode ~~old~~ old solns to problems (so solns are f. same) — It is always f. of cones in the Search PD for L&R, but however it.

evaln. PD invariant.

Perhaps a very impt point: When I get to **T&Q's** ^{f. "final problem"} I get into trouble, I want to be aware of the various options possl. in all different aspects of 4-TM project, Ref could lead to a set of diff.

One possy would be to start on a **T&Q Now** & see how much it messes diff. parts are!

165.12 : Why ~~we~~ would we ever want to look at v. Data & use it to control such, by path selection?

One Big Reason: By looking at data, we enable our Biologically derived system to utilize a pri info Ref we are not consciously aware of. This can be a Big Deal! i.e. a very important heuristic. Is there a way that we can enable our heuristic & networks still have the advantages of v. U.D.?

Not so easy! Since we are modifying the retentions a priori, after we see the data ... !

Hvr, this "modified A priori" ~~is~~ need not be v.g. — it can be very bad!

Cross Validation is designed to deal w. this Looking at the data a priori. It may not be the Best way, hvr.

Also it is often diff. to apply to many kinds of induction.

Also Re: Cross Val: Consider the "leave one data pt out" method. Even it can be made to have spurious reports if we make a soft A.H. hypth! A soft A.H. Hypth can give us any result we like!

Proof that $T \Leftarrow 3T$ is Best for Lurch -20
(rather than $T \Leftarrow 2T$, say)

Work GP/GA problems.

What if "state of Art" is now, is unclear: 188-32 show (2003 Dec) may be up to date.

Stolker (188-30) has web site w. ~~the~~ applic. of his ideas (on PSC-discy, among others)

to various problems: I would be well to see what progress has made since 1994 (10yrs!).

Stolker: Now at SRI in Ind Comp Sci Ctr. Berkeley Calif.

2002 SRI L.M.: an extensible lang modeling toolkit.

Stol... 95 ... computing principles

98 Entropy Based Theory of "Backoff lang Models".

"BOOGIE" was sw. in Lisp/CLOS to experiment w. Probabilistic Grammars & lang. Atoms.

[SN] It may be v.g. to try to Model a "G" function of a set of strings by a PC-

function of the set of strings. I wrote a lot about this, but I think ran into

a dead end because of an error in my thinking about lang. Models.

This recent confusion PSC discy may be what is needed!

Along w. the G function I will want to optimize (at least (param model) of a

G \rightarrow PC function like $PC \propto G^{\lambda}$ or $PC = A G^{\lambda}$ ($A = \text{Normal const}$).

(Yrr, we are mainly interested in the Goodness of fit of G to PC, in v. hy

G-rapous.

In $(T \Leftarrow \frac{1}{2}T)$ Lurch; What is optimum value of α ? : There are 2 factors of ~~math~~ efficiency!

$\frac{1}{1-\frac{1}{2}}$ factor = $1 + \frac{1}{2} + \frac{1}{4} + \dots$: its the amt. of time wasted in repeating trials.

T. after factor α is the amt. of time ^{wasted} in trials of α final round: If we assume that

the probab. of success is uniform in log domain. Then α final round will be between α & 2α .

Success can occur at any point. T. expected value will be $\alpha \cdot \sqrt{2}$.

So what value of α gives min for $\frac{1}{1-\frac{1}{2}} \cdot \sqrt{2}$? = $\frac{2^{\frac{1}{2}} \cdot 2}{2-1} = \max(\frac{1}{2^{\frac{1}{2}}} - \frac{1}{2^{\frac{3}{2}}})$

$-\frac{1}{2 \cdot 2^{\frac{1}{2}}} + \frac{3}{2} \cdot \frac{1}{2^{\frac{3}{2}}} = 0$ ~~for~~ $1 = \frac{3}{2}$ so $\boxed{\alpha = 3}$

Just stopping time should be uniform PC in ln domain: Common situation, related to idea that "1" occurs most often as leading digit in lots of sizes of objects like lakes or towns, etc.

That is, def. should be uniform in ln domain is neg. if we expect rhvance w. multiplicative? change of scale (dimension ("stability"))

Actually a Grammar Grammar is a v.g. way to get a prys of Grammars.

It enables us to use statistical info on previously successful Grammars.

This is what I did in Sol G4b... Also J.S. Hornby. But Hornby used a mathematician's

definition of a soln. to the Gramm. discovery problem. He way he used heur

such also, - but not enuf!

The PSC is already correct, it is not "obvious".

~~Handwritten notes and scribbles on the right side of the page.~~

PSG Discovery 20 (very specific).

[SN] Gen. remarks ① Does Chou-Liu paper just get a better than "Greedy" soln., or does it get an optimum soln. ② My impression was they implied methods were reluctant to "move complex dependencies"

③ On idea of "Order of Choice" If order of choice is known, one can always just comp. data to get pc's for each token: perhaps by PPM. (In Norm's PPM, order of choice is known: its simply sequential choice). Anyway in generating functions, what are "reasonable dependencies"? How is that related to idea of "Sub-functions"?

Each App. of a function must be a complete "Branch". Each branch has (historical) "Context". This can be "external" context" or "Generalized context" (= ALP).

In funct. Generation, we may know order (or partial order) of choices, but we don't have clear idea as to which tokens in text, our present choice is miss controlled by.

HA! In ^{function} text generation, a Branch occupies "One Place", (in terms of distance from ~~present~~ ^{or Polish} token to be chosen.) So in {add Branch, x}: T. context of x is "add"; (by commutativity we could rewrite it as add, x, branch... which would be ~~clear~~ ^{wrong} more obvious: ... consider sub branch, k. ; x is ~~an~~ ^{the} argument of "sub branch"

0: : PSG-D A nice way to discover/construct heuristics: say a, b ^{wrong} concatenation (norm's set) "a". a & b have pc's within α . — So we have rules for producing α , then we use p_a & p_b to deriv a & b within α . Doing this gives a way to deriv corpus that is checker then using a & b directly, r.o.d.

we can tentatively add (or subtract) a gen from α & scan of PC's all over corpus PC A on α ; I fit A we add that gen to α . Similarly w. subtracting an gen from α .

We can ~~use~~ tentatively construction heuristic by looking at all words first ^(follow per coda) & confirm by drag word. This usage is often "framed" (i.e. 25-26, modification types) Also Note: The Merge operation & pc of generation: T. "chunk" & pc of model: we can add 1000 new characters of those 2 kinds.

① About Discovery of CFG's by Stolcke etc. After Stolcke's Notes 1999 he wrote applied Pascal ideas to MUCH stuff in "A.I." type problems (See his website for enormous no. of downloaded papers). Shan, McKay, Barker, Abbass, Essam, Nguyen AUSTRALIA Dec 2003 is review of applic. of published models to speech. — But they sat on Stolcke's Pascal for their model. These guys seem to be saying that Chou (98) & Stolcke (99) used a poorly defined Apsip function! — That point is much better!

My present impression: that Stolcke probably had a Apsip eval. func that may have been "good enough" to begin with. PSG discovery problem, & that he had 2 mutations per PC (maybe) that were a very good beginning for PSG-D. ← This was in (1994)
In (37) Shan et al 6 yrs later may have corrected Stolcke's eval. function & for tests use improved methods ~~PPM~~ ~~Model~~ ~~production~~ What they do instead now, is to use PSG-D to

ATM

0) Baluja 95 compares GA vs. other methods for 27 or prob. Other methods "Better".

An ordering of improvement in Prob Models for "GA".
1) PBIL ~~Baluja~~ Baluja ~~96~~ ⁹⁶ - extends scope of PBIL

2) MINIE De Bonet et al 1997 MIT: "second order pc"

3) ~~Baluja~~ ^{Davies} Baluja 1997: Improved PE of 2) This uses Chow, Liu (1968) for fast improved P.D. They say it's significant improvement.

4) Baluja, Davies 98 ^{Optimization problems} (Chow, Liu method of speeding up both PBIL & "Hill Climbing") - This is "epidemic" they don't use probabilistic Model! It's correct w. "Which Groups of Params are more important for GA" (which GA is not concerned w. - i.e. all crossovers (2nd & 3rd)).

This paper reviews PBIL from ATMIC. (Jan 98)

Chow, C | 1968
Liu, C
ST 14 462-467

4) Baluja, Davies 97a (or 1998?) Fast Probabilistic Models for Combinatorial Optim.

This uses COMIT: How this differs from Baluja 1997a (info 1) is not particular.

[P2 col 2 of My paper: 1st part] supposed to tell what COMIT does. P3 section 2 is 2nd element.

I think COMIT uses 3rd or possibly 2) [which are by CC methods but v.g.] alternatively w.

PBIL ~~CC~~ CC & fast search methods. (or alternatively w. a Hill climbing (fast) search.

3) is used to reinitialize PBIL (or Hill Climbing) every once in a while.

My formula for how often to use by cc update routines of ATMS 12.4.33 is relevant here: Using methods that give fast formula would enable one to tell how much time to allow for a slow search like (3) or (2) & how much cc for PBIL (or Hill climbing). The, at present, I don't quite understand why COMIT is a good idea (if, indeed, it is!).

5) Pateman Goldberg survey 2000 | 1999 was latest read. (v.11 pp: This is "literature" review. PS/7/26/04)
of work in this area. So refers to General work on "Building Blocks".

6) Baluja A priori Data for Prob. Models & Optim: 2002 26 pp This looks like a PS/7-29-04
more detailed review (2 more relevant than) 7-29-04
Gives list of survey papers: 18 31 ²⁰⁰² 39 45 41 ²⁰⁰¹
1996 2001 } Pat. Goldberg....

7) Pateman Goldberg Tsutsui 2002: Comparing Bayesian Optim & Adaptive Search.
I think this works problems w. smooth. decreasing, continuous Optim.

4TM

Re: may be case counts for corpus

00 : How, How do we generate counts? : T. ~~next~~ array tells how often each pair of digit values occurred in 1. corpus. How can we use this to generate counts?

02 One (slow) way: ~~to~~ (obtain) ~~next~~ array: pick element at random, or its total pc. On basis of that 1. choice, put ϕ 's into those elements that are now (imposs.). The single chosen element, "a" gives ϕ 's for all other elements: choose one on this basis. — "b"

Now, w. 2 chosen elements, use conditional dist. on elements a & b: sum the ϕ 's (put more ϕ 's in illegal bins). When "c" is chosen, we now have 3 comb. probab. on legal choices. We add them (if normalize) & choose "d", 4. next choice. We continue until all elements for each of 1. n positions is chosen.

07 100 ft: 1) is much slower than original 2) It is not clear that we get to write ϕ pd on counts this way — or even that it is better than original (slow) way.

0 A big Advantage of MIMIC over PPM, is that MIMIC can use correlations between non-adjacent elements.

Its Disadvantages over PPM: PPM can use long sequences for induction, w. high SSZ, a good indication. PPM can deal w. variable length strings. Mimic may be able to be modified to do this. I don't know how, yet. See .18-20

15 Some "improvements" of Mimic: ① use last 2 choices to get ϕ 's for third choice. Needs more memory, but otherwise, it can be designed w. close "analogy" to Mimic. Initial use of set & permutation of 2 positions. use of Max likelihood for choice of permutation. → 0.27

18 ② In stead of fixed vector of length n with fixed coordinates; ~~we~~ make coded system "moving along" as in t.s. produ. The "case no. matrix" those "n" x "n" elements.

1. as probab. may be as large as feasible. At (i, j, k) element gives case count for digit i following digit j by i spaces. How do we use Best Matrix for ① produ ~~to~~ generate counts? (This can (perhaps) be done using the ideas of (0.02-0.07) which is of uncertain accuracy. Also, since we have no ϕ 's in 1. count, we can generate as long as we like w.o. filling it up! (Unless we have a "stop" symbol). So we can continue generating ϕ 's until all of 1. column before "stop" symbol is filled in.

26 So we have a way of randomized length.

27 ③ We can extend .15 to 4. last K choices: we can vary k, depending on current SSZ. (SSZ can \downarrow as we close in on "minimum")

0 While PPM & MIMIC & its Modulas may do no so bad w. function Generator, what I really want is a method that readily recognizes sub functions

A ~~sub~~ subfunction is obtainable by deleting any complete branches from a count. But I'd also like to ~~know~~ make a funct. by taking any branch & deleting ~~the~~ sub-branches from it

See if there have been recent developments in MIMIC 37 . 184.11 12 1998 ref. "improves" Mimic" (By Boloje)

47M

(1997)
R. Szustrowicz
in J. Sankh.
in q. vol. 13
Inc. p. 84
"PIPE"

2 11 17

[SN] RE: Juergen's paper! I. Evans he fixes the PC's at 4. nodes: How does a run of trials w. fixed PC's. The best could be selected from a run & used to modify the PD's at 4. nodes - (Hrr. check this for correctness) } See [low pass analysis] in STM 9.11-ff
He also uses several update methods.

[SN] Definition of functions: Several approaches!! Koza: separate random generated functions!!
(incoherent) (or exactly?) (Approximately) to same!! Soli: No definition made until evidence that it is useful...
Takes more cc, but they have more pc. (?)

.04 may be one of most critical Q's in AI!
Also consider how Baluja & others deal w. these PC's!
Also T. idea of finding common substrings in DNA, Proteins (M. Li; Garry Wolf).

Atta "hyper level": I was assuming that for more complex variants would be "expressed" "noticed" at a "hyper level" - To what extent is this true? Exactly how does it work? How expensive (in cc & ops) is it? What is best way to balance cc & (z pc) of discoveries into a "lower" vs. "hyper" levels?

[SU] (Re) Notable PPM is as ~~least~~ least as good as Koza branch sampling, because PPM does always recognize complete & partial branches.
Hrr. "Automatically defined functions" (.04) By Koza seems like way to recognize functions that I implement in .054 - at more or less c/c/pc. T. Q is gain or loss of cc/pc.

Comments on MIMIC 1997 by Bonal, Iszball, Viola.

I assume the model is a k component vector w. index r.
In this case, there are r^2 components of P(x_i | x_j). (i, j = 1 to r)

Mean we want to choose an ordering of the x_j such that $\hat{P}(x_{i_n}) \cdot \prod_{j=1}^{n-1} \hat{P}(x_{j_i} | x_{j_{j+1}}) = \max.$

The $\hat{P}(x_{i_n})$ is $\hat{P}(x_{i_n} | x_{j_{j+1}})$ are empirical values

So ~~many~~ we would need r^2 * n parameters to describe distribution. Hrr. we end up w. a model having only r * n params (still, not so small). Most ~~params~~ would be ≈ 0 (using straight rule) & $\frac{1}{200}$ using kops rule (200 is 1/22 of 2^n "iterations")

We don't want pc's to = 0 since we then can't reach certain values.

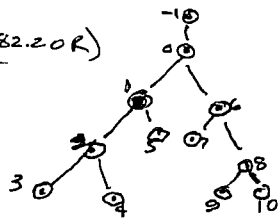
Woops! I've been mixing 2 models! T. model has n^2 params

A complete model would have $\sim (nr)^2$ params - ~~is a pretty~~ there are nr "states" - each given by n x values at each of n positions, so $\sim nr$ transition probabilities (from fact e. "n states" values need not be sequential).

By choosing a specific "order" (= parameter), we only need nr^2 instead of n^2 * r^2 params. While this RAM heads, it doesn't help w. ~~is~~ there is no "pooling", so it might be better to store all n^2 params. Hrr., updating takes ~~times~~ as long also - (but we don't update ~~after~~ each cond.)

4TM

o: : "Repn" of (182.20R)



Look at node 5: while in ^{regular} Polish, 5 has context 4 (which is ~~not~~ very useful), T. context 1 2 3 4 ^{looks like} a useful subfunction

From t. approach of (8).00 - 13 + immediate context of 5 is 4; next 0 1 2 further out -1 0 1 2 (3004)

I could code 5 w. a simple "linear" (= sequential) prefix: 1 is immediate, 4 parents

sp. node 3 (0 or 2) in the convention so choice is unique then (2 or 0) next (4) ³ -1 or -1 2 (4) ₃ etc., so we will have unique context for each node, of length that can include entire cond.

If we get a mismatch at any pt, t. context stops.

One objection is ^{for example} that we have to choose arbitrary data: 0 is 2 as context.

say 0 matched but 2 did not: I'd like to use 2 as context, ~~but~~ but if convention said "2 was to be matched first, I'd lose ϕ ". Also, if 2 didn't match, I'd like to

use 1, 0, -1, 6, 7... as a context for 5.

It might be possible to ~~index~~ "index" several // contexts for 5:

1, 2, (3) and 1, ϕ , (6)

~~Other~~ Other than need for extra RAM (maybe available $\frac{1}{2}$) 18 May not be a fair cost to R.

~~Prodn. of 2 // contexts; If both contexts were of length k, we get~~

If we have 2 ^{parallel} contexts, ~~both~~ both of length k, then each context will have a certain "breadth" of its P.D. - The 2 contexts m/l will have a "breadth" of $\frac{1}{2}$ or $\frac{1}{\sqrt{2}}$ as wide (∴ better): But a context of length $\geq k$ would probably have a P.D. breadth of much less than that of a k length context. So much hyper for prodn 2 // contexts of length $\geq k$ (vs opposite)

2 // contexts of length k

[SN] Areas of induction relevant to comp. Problem: Recent Statistical Mac. Xitu.

1 DNA analysis: Protein prediction... (is Wolff variable?) (Muzic).

Other ideas: 3 PPM with straight RPN or (Puz Polish), 4 Modifying of RPN like 183off or 181.0off very only ~~parent~~ (parent) 5 strings. 5 Iurgens, Salomonitz (1997) method (what is this?)

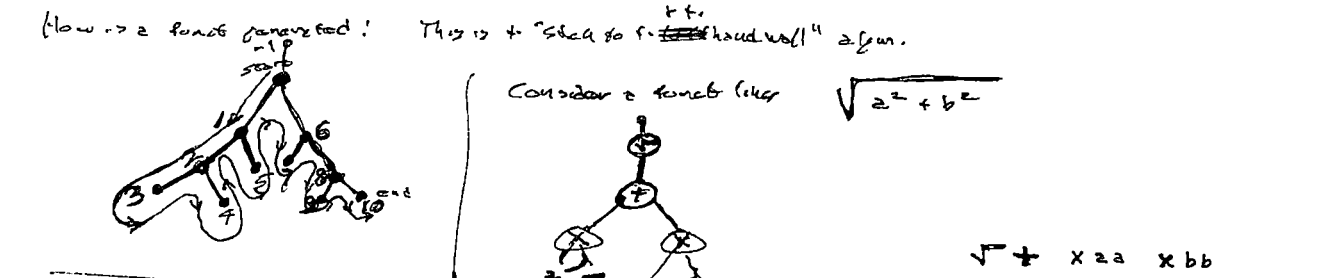
4TH

So every finite string starts w. (0 or more ~~characters~~) B 's followed by 0 or 1 f string:
 A formula is a B .
 A string needs $<$ ~~characters~~ f formulas to complete it as x form it into a B .
 So, for every string we have 2 integers b and c .
 If $b=c=0$, the string is a null string. (If $b=c=0$, the string is a B (only $b=c=0$ gives a B)).
 If $b=c > 0$, the string is a f string (function string): and then if $f > 0$, it's a ~~complete~~ incomplete function f string.
 If $b=c=0$ it's a null string. || If $c=0$, the string is a finite seq. of B 's.

So, a string is either a seq. of B 's (set of arguments) or a set of B 's followed by a f string (incomplete function).
 The inc. func. will have 2 "c" characteristics ≥ 1 . (\equiv no. of B 's needed to ~~complete~~ f string).

Under it 2 B 's ~~needed~~
 f string is actually a function w. c arguments.
 A string that is a seq. of B 's is a seq. of arguments.
 by a f string of c args. ($c \geq 1$).
 So any string is a seq. of args. followed by f .

So: ~~the~~ ~~string~~ ~~utility~~ depends on what end one wants to grow.
 If one starts at the root: So all incomplete functions will be "f" type strings: a ~~branch~~ point, having a certain "c" value needed for completion. Hvr, success will be usually easy "partial" matching, the string to be matched ~~root~~ start w. a few "B"s.



Also would utility of x^2 ... a unary function be easily found?
 x^2 : perhaps ~~leaves~~ ~~could~~ be addresses ~~for~~ ~~args~~.
 One could use ~~both~~ ~~facility~~ (but this misses the point of using PPM).
 Try finding what kinds of functions are found by PPM.
 Recognizing that x^2 (as x^2) is a unary function may be a "higher order" intelligence.

To some extent, x^2 is "recognized". When x^2 occurs, a c is more likely to follow.
 — But if x (B string) occurs, would some B string be likely to follow?
 I guess I want to find out just what kinds of regex PPM could find in "RPN" or "PM".
 Could PPM find useful regex in "infix" notation?
 — Using just the (parent) string relationship ~~of~~ ~~function~~ ~~tree~~ PPM could find useful stuff.

There are probably many ways to code function trees: Certainly some (w. PPM) will be more better than others.

Actually, 181.00 \rightarrow 13 may not be so bad / seen \rightarrow 183.00 \rightarrow 181.00
 q.v. \rightarrow spec

4TH

In Application of PPM to Generation of function trees:



To start off, each possible funct or ~~leaf~~ terminal (= leaf) node has an uncond pc. First uniform, then Laplace formula.

Next, it depends ~~unconditionally~~ on function of a node.

So sprouts of A have ~~2~~ 2 pc d.f. that depends on A.

Next, it also depends on ~~the~~ parent of A ($\equiv B$) or on other sprouts of A.

So in general, the dependency works back wards & sideways. ~~Just~~ for each node we can list nodes that are "distances back" of 1, 2, 3, 4, etc.

There is some Q about whether "intensity" deeply in 1 direction is more valuable

than a shallow match in ~~the~~ 2 & 3 out direction. We do have a priority (ordering) of matches.

At a given node, ~~we can have~~ can we have > 1 post. "history"? If so

use ~~the~~ suitable wtd mean pc. for children. $\rightarrow 182.37$
 183.00

If "context" is defined only by parents, & ~~parents~~ parents... then each node has a unique seq. of parents leading to root. So we can put each node in a "lexical order" & use PPM & its variants.

Look at ^{Jürgen} (Schwartz) paper: Compare his method w/ $.00 ff$ & $.14 ff$.

Also, compare w/ use of PPM on RPN ~~string~~ ^{of .00 ff} expression of functions that ~~was~~ originally considering. Do the new ~~code~~ ^{code} really hold the ~~van~~ ^{any} keys?

In RPN any complete branch is expressed as a coherent string (w/ ~~the~~ ^{any}) — So PPM may be at least as good as Koza's "random branch exchanging"

In RPN, a string can represent a complete branch (as $.20 = 2()$). Also, if it starts out w/ a function ~~it's not complete~~.

It can represent a ~~partial~~ ^{zero} branch plus 0 or more branches plus a partial branch.

Recursively define: a string can represent a ~~sequence~~ ^{represents} sequence of 0 or more branches ~~as~~ ^{as} terminal or a "string".

A string ~~starting~~ ^{starting} w/ a function, can represent ~~2 or more~~ 0 or more branches ~~if~~ ^{if} it starts w/ a function

A "string" starts w/ a terminal (\equiv leaf), ~~is~~ ^{is} ~~not~~ followed by a ~~major~~ ^{major} may not. It can be followed by 0 or more leaves, ~~is~~ ^{is} a ~~major~~ ^{major} may not be followed by a "string".

F T F F T F

~~can~~ ^{can} start w/ a funct. — If a complete t. funct is become a "B" (\equiv complete branch)

or, it can be a part function so all args ~~are~~ ^{are} aren't complete.

A "P-funct" is a string starting w/ a function, but not having ≥ 1 args of that funct.

Partial defn. \rightarrow A P-funct consists of a funct. followed by (or more B's (\equiv complete branches)) followed by a P-funct.

~~is~~ ^{is} ~~for~~ ^{for} every P-funct, ~~if~~ ^{if} ~~it is~~ ^{it is} possible to convert it to a "B" by concating on one or more B's.

~~Every~~ ^{Every} P-funct has an ~~associated~~ ^{integer} number c , ~~that~~ ^{that} (≤ 1) that tells how many ~~concatenated~~ ^{concatenated} B's are needed to make it a "B" \equiv "fs"

Every string starting w/ a funct., has a $c \geq 0$. ($c \geq 0$) B's falls thru many B's are needed in concat., to make it a (complete) "B".

Every string not on fs, must start w/ a * T (terminal). It will have b terminals before first ~~of~~ ^{of} funct. occurs.

Advantages of FS over normal GA:

- 1) More efficient: less LC wasted on re-parsing by pc trials
- 2) Ability to deal w. ~~the~~ ^{the} trials in a good way
- So it can deal w. recursive funct. easily

179.29
00:179.90 : say we have strings of length l , from analysis of ~~the~~ F_{200} d.

Every time an example occurs, we have l^2 updates of perms. Each ~~of~~ t . l^2 pts has d^2 possys to add I to.

between At the end, of all updating, we have $l^2 d^2$ bins that have integers in them, that represent the population.

~~the~~ Given Perm $(ld)^2$ array, I don't know how to use it to generate M.C. Cords ~~or~~ to test trials in l pc order.

07 One possl. guess at M.C. trial generation: By suitable summy over $(ld)^2$ array, we get pc's of each symbol at each l points. We choose 2 first symbols ~~at~~ t .

10 First l pt. w. pc. = it ~~is~~ marginal pc "at that pt. Given that symbol at first l pts, we get t . pc's of all symbols ~~at~~ other $l-1$ pts. directly from t . ~~using~~ $(ld)^2$ array $l-1$ data pts ~~at~~

11 ~~at~~ t to a Assoc Matrix Element. This ~~will~~ give some kind of M.C. Cord approx. ...

12 How good it is ... I'm very uncertain!

13 Assuming .07-.12 is ok. , to generate cords in l pc order!

First pick "a" by pc cond. in only way: first pick t . most likely symbol for each of t . l positions. Then to product all pc's of Row l symbols is P_0 . This num will be used ~~to~~ to calculate initial limits for l cord.

Another way to get P_0 : pick t . most likely ~~at~~ symbol at ~~all~~ positions. Given that symbol, pick most likely ~~or~~ pick symbols at all other $l-1$ pts w. most likely pc's given ~~initially chosen~~ t . ~~initially chosen~~ symbol. These $l+ (l-1)$ choices give another way to get a "P0".

20 Heavy procedure to obtain a P_0 ; we try all cords w. ~~the~~ $P_0 > P_0/2$.
then random $P > P_0/4$; $P_0/6$... (+ use of a factor of 2 is not best; factor of 3 may be bit better.

These pc cords assoc w. cords viz .07-.11 ... Notes .12!!

Other selection ways to associate pc's w. cords can probably be devised.

Since $(ld)^2$ can be very large: say $d=20$, $l=10$; 200^2 array: Means our SSZ for each $vs.$ is quite small. Suggests we use a variant of PPM.

General Remarks about G-P, PPM: My impression is that PPM (or similar Assoc Methods)

are not much good unless SSZ is very large. It will improve G-P, but not by an enormous AMT.

The "quick / rug" of Housar in certain situations may ~~not~~ involve a fairly large amt. of

(internal to human) such. The idea of writing good TSO's is to try to find ways in

which to part to a soln - is, instead of types, a design (a process & assoc. such techniques

that make these parts feasible.

ATM u's sm to ALP.

Cross Validation:

we have a data set $\{x_i, y_i\}_1^n$ we want $y = f(x)$ as a good approx.
 α is a param set. α_j is the best fit α for r -data set $\{x_i, y_i\}_{i=1}^n$ w. j data used.
 So error estimate for test n is $\frac{1}{n} \sum_{j=1}^n (f_{\alpha_j}(x_j) - y_j)^2$
 (e.g. we use param α_j to predict y_j & look at error².)

The variance of this mean is not easy to obtain, since estimates are highly correlated.
 To apply to sm probab. strat's is not so easy, since eliminating a data pt. may not be clearly defined. I.e. eliminating "1 trade" depends on which α was used since each α falls when we go on to out.

We might just barely cut out sections that seem to be of average "trade length & time".
 We terminate a trading seq. out trade before it would have to continue beyond the boundary.
 However, we can always start a trade at the edge of a section.

HV: considers A.H. hypoth: $[IF(x_i \text{ occurs, guess } y_i)]_{i=1}^n$: \leftarrow No adjustable params.
 If all x_i are distinct, this will get a perfect score via cross val. yet it will not extrapolate well at all.
 So cases beyond n data. We could make rule of (1) more general by choosing y_i from x_j that are close to x_i input data, x_j as poss. - so a very A.H.
 "Nearest Neighbor". If we have > 1 nearest neighbor n . default y_i 's ... use a uniform prob. over those y_i 's for i predn.

So it looks like Cross Val. does not deal w. A.H. hypoths.

Well, ALP would not do better, if this were the only hypoth. being tested! HV: ALP
 Normally does include the "prohibited" param (as well as the other A.H. param). - Both have about same wts.
 (Perhaps Cross Val does this too!)

To A.H. hypoth. ~~uses~~ "nearest neighbor" & probab. better than A.H. hypoth. for predn. Also, ALP would give a warning that little if any comp. was obtained from the A.H. hypoth. so we should expect poor predn.

Remember: ALP doesn't do anything but tell how much better model
 If one is doing LS, then to use "simple length" for a par. is not biased. But for (practical) ALP
 LS, we use a guiding p.d. based on previous probab. solns, which which is certainly "biased",
 but perhaps not bad! T. Q is: would it give unbiased p.d.'s?

Another cross val method applicable to sequential data (like sm strat's):

Given data seq: x_j ($x_1, x_2, x_3, \dots, x_n$). For $j=1$ to n do predn. off x_j , using
~~any data~~ model made on x_1, \dots, x_{j-1} . Then error will be e_j . e_j^2 will be
 a G function of j . Since it's to extrapolate into future. As a smoothing model use
 $\sum e_j^2$ constant. (C). This may not be right. In fact, it's wrong, even if
 we used full (nonrecursive) ALP! Make better Model: say $e_n^2 \approx \frac{n}{k+n}$? or $\frac{n^2}{k+n^2}$ (d.c.).

See 179.00 for a paper that discusses cross val. and "dangly estimates".

7/19/04
4TM

I have copy of all of these: usually in full hardcopy, some only partly in hard copy
21 are in PS directory (Bib1 computers), Bib1: see 184 for next Bib1.

Bib1. of papers on PPM & assoc ideas:

- 1) Unbounded Context Length Contexts for PPM. Cleary, Teahan, written ~ 1995 (has 1994 ref)
Reasonable explanation of PPM. Shows how related to Burrows-Wheeler (Burroughs-Wheeler) → 1994
- 2) A Block Sorting Lossless Coding Data Compression Algorithm 1994
Burroughs and Wheeler.
- 3) The Zero Frequency Problem. written, Bell IES IT July 91.
- 4) Solving the Problem of Context Modeling: Charles Bloom ~ 1999
- 5) Modeling English text: PhD thesis w. J. Teahan 2002; 7/24/02 - PS1: HP

PS1 - P28
6/6/03:
trace comp code
Cleary

6) B = (Kashkol, Kurtz, Shkarin 1999) Improvements on Burrows. Is concerned.
speed & RAM needs.

7) PPM one step to practicality ~~Shkarin~~ Dmitry Shkarin ~ 2002: High class.

8) Bell, Moffat, Witten 1995 (in PS/1.7.04): Hard copy of P28:
Reviews how good compression was at that time. ~ 1.3 bits/char for English for better compressors.
we now get ~ 2 bits/char for (good) compressors ~ PPM (Burrows).

9) Bunton 1996: PhD thesis (PS/CHB) 7/19/04 → PDF: Buggy

step changed 5 to 12 improvement
only 18 pp: Computer Journal 1997
PS version with loss bugs
- but 2000 version
* Sincerely Motivated Improvement
of PPM Variants. I Made Hard Copy
- Mike says go to front up

00

10

20

30

4:14 3:14

1879 to 1910 or from 1849?
or 1869 " " ?

Here,
We will ~~be~~ be interested only in the symbol probabilities and in how they were obtained.

SN In the processing of \vec{L} v. z. f. kernel k : T. wt. of ~~it~~ depends on ~~kernel~~ $|X|$, $|X|$ being ~~the~~ distance from t : "center" of \vec{L} -kernel. Hvr, in fact, ~~it~~ t -wt. should depend only on the ^{number} of matched symbols betw. t : "Matchee" and the context considered (e.g. "Match Length").

This would change discretely, integrally, not continuously like $\frac{1}{|X|}$ or $e^{-|X|}$.

Would this make much difference in result?

Instead of a ~~kernel~~ continuous kernel, one would have a function of integers only, ~~the~~ integer being t : "Match length $M_t(\vec{L}, \text{cont})$ ".

7/23/04 ~ A way to do this: For each "Length of match" (w. present ~~kernel~~ corpus to be predicted)

We have an assoc. wt. This set of wts was obtained by optimizing over "disturb past" - it is a kind of "updated" a prop. There may be a relatively simple (approx perhaps) algo for updating t : "a prop" ($\equiv t$: wts). Every so often we will have a long R of match > any ever before - so no "a prop" data. We could, just not use that data until a suitable wt for it had been found ...

Alternatively, we might "smooth" t : known wts. ~~of~~ extrapolate far to "novel" "Length of Match".

Now into into into cont. of info in t priori - or, more ~~exactly~~, how best best t : info & w.

t : particular symbols that are related from each context.

So we have t : longest context "abcde"; t : next longest is "bcde". T. ~~data~~ ratio of a props would be t : pc of symbol "a" Also relevant is size of each context.

Say ~~the~~ prefix abcde has occurred k times. Using weights Lap's rule for pc's of ~~predicted~~ symbols; compute ~~the~~ product of those k products by ~~the~~ context abcde.

Then compute same product of k pc's for ~~same~~ same predictions by context "bcde".

t : ratio of these 2 products is part of t : relative wts of t : 2 contexts.

Another factor is t : pc of symbol "a" which is t : ratio of a props of t : 2 contexts.

w. ratio of $\frac{abcde}{bcde} = \frac{.21 \cdot (\text{pc of symbol "a"})}{.22}$ } I'm not sure all sum of this! how real "a" should enter is unclear!

$\# \approx (pc of a)^k$ proper factor?

Hvr, I suspect that $\#$ may not be very imp. It may be $\#$ & flow $\#$ to $\#$ in entropy/symbol.

- but what is of more import in GP, is finding reasonable pc's to n grams, n -grams, etc.
- T. Sort of thing discussed in 179.07 paper.

4TM

0:17:40 : (See section 1, p7, eqs 3 and 4 for fuller discussion).

z_0^i is the a priori probability of O^i and is associated with the language used to describe the function, O^i .

In Phase 1, the system looks at the set of O^i functions that have been successful in the past and obtains from them a probability distribution on candidate O^i 's for new problems. ~~the~~ PPM is used in the present

implementation to obtain this probability distribution, which is, in turn, used to guide

the search ~~for~~ O^i 's for which eq 1 is large. This probability distribution ~~will be quite different~~ on O^i will be quite different from z_0^i , which is a function of O^i only, and quite independent of any information about past successes of various O^i 's.

A telling criticism of the Phase 1 methodology, is that it does not really understand that it is trying to maximize eq 1. It just makes trial O^i 's

"similar" to previously successful O^i 's and uses eq 1 to evaluate them.

~~It has no way to take advantage of~~ and ~~has no way to use information in~~ It does not understand why certain trials are better than others ~~and~~ has no way to use information in

unsuccessful trials. In spite of these deficiencies, we expect that Phase 1

will be good enough at induction ~~to bring us to~~ to bring us to

"Phase 2", which does have a good working understanding of optimization

and can use information generated by ~~both good and bad trials~~ both good and bad trials.

Before telling how PPM is used to generate trials for Lsearch, let us first ~~look at~~ look at how

it is used in text compression. To compress text, PPM ~~works by~~ tries

to predict each symbol of the text, based on ~~the~~ the sequence

the "context" of that symbol; i.e. the sequence of symbols immediately preceding it. Here, we mean by "prediction", that PPM is able to assign a probability

to each possible symbol that might occur at a particular point in the sequence.

The product of the assigned probabilities of all of the symbols that did occur, is a measure

of the compression of the system. ~~The~~ ^{negative of the} \log to base 2 of the product will give

the number of bits needed to code the sequence.

We can use PPM for full climbing if we can update using enough (Good) cond. Hvr. older takes pos. can't be used otherwise. In T & ZT, we can only update fully, otherwise T & ZT "Round"

4IM

dep

This is 2 rewrite of 159.20ff. T. 2/26

The Application of PPM to Phase 1 of Alpha.

Introduction

We will ~~describe~~ show how to use the PPM text compression program in the implementation of Phase 1 of ~~Alpha~~ the Alpha, machine learning program.

First we will ~~describe~~ explain how Phase 1 works and how it differs from the more advanced, more general, Phase 2.

The PPM ~~program~~ predictor/compressor is then described and we show how it is used in the L search part of Phase 1. There are several ways in which PPM can be implemented — we will discuss the trade offs for some of them.

L search can be done by either limit doubling, parallel search or Markov search.

In the present trials, we will be using limit doubling & we will explain its advantages for the present application.

4IM

one

00: This is 2 rewrite of 159.2011

of T126

Application of PPM to Phase 1 of Alpha.

Introduction

We will show how to use the PPM text compression program in the implementation of Phase 1 of the Alpha, machine learning program.

First we will explain how Phase 1 works and how it differs from the more advanced, more general, Phase 2.

The PPM predictor/compressor is then described and we show how it is used in the L search part of Phase 1. There are several ways in which PPM can be implemented - we will discuss the trade offs for some of them.

L search can be done by either limit doubling, parallel search or Markov search.

In the present trials we will use limit doubling. We will explain its advantages for the present application. } changed.

A critical part of the Alpha System is the training sequences. We will describe various methods of implementing them.

Section: What is Phase 1?

Phase 1 is described in section 7.2 p. 294 of 50103b (NIP's paper). The AZ language is used to better describe candidate functions, and assign a priori probabilities to them. In the present implementation of phase 1 that we will describe, we may use AZ to describe functions or some other language to describe functions.

~~We will use PPM to assign probabilities to them.~~

Associated with each language, there will be a natural a priori probability distribution. This a priori probability will be used to help evaluate the "fitness" of candidate functions.

There is a different probability distribution on candidates that tells how close the candidate is to successful candidates of the past. PPM will be used to implement this distribution.

The problems of phase 1 are all "Q, A" type problems:

(see section 1, p 6 of 5105b for description of Q, A problems). Q, A problems are a fairly general form of induction problem. The main problem in phase 1 is always to find functions O^j such that

our (no) hypothesis $\Rightarrow \prod_{i=1}^n O^j(A_i | Q_i)$ is as large as possible } difficult task

Note: PPM almost always is of 50103b. (3) of 20103b. 1 page?

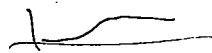
30 : My conjecture about Barron's result on (∞) ANN-like models: That while they need fewer params to desc. a model w. given/precision, they ~~require~~ require ~~more~~ more precision, so perhaps same total no. of bits.

Having fewer params into ↓ cost of finding soln. ... but maybe not: If one looks at random, search finds depends on no. of bits in soln., indep. of no. of dims.

If I'll climb is being used, ϵ of soln will depend on hyper-volume of "Basin" that contains a soln. (T "Basin" is a region in which ^{all} slopes is toward desired soln.)

Maybe Basins are smaller for high precision solns.

For induction, it's just ϵ no. of bits in ϵ soln, not ϵ no. of params, that determines ϵ w. of ϵ particular Model.

How, this Barron result is wild! The log-sigmoid  used has no params, so ϵ w. is not have all ms. That ϵ no. is able to desc. functions

w. few params is amazing! T. Kolmogorov permits expressing a n dim function several (dim) func's, uses very complicated (dim) func's (as "Basis"). — Barron does not use that trick! T. Q is: is it possible? Increases in precision would not seem to be adequate!

At any rate, it may be possible to show not one needs lots of bits to desc. an n dim function w. precision ϵ using logistic func's (of no params). This could be done by counting no. of n dim func's w. precision ϵ (satisfying a certain constraint) that is ϵ x ϵ .

30

∴ Google: To day I was surprised to find that Google (which says it has 2^{32} web pp ($\approx 4 \times 10^9$)) can do PPM! I.e. if I write a sequence "a c a b d f g" for Google search, it would scan read all 2^{32} pp to find that sequence & present results to me.

Note quotes on a c a b d f g ... it reads all 2^{32} pp. to get hit. Takes < 1 second.

Maybe we could request Google to generate "Google random" sequences.

Maybe randomly generate. Also Google random functions — to be used to solve problems.

I don't know if Google already will read pages random pure LISP, or "C" or Fortran.

— Try this out. For certain particular sequences, Praxis will be too

many vars. Google will ask for more further consistency. We can write:

{ "a c a b d f g" LISP } to constraint (initially) — Other ways are possible... set by

debug Language, Google has other constraints by new "Advanced search", "Macro Contexts".

Also, look at search engines BLINKX (7/16/64)

Maybe get Lew Minton to write program to do this on Google

[SN] On improving PPM w. definitions: for an ordinary (linear corpus we make up

definitions, it's either so often, we reverse corpus w/ Praxis debug. [166.11 reverse it desc. a possibility way to find usable vars.]

How, when we define function (subtrees), we also must periodically reverse ϵ corpus w/ Praxis

debug. My own process, is that finding Praxis / subtrees is not so easy [165.20] — Praxis way to find.

to reverse something (the PPM search, because PPM search can be desc'd in terms of Praxis

is subtree search praxis to be a Praxis search. [See 167.375 ff for stuff in Praxis spirit.]

Also note that subtrees often can be reversed by Praxis! (167.18 discuss)

On Pragmatism (Egoism / Emerging / Awareness of outside world) in TM; ^{actually just about any large enough, well defined problem}

Any open ended problem - like OZ probs - induction - or even INV probs, are potentially capable of having "self-improvement" (= SI) as a byproduct.

This will occur in any "large" INV problem (i.e. To P goal is soln. in M in time) or in an OZ prob. w. large enough CB. - In such cases, it will be worth while for TM to spend a fraction of its available rc on SI.

Working on SI in a TM that can use resources of "outside world" - such as Internet is potentially capable of realizing outside world can be "Manipulated" (as can the USSR). = "Warrior".

We may be able to limit TM's internet access so that it can't manipulate outside world.

One way would be to give it a "copy" of the internet that is not connected to the outside world -

say one of Bostrom's backups or archival copies of the internet. Unfortunately, it would not be able to use Google, which is similar services ... which cuts down on its considerable

bit eventually, it would be able to simulate Google's functions (if it had enough cc ... Google

has made 100k computers) (20\$ 300 per computer, 10M bytes $\frac{10 \text{ k}^2}{300} = 30k \text{ computers}$.)

But, Google may be "effective" w. much less cc.

We might give TM some head of "limited access" to Googleish services ... otherwise, it could easily surf the web and manipulate it in other ways.

→ [It may be that it is quite difficult for TM to do much in "outside world" by access to Google.]

: Getting back to the sheep; review of 153.00 -> 40 -> 160.30 -> 40

Best archive review, w.o. much bibli refs.: then put in bibli refs later by scanning back thru notes: Review will give places to Guy references to relevant discussions.

O: 170.90 : not codes, fast. corpus, but over PEMs. : This makes biased searches more difficult to design, but still not impossible. Is there away to design "unbiased search" that will not discard some possibly imp. (good) PEMs? Answer: search over partial recursive forms.

Perhaps what we want is a search technique that has a "short down" (e.g. prep).
The "bias" in the search can correspond (perhaps) to a bias in the U.D. itself — i.e. like penj in B.

05:170.98 **SN** on objections to Occam! PS (6/13/02 "Against Occam")

If these guys found a short code for corpus, Pen found a longer code for ~~corpus~~ ^{model} w. better fit,
Pen is likely that BC of model + BC of data was shorter (for Pen's best) Pen's best model w. "short code".
— But look carefully at just what Pen claim actually did!

O: 166.90 : **SN** **SM** : Turn on Sching's Max (fitness for ANN's SM): For k different interacting stocks, $\{S_i\}^k$
Find a ANN function $\Rightarrow F(\vec{S}) \approx 0$: To do prediction, fix $k-1$ of k 's ^{known} ~~known~~ parameters as
inputs to $F(\vec{S})$, then try various inputs to guess S_i : It gives a D.F. (not necessarily 0)
but this D.F. costs give ϵ (inverse) d.f. of predicted S_i .

This uses fewer params than having $k-1$ drivers for each S_i ! — But τ transfer,
we usually have $k-1$ S_i values from ~~the~~ yesterday & we want to guess S_i from today —
So Res. A will not be useful!

- Also, we can have one ANN w. k inputs of yesterday's values, & k outputs of predicted values. We can then minimize (\leq sq error) ~~error over data set~~
- We should be able to do better than Jayram, since we don't need the set's test set.
- Also, we could do k models of ANN, which will get significantly better results.

— i.e. J. used single central "best point" in param space, having max ~~min~~ fitness.

I would expect ~~of~~ of k param space to have a "best pt." — parallel models.

→ Also, we don't have to keep our drivers in the set of S_i to be predicted ...
[So if we used some set of drivers for all drivers, this does restrict us: —
In general, ϵ . no. of drivers will be \gg τ . no. of drivers: This was ~~just~~
use a single ANN for all products.

→ Also, we would ~~be~~ change the set of drivers, from time to time. (See soy problem
enters here). Also we would change the set of drivers from time to time!

Each day, we update (k windows) the ~~set~~ {set of drivers \times set of drivers} array, giving affinities
of each driver towards driver. τ . driver set is \gg δ . driver set. Any buy can be a driver, but ϵ . driver
have to be "subtle" (large time, large sums of money, large volume).
The updated array gives good clues for modification of yesterday's ANN.

7/13/04

4th

00: : types of Mathematical problems (like in Phase 2 of α)?

16 : **SN** Opt. value of using $\left(\begin{matrix} \text{random error}^2 \\ \text{K-L error} \end{matrix} \right)$ in Thm 3 of Sol 78: first, this larger than $\left(\begin{matrix} \text{error} \\ \text{error} \end{matrix} \right)$ so more significant, but $\textcircled{2}$ If α alphabet is very large so pc's are often very small, $\frac{1}{K-L}$ then for $\sum \left(\frac{p_i^2}{\alpha_i} \right)$ to \rightarrow sum to $\ll \alpha$, say, is not such a big deal, since p_i 's are \ll very small. (Hm, that $\sum p_i \ln \frac{p_i}{\alpha_i}$ converges in this case is of Much more interest!

In fact, in "Bay induction" produce w. corpus of recorded, finite set of finite strings, we can use R_c subset strings as an alphabet - in which case the true generator is ALP, give p_i 's resp. , in which case $\sum p_i \ln \frac{p_i}{\alpha_i} < \ln \frac{A(\text{Corpus})}{P(\text{Corpus})}$.
- In 19 we have $\sum p_i = 1$, so this eq. has to be modified a bit.
- Also in 19 fact eq. corresponds to predicting only one "token": unclear how Rens of Sol 78 applies.

13.165.18. **SN** On Occam's Razor? If $\alpha \neq \beta$ are 2 PEMS α P_α , P_β are both codes for corpus, including some of α & β . If $P_\alpha > P_\beta$ then by corollary of Th 3 of Sol 78, P_α has smaller expected ms error for corpus.

My impression is that Bayes (because of the peculiar properties of ALP) implies that P_α will probably do better in the immediate future. See 164.30 - (85.18 for poss. ideas) on this. \rightarrow Perhaps if we assign a priori α & β properly, the discussion will be as if α & β were chosen before data was seen. If $\alpha \neq \beta$ are approx of α & β resp, then, w.o. seeing corpus, we would (perhaps) evaluate α , β to be 2 codes.
At present, this is a simple "Unsolved Problem".

See 164.30 : It seems rather odd, perhaps identical!
So 3 probs : maybe identical solns!
1) In approaching to U.D., how does know our model is closer than our Peer to U.D.?
2) 164.30 : How does ALP get rid of "grey sea, rest sea" diffy?
3) Occam's razor 170.23, 165.17
[4] (Not said) In induction, we need ratio of pc's of (2) future corpus. How can we find this in an "honest" way be sure... A different? Recall, normally area covered... \rightarrow 171.90

37
38
39
40

4/11

"continuous" in context of recursion in PPM is anti of DG expected from PPM.

20:168.40 Methods that are "diminishing returns". We can spend a very large amt. of time, but, reasonably w. much time to spend on updating.

The "Mixed update techniques" may not work together. Suppose we use slow method to derive various sub-functions & ways to put PPM together. If this ^{slow} update is followed by a "fast" update... like PPM, we take all into in 1 slow update. PPM means slow updates can build upon previous slow updates.

A way to Usefully Mix the slow & fast update (PERHAPS!) After "slow update"

derives new "grammar element" & compresses corpus by recording it using these grammar elements. PPM can be used to further compress PPM code (sometimes?) - Anyway, we might use PPM to "interpolate" between the "slow updates".

Or, we might have the sigms defined by the previous slow updates; code is corpus in terms of these defined inputs. So each slow update creates a new input alphabet for PPM updates.

20:168.29 A life! Problems in TS "identifi". We are given a spaced sequence generator containing random words & several unknown phrases (some non-linear, say). To identify the seq. generator in PPM & confirm prediction, { the phrases can be chosen by each contestant & the other contestant... or chosen at random for both } (a common prob. for both).

20 is perhaps closer to the pseudo. prob. to unit PPM to solve. The "contest" to get max compression for "Colony Corpus" is a problem to 20. So in 20 ft, the winner could be the one that best predicts the sequence. The space of models for the Corpus can be quite large, but it has to be stated. The number of phrases. The size & nature of the space may be such that identifiability of the generator

is impractical, & that only approximations are feasible - it can be used for predn.

which is much like the General Predn. problem. [We compare stochastic T.S. using random discrete phrases & random continuous phrases. T. "God" could give PPM as part of TSD]

It might be (possi., feasible) to use 2 or more techniques to pool their forecasts in trying to predict a sequence. Could there be a way for PPM to show "Prior Sequence tracks" to create a prediction scheme better than a simple linear combination?

There has been much work on optimum way to combine forecasts of the different experts predictors... I'm not interested in that... I want the kinds of strong techniques, tricks, sub techniques.

well, what about a way for deriving soln. techniques for time series or forecasting

AT

ALIFE → "Micro Cosmic God" 100

Spec
149-19

SN on Alife w. "Micro Cosmic God" ~~is~~ How do we structure the Alife Society so

it can usefully work toward solving diff't problems? [that of. Middle Ages; in which Mathematicians would pose problems to one another: Trick is its easy to make a diff't problem that the proposer can solve easily but the "acceptor" of the problem finds diff't. e.g. factor product of large primes; solving a f.c. eqns.; perhaps combinatorial problems, etc.]

Hvr., if a problem is of that type (≡ "trapdoor function") no "acceptors" will be found! — so: Example: Proposer says: "I have the product of 2 primes and the result is 200 digits long. I will give you \$1.00 if you factor it in < 1 day you get \$100; if not, I get \$100."

Consider eqs $e^{2x} + \sin by = 0$; $\ln(x) \neq \ln(y) = 1$

If I give you $a, b, c, \in \{ \text{limits on range of these params} \}$ can you get soln. for x, y in 1 day to .001 accuracy? or, a bot is selected at random; Both Mathians try to solve it; first one wins. or, each is allowed to select params of problem for opponent Mathian → 125

To push tech skills of this culture in certain directions, we have to frequent funding TSO of the "Micro Cosmic God" — This gives very large prizes for solving old problems; Individuals may solve them; Or Coalitions of individuals may organize to solve these problems, & divide up the profits in some predetermined way.

15-18 could be the only source of money in future economy; winners of 15-18 may engage in bets w/ their community citizens. Other citizens could go into debt if they lose a bet, & have little or no "Money".

If a proposed problem is too hard, no one will accept. If challenger so

D. Fogel's checkerboards are very similar to this idea of competing Mathematicians

if time limits must be set or the problem must be made easier (Gives hint!).

Another kind of problem: A math task w/ 4 params: 2 Mathians are allowed to select 4 param values for the other Mathian to solve. First one to solve, wins all of the money

One BIG Q: Would this A-life approach be any better than doing about something directly, w/o. "Alife"? i.e. simply try to develop population of competing, cooperating pgms working toward common goal of "Macro TSO". Well try both approaches: See whether can suggest ways to better the other. The A-life ~~may~~ write

biology transparent as Trojan Viruses on the Net of computers. → 169.30

Spec 167.40 "Contexts" situation has been dealt with the party no object a pd. of which (to separation of the cloud) does not. It could start a new sub-trace: perhaps "barren" 2. substance from part of the corpus; perhaps present a final output, 6. output of one or two "sub-traces".

167.375-168.33 write the proposal as "very close to" or "suggesting" a ~~particular~~ general outline of a crossover scheme.

The large stuff will result in some "guiding pd" updates models that are very expensive? So what we do is/ use the fast-update models, but even so often, use the expensive models.

The frequency of updates given by eq 124.33: we can use up date

TSQ-10

00 (166.10): Note that PPM usually ~~derives~~ uses only short context for prediction - not f. context branch.
 This is at once better than worse than k-gram. "Better than" because it always detects common functions. "Worse than" because while it can detect large common trees, it usually assigns a rather low pc. And then perhaps lower than they should get.
 On the other hand k-gram may assign too high pc to them, or just assign low pc's to ~~some~~ common short functions. 18

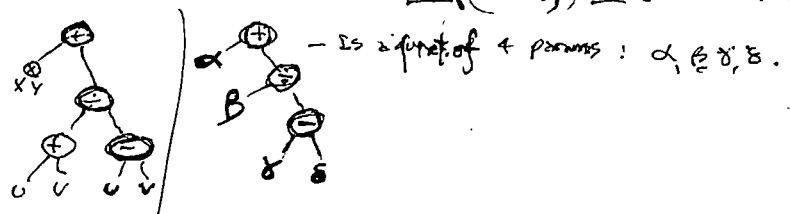
01 [SN] **TSQ design:** Start out w. TSQ in English: but be sure examples are "well defined PA problems": E.g. I could do the ANL of SAARB. Keep paradigm of TM in xy level English, and Express ~~the~~ Hours in English. When 1. TSQ seems adequate in English then try to formalize it. But do so in way that ~~is~~ is clear & easy to work. "English" gives a "fuzzy" idea of the (ing) process & enables us to transfer it ("Material") in various ways.

167.05
2.166.10 [SN] **In PPM used for "Guiding PD in LSPH":** Note that PPM uses contents of various layers for pred. - so if ~~some~~ subfunction is represented as a sub-string, PPM could "recognize it" as a useful, predictive context.

Trouble is, all (or perhaps not even most) functions in an empirical context, will not be representable as substrings. The tree representation seems much better (in some respects)

If we put what we want to recognize is sequences of signs in "EPP" in them.
 Say A...Z are news: Then ABZ BEZ ZD AB BE XX D are "similar"
 Since they have A, B, Z occur in that order; ~~but~~ This sort of "Material" is done in DNA analysis. → 19.25 or Ming Li's ppm to do things like that

What I have in mind: Sum (mult x, y) Dir (B Sum u, v), Sub (u, v)



→ [SN] **Formal modify notation** (with use of list) to make it more clear (likely to be discussed later).
 The use of Trees to describe functions is one kind of "Notation", But we have to map it into more simple form we can use PPM (perhaps not so far we are generating a code is at that point).

We have generated one or more subtrees. We can look at corpus for previous occurrences of one or more of these subtrees or of subtrees "crossed" to them. From the way they (16830) →

FTM

Ockham's Razor → v. a. exp. data? (Also note 170.23)

164.40 badly distorts the merit of fit of hypth. to test data. In using U.D. this problem does not occur. After we see test data, we can make any hypth we want. If it is A.H. it will automatically get small WT.

T. Arg. of 164.30 ft is an impt. part of t. Q: "Why is the Incomparability of D.U. irrelevant to its use in practical induction?"

06 It is still possible to look at test set, propose a single A.H. hypth, & stop. This Ritz seems like a stupid way to do predn., it does seem legit. Is there any such method or general principle, that prevents use of 106? Well there are the set of k A.H. codes, that each predict a different set of t. (different alphabet symbols (w. = R)). Also f. Bern codes. A codes should at least do as well as Rizes.

0 There does seem to be a problem here, in that one could look at test set, & devise a set of codes that would bias the prob. estimate considerably.

12 → Such using a principle UMC is a specific unbiased way to do predn using U.D. (192.30)

13 Another assoc. idea: People regard the choice of UMC as arbitrary, so change in code length of a corpus of 10 or 20 bits/bytes over not imp... "we can get that with change of UMC choice". In fact one can't change UMC so easily. It represents our apri info. — Changing

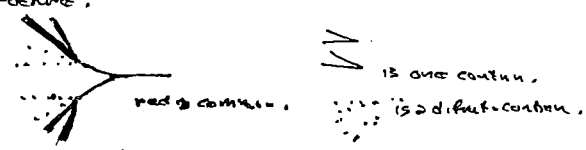
16 UMC may make Ockham's razor unusable (illegal — since all info is now in a priori)

17 [Ritz Ock raz: Use convergence Penn: T. Shorter codes, (smaller is rms error wrt converge of data. $\frac{PE(\text{corpus})}{P_{true}(\text{corpus})}$ → But what about future?

0: 164.15: In discussing Coza's Lazy Trees v.s. PAM PAM is usually described in terms of "context trees".

Can we somehow use Ritz to dec. "(Program & function) trees as well? One trouble is that I'm not really familiar with Trie concept. But just as contexts can be put on a tree,

branches can be put into some kind of tree structure. Superficially, it's not clear how to do this! Seq based tree amounts to 2 PAMs



Another tack: Consider 2. PAM (tree): Each node in tree can be considered to be a "context" — (which is + branch growing out of that node): How 2 contexts can "match" to some extent, but if non matching parts can have different sub-branches — so its much more

complex PAM (more context. A big Q: A matches B to some extent, A matches C to some extent. If A B & C are branch contexts, its hard to say whether B or C is a stronger match (as opposed to linear context, where a match of first 5 tokens is better than match of only 4 tokens.) In branches we can also count word matching tokens, but if simple number of matches is not a compelling "evidence" criterion

Another big difference is that linear context has a lexical order, because of lex order of tokens. But if a branch of trees match in k tokens, there is no order way to order that bunch of trees. This lack of order makes it hard to "sort" the trees out.

Let's look at Lazy Expressions: (RPN)

7/2/04

4PM

Try set / Test set problems how Universal D.F. solves it. 3:30 ft

16:53

164

13.52P	390V
14.03	400V - 30mm
14.04	4.01V
14.06	4.01
14.09	4.01
14.16	4.02V 61.8 MA
14.16	4.02V
14.19P	4.06V
14.39P	4.09V
14.48P	4.01V
14.49P	4.08

00:163.40 : Ont. problem of $2+2+2 \dots$ K times being uncodable to Lisp in many ways
ways, so stop to pc of e-nave "2" int. sequence.

However, an experienced TM will have found, empirically, that the prefix "2+2+2..."
is likely to be followed by "2". Since we are not using the univ. d.f. in Lisp, but
rather a "MBC", we don't ~~add~~ add all those pc's together (as in ACP).

So maybe PPM would do well here!

06 (163.40 spec) : ^{rigid} "rigid" strict rules here, are for creating the Guiding P.D. - They are like strict rules that
enable PPM. The program generated is then used to create code in Lisp. [We must, here,
consider one shot (long), otherwise we are too restricted in our search (too "slow"). OSL corresponds
to GP picking up a random branch from a random cloud --- this branch having occurred
only once in the corpus.]

09

It may be that PPM can do as well as a better than Koza in 09-10 in general

with swapping branches (no PPM doesn't actually swap branches: it borrows branches).

Koza's usurpation song! are they expressible in PPM? If Koza actually assigns a
Tofra to song, then, perhaps PPM is "left behind"!

15

Also Koza has a way to do recursion...

It may be poss. to get all of that w. Forth, since Forth is universal, but Forth may not
assign near to some pieces as Lisp does. (Hvr, we start w. all loop functions
in Forth's "dictionary")

20

One way to get definitions is by my PS-discy routine. - I also get some recursion by $A \rightarrow AB$

Hvr, my PS-discy routine is probably very slow relative to PPM. perhaps start with PPM

when you get on a good slopy part of it, go to PS-discy. When to switch will depend on:

- 1) Time to execute code (\approx "fitness cost")
- 2) Expected \uparrow in npc per iteration due to PS-discy.
Look at up to
- 3) 124.33 for way to do analysis.

25

30 HA: (Cyclic Great Breakdown): In using U.Dist for probus: How it avoids "Try set, test set" mess.

Say we have some data: before we see data we use U.D. to make a seq. of models ^{Meaning?}
for data. When data comes in we start w. first model. If it's "adequate" we stop: if not,
we continue using new models, (which always in comparison to (ones) & seq.

In a "Try set, test set" method we look at try set, hypothesize model, then
test it on the seen test set. If it is inadequate, too bad, we need new data.

U.D. method enables us to test an arbitr. no. of models on the same data - i.e. it uses all data
for the test set, since the try set is empty.

In the 2 data set method, we can hypothesize models that fit very well, so, it is not possible
for "model cost" same data method to produce well fitting models. So, having seen data

6/30/04
4T M


Phase 1 Induction Methods. ¹⁶³ 19th → 164.25

2205 : 6/30/04
β

00:162.002 1. λ (e^{-λ} order) param in weighting have to be decided on.

In regular Phase I Q.A.'s we just saw the past "acceptable solns" & use them w/ = wts as "Th. Corpus" for updating. For each new problem, we may try for several solns - the Pop will usually have different G's (≡ PC's) - but we will be using all solns. in predn. In PPM in Phase I, we may want more wts. on solns to reconstruct problems. What we could do is take some fixed population of some G funct., & adjust params of both PPM & conventional GP, & see how good the results of the 2 methods are. We can try this w/ various corpus's.

Note 12.8.83: eq. for freq. of coding.

Various G functs.  Another way to think about it: for U.D., the set is all data from better + corpus (5 test sets) is seen.

→ (SM) Take all pairs of each interesting / papers / paper on GP. T. 2-3 times of some measures, but

it's quite possible I'll find a better way to do w/ PPM. Do this soon: Do important job, then add genes, papers, labor. T. COS are "add to table": Perhaps try to do some later reviews, leaving imp. corpus, derivations.

- 7/6/04 So list various things I want to try PPM on so I can see which forms of results to a degree.
- 1) GP, 2) Compression of E.g. of Lisp genes; compare w/ other variants of PPM.
 - 3) See if PPM can compress L. 4) Try writing TS: use GP 2/0 PAM to form it.

Q: How could GP work on Phase I problems? In Phase I we try to find short codes for F.

Corpus. ~~...~~ The fitness function is a MESA ∴ No "hill" just "yes/no".

→ Unless we start w. a.H. + promising primitive codes & try to improve them.

What kind of fitness funct. did Koza use in his electronic circuit design? At first, start with a program, for a very long time, there would be no identifiable peak in "fitness" (i.e. "MESA") In general, it would be instructive to see how GP solns. work by conquering MESA (LNU) problemists or prob.

In 19 Most (non-universal) methods of induction can be "hillclimbers", e.g. finding drivers, & domain cons. Optimizing the params, trying various functional forms (w/ optimizable params) Universal codes can be found by making trial codes then "corrections" to trial, (MS error produ is a special case of this).

In my primary about Q.A., I found to think of a search for $O^i(A_i/P_i)$ in which ~~...~~ does not preclude the finding of O^i : This is probably a bad mistake! Use of Correlations, random driven (drivers) act do involve "looking at A_i's".

It would be great if I could find a (heuristic) "path" from, say A.H./promisc. codes to a final single short code. My work on P56 Discy has been in that direction - Discovery of users, then users, then ~~...~~ concat rules - but, to forget stuff is all rather strictly determinist non-heuristic such. - Derivative rules for study these regys.

(164.06 spec)

6/27/04

161

4TM

SM .00

>0: From Φ & Δ (29.20): ~~the~~ A criticism is that Φ & Δ drivers \rightarrow Δ drivers sets are Φ same. It would be better to have them separate. We want ~~the~~ large Δ for Δ driver, but Δ of Δ drivers is irrelevant.

If we have several drivers for a given driver, we can study recent contexts of drivers (or context). e.g. Φ context of drivers on previous day. How to do "Partial Matching" is unclear!

If we have N drivers; each context is a pt. in N space, so we can use a machine (conditional or prior) to tell which today is Φ to other days in past. Certain of Φ N components will be more imp.

Run others \rightarrow bringing Φ ~~into~~ matrix: — Φ Φ is beginning to look like linear product!

Not quite: ~~is~~ is more like "Case Based ~~product~~ product" ... we look for sub-features in past.

Anyway, Δ BIG is ~~is~~ as contexts of stock prices good drivers (\equiv predictors). \rightarrow

Since single stocks are not bad drivers, contexts can't be worse, but are they really ~~more~~ significantly better? (considering SOY diffy!)

00 (SRAC) 158.30

SN on Adequacy of PPM for Phase 1: W. an adequate TSC, PPM has to work! - (presumably to language used is adequate). While PPM ~~was~~ used for Compression. Needs large corpus (E.g. 100k) it shouldn't need this for Phase 1. - I think I have a lot of to be needed in machine mechanisms in Phase 1 - as is, PPM is pretty much kind of Bernoulli's or "Laplace's rule" type of "duccor" - which should be a decent for Phase 1 (perhaps for all induction!). Try to find/recast to SAAB TSC for ANL. See PPM fails about same P's as original PPM. - Perhaps no using STORAGE instruction! - since PPM is supposed to "simulate" "Storage" (or actually Emulate it!). It a conc. seems to have its sub-concs in it. but if they have been .11

In Phase 1 Notes (last 2/100 P in particular), Part 2: Secure/empt. ideas on TSC construction/design.

~~Phase 1~~ Enter these remarks & try to summarize next empt. ideas. ABCdefg ABCDEFghijk $\int_{-100}^{100} e^{-\frac{x^2}{2}} dx$

11 .09 -> Irnd - yet to composite concs is not bearable in reasonable time, then does it one or more

How's messy from Conc. net - - - that have to be Irnd first: This can be Irnd directly (perhaps at by cc) or w. "hints" to various degrees, or, to ultimate "hint" is to

16: PPM into how directly. SN On using 2-3 concs for PPM, for conc generation! All previous concs are represented by strings that start w. a start symbol (S) But we insert a new content (e.g. the front part of a conc) into L, we change it (S) to (N) (N) ~~is~~ is a unique symbol that occurs nowhere else, so it always causes mismatch at that pt. After insertion (N-S) of T. ferris, this is to prevent overfitting from every past Z, to previous conc - which is irrelevant.

0: 159.40 : One Limitation of Phase One is that it makes new trial O's by looking at previously successful O's, and making trials that are "similar" to the successes of the past. Though it does not understand ~~that~~ optimization is what is wanted, ~~at least~~ it is probably capable of finding reasonably good O's. Its ineptitude may be corrected in Phase 2. ~~It is not clear that it is possible to find a better O than the one found by the current method.~~

25: Using a corpus of ~~new~~ successful O's of the past, Phase 1 induces a distribution function of over/candidate (trial) O's. ²⁵⁵ It then uses Lsearch to find concs of Max a posteriori probability. ^{255R} ²⁵⁵ In the ID report, the distribution ~~is~~ is generated by the language A-Z (see Appendix A ~~of the ID report~~, p. 27 of ID report) ^{255R} ²⁵⁵ In our system to be developed we will be using a PPM-like method to generate new O's from a corpus of successful O's of the past.

30: I ~~had~~ had a modification of PPM in mind and I wanted to test it on English text (perhaps computer programs (Lisp, say) - ^{← The Calgary Corpus has some Lisp.} modifications with other variants of PPM. ~~to see if my modifications were better to company~~ ~~more or comparable to other modules of PPM.~~

[Follow w. a detailed explain of how my version of PPM works]

35: ²⁵⁵ ^{255R} When an acceptable ^{solution} ~~found~~ has been found < T. criterion for acceptability can be various: ① Max search time is exceeded or reached: Pick best set of concs (log for ② Total ~~probability~~ a posteriori probability of concs is > a certain threshold ③ Trainer may decide when search should terminate. - on heuristic Basis > This ~~new~~ ^{solution} is added to the corpus of solutions. Periodically, the ~~probability~~ probability distribution used in Lsearch is updated using Rec. latest (162.00)

Max search time
Version of the corpus
Spec.
(162.00)

10: So: Abstract/summary of Revised "Sorted" on 139.00

What I want to do is program a program version of "Phase 1" (Phase 1 is described in ... of IDP/IA report) I will be using PPM (or a variant of it) to assign pc's to tokens ... in the generation of candidate solutions to problems. The sequence ... discusses why I thought PPM might be very good for this task: how it could do some things that deterministic do. ... PPM seems unable to do.

Wanted to use PPM, 131.05, 130.06, 130.27 R, 133.05 conclusion of problem, 132.27 ff, 133.02, Also not 12.23

SN A (perhaps) different way to write this review: just write it up in "sequence/order" ... Then gradually add comments & refs to comments. ... ideas found in this region of TAI manuscript: 1) What "Phase 1" is: how to just first make candidates ... 2) Why PPM ... 3) How to apply ... to Phase 1 problems: 3 kinds of L search: RETET; U; Random ... Advantages, Disadvantages: why ... was chosen: ... work out proof again.

Vertical text on right margin: PPM to assign pc's to success ...

20 This Macro will describe/retrace the use of PPM-like induction to realize Phase 1 of the Alpha, AS program. We will first describe Phase 1, how it works and how it differs from Phase 2. The PPM compressor and Prediction programs ... are described, and we show how it is to be used in the L search part of Phase 1. There are several ways to realize PPM and ... related programs. We will describe some of them and explain why they were chosen. There are also three commonly used methods of doing L search ... and explain why we selected a particular one for Phase 1 induction.

Vertical text on right margin: ... symbols in fr. Cards. ...

30 Section: What is Phase 1? Phase 1 is described in § 7.2, p 24 of IDP/IA report & Rev 2.0 Oct 2003. Here, instead of ... we will be using a PPM-like system to obtain a probability distribution over trial solutions to problems by generating tokens, one by one. The problems of Phase 1 are all "QA" type problems. This is a fairly general (but not completely general) form of induction problem. The main problem in Phase 1 is always to find a function, O^j such that $O^j(A_i) = Q_i$... (see ... for more detail) ... discussion.)

Vertical text on right margin: ... need not ... after ...

6/22/04

4IM

Mur251 ^{power} user
Vessal LOVELO & all etc. } Muray at U2 hood.
love20

That PPM May be Adequate for Phase 1 & Beyond! -11-30

00:157.40 ::::: I want to see how f. parp. idox (13) related to actual TM problems in which PPM is normally used.
Well, we use PPM to get pd for next token of a string; based on either statistics of a long string or at a bunch of shorter strings (say strings that represent somewhat relatively successful candidates for solving some problem. So consider we are working on a QA problem in Phase 1!

We have this historical [O']_{1,N} Corpus: PPM gives us a distribution based on recent p bet for each token trap, during Cond Generation.

From sense, the PPM is "unbiased": for a suitable corpus (≡ TSQ) we are of any pd. we like; in f: sense that for each a certain set of "features" (≡ near past data strings) we have any pd we like for next token. PPM is probly "consi", for any of a certain large set of stochastic data generators - maybe Finite State Grammars.

6/23/04 A Major Q: If we put in a part of a TSQ that is supposed to embody a certain copy that we want TM to learn: Will PPM enable that copy w. reasonable ESS & cc?

N.B For certain theories, the "Corpus" must include "traces" of previous token solus. - not just a post probs like solus - How PPM would do this: I haven't worked out yet.

Result: It works seem that answer must be "Yes" ; It's component concs. have acceptable p's that pe of b. come to be discovered will be "acceptable".

T. Trouble w. 15-16 is that in PPM all p's of tokens or concs are "conditional" previous token sup. in b cond. In concepts (as in Sol 89) f. p's of concs were "unconditional" using PPM can vastly (or ↓) those p's. I guess the "connet" idea really is incomplete - it's an "outline" of an idea; but "context" is very imp. in assigning p's to tokens.

Which is a Serious Error in Sol 89. I don't know whether importance of "context"

From f. sense, my impression is that PPM may be good enuf to solve very diff't concs in Phase 1.

What "apt. meaning" (3) Get beyond L-shaped & family (no comp factor) optimum p'ch.

Hvr, using PPM, it will not be obvious that a certain cond has by p's & will have to be evaluated by the ppm. Whichever TSQ will be very empirical & intuitive.

Still, PPM seems to be capable of recognizing only a kind of "context". E.g. Source of a problem is an imp't kind of context, yet it's invisible to PPM. We could indicate source by an index on the problem. How could PPM deal w. indices of this sort? → 160.00

4TM

20:15:90

Note 166.11ff on how to use BW when to discover "words" in a corpus.

run into trouble when we define > 1 agent, because of overlap of agents - so we have several poss. parsing methods for a corpus. (In recent 42001 we had parsing problems as well... we wanted to ~~use~~ ~~use~~ ~~use~~ in p.c due to alternative processes) in 'Pur' as well as recent ideas on newest derivations, we do as G. Wolff: reparse corpus after every or every several new derivations.

AZ 191
DAP

57.33 ff tells how to generate "taken sets" - So we first have to have defined some tokens (which, typically, are nouns). We can also use the BW when to get freq of various nouns, ~~and~~ quickly, some can use AZ 191 to derive them.

I'm using AZ 191 in the sense of ~~the~~ ^{the} way to assign PC's of single nouns, that I described in my long letter to G. Wolff.

Is it poss. to use the ideas of 166.33 - 157.08 to 1) define nouns of longer words incrementally (as envisaged in original ~~part~~ ~~part~~ ~~part~~) plus derive nouns? If we could do that we'd have well on way to CFG discy, or even wd. CFG discy: NGMDS and good enough to do rather good probn.

T. main "New" ideas 1) use of BW when to speed up search for tokens ~~and~~ taken sets.

2) frequent re-parsing to get at least at least one good parse. 3) T. idea of hillclimbing for incremental changes in ~~(word~~ ^{to keep} ~~tokens~~ & ^{to keep} ~~agent~~ ~~tokens~~).

An add. formal heuristic: BW when enables us to find nouns of by freq. That have a very idiosyncratic Agent set of tokens following them. ~~Some~~ nouns of this sort are ~~also~~ ~~a~~ ~~candidate~~ for a newst data. Also d.f. of tokens following such a noun is also a "symbol set" candidate.

I've written various times on ways to get newst-candidates for PSB-discy: these methods would be speeded up by using the BW when.

204: 156.29: T. reason this is very imp.: that I may find a way to use/modify PPM so it's a much more complete/universal system. (Also, I will understand better what's going on.)

Say, the operator α reverses order of ~~tokens~~ ~~tokens~~ ~~tokens~~ ~~tokens~~ less than tokens, so that if we use RPN:

$$a, b, \alpha, F \quad [F \text{ has 2 args; } \alpha \text{ has 2 args.}] = b, a, F.$$
$$= (a, b, \alpha), F = (b, a)F.$$

So say we had sequence a, b & want to

to discover that a, G (RPN) is v.g. - even when (Drop that exact line)

Say PPM observes given a certain PD, based on observed covariates betw certain contexts (≡ nouns in P.c. case) & following tokens. For each context C_i , ~~the~~ PPM gives

a d.f. on tokens = $P_x(\frac{1}{N} T_k | C_i)$: Say G is a function on tokens G_i , that can permute them - thus changing the P.D. in this way - i.e. we get $P(T_k | C_i, G)$.

Which is a quite different d.f. on P.c. T_k than given by

We may think of G as a "d.f. modifier": a common example is a Time:

~~Input~~ P.D. on input \rightarrow P.D. on output.

(Spec 158.00) \rightarrow

30

TM

PSG DISCY.30 to William O. Feldman. 157.23 Context

10: : **[SN]** On Phase 1, PPM Phase 2: ... Intuitively, it would seem troubling kinds or troubling, cases, should be learnable in phase 1: It may be expressible in U (Universal) lang that Phase 1 uses, but may not be learnable ... via PPM.

On the other hand, PPM ~~can~~ can only observe recent context, clearly not all kinds of "contexts". — which is ^{main reason} why we have to go to Phase 2 for very subtle behavior.

10: **[Also, phase 1 seems unable to really understand what optimization means]**

So .00-.05 suggests some confusion in what Phase 1 w. PPM can & cannot do. & whether Phase 2 is necc: or perhaps number 2 is — is Phase 1 very? { Could it I just derive an (universal) grammar for PSM's & try to correlate/relate problem sent to f-set of PSM's? (This task is a kind of QA induction ... to Phase 1 for first approx.) }

Ind. intercepts:

145.17 PB → PBCC
Not solved. 200760

145.16, 145.18
146.17

Speeds for 2-3 trees
Many sizes
147.02-14

[Alike] 147.28ff
ANN (5)

142.28, 143.07, 12
[EQ] 152.30 good!
[53.38 - 54.12]
137.26

Let prob copy 154.30
IND Component Analy 155.00

To argue against PPM: Consider a set of probs that are similar (in + sources of lang) solvable by 2 or f. Given $z(p_m)$ — say these probs are of f. form $x_i k y_i$ where k is common to all one of f. problems, but x_i & y_i may vary between problems. As is, PPM couldn't notice f. "similarity" k , because of f. barriers to z 's adjacency" x_i, y_i . I.e. ⁱⁿ the problem, soln, pair $(x_i k y_i, soln)$ ~~is~~ **T. context "K"** has no immediate following in an of by likelihood. K and soln are linked, but not adjacent.

Note. If we could somehow generate f. next $x_i k y_i$ then I think this would Solve Diffy don't think

I certainly had been thinking of nexts (e.g. days → slow → p'd's) — But I thought that I even found a way ^{to} extend PPM to nexts or p'd's or d'leups or slow. → See .30ff on latest "Next Rev".

PPM is already a kind of next — albeit more like a slow, because f. elements of shorter length have less/more wt. One way I considered: Consider all contexts (preceding) f. next α , ~~is~~ This is a kind of slow. We can context cont. r(ph) a next of random strings, so we get f. next: nexts, ~~to~~ ^{to} nexts — which means next can be followed by some random string. If the next next can be followed by some fixed predn, with a distance next a P.D.

Apr. .00-.10 is f. Main Argument: I poorly know how to get it straightened out!

Just "Thinking Out loud": We can perhaps simulate a Univl. D.f. using PPM if we consider PPM weighted sample of universal instructions. → (157.24)

[SN] As SS → ∞ perhaps PPM "must" (2) approach f. universal D.f. is in this sense it would be "Consi".

I think this is what I've been mulling about: PPM defaults a convention of reps on "passive" symbols. Here, "active" symbols, makes it poss. for certain symbols to modify f. meaning of previous symbols. Say "α" means "active symbols ≤ next". We can take: corpus using α followed by selection of which of f. it is (PC will be α ~~at~~ freq. of next subseq within α accurrences). So when context corpus α get its PC f. by using next α . Next we can kill (delete) on additions & deletions of symbols to α . This mile to do very rapidly if we first allow input corpus! — since easy to find all occurrences of r, v, s, t. We can evaluate corpus PC's next (4 f. families) while not spends up at least 141, we still 157.00

4TM

→ Several Tutorial Plans of Web: also see PS D: PS 6/17/04

00

PS ICA: Indip Component Analysis May be of much interest. It's normally "Linear" (so normal ICA is "Linear")
 Joergen Sutterlin (i papers books...) and "non-linear" ICA w. ANN.
 Also look: J. & H's paper on "Fast Minima" for ANN. ... They got some nice results.
 Looking: Joergen's & H's website web sites for papers. There is a report on "Fast Minima" in pseudo code! I don't know which one, kvv.
 — crazy ask J. : It may be a super up "Factor Analysis" ? Est like Zato code detection?

06

Each set of points is 2 pt in N space: we could want to power a reduced no. of dims.
 The N. no. of channels "diagram" is large, it's still dimensionable to PPM. We may have to do modified "escape" since changing PPM ~~etc~~ we will often observe chars that have not been seen before.
 We do, however, have a (arbitrary) metric between chars, so if a char has not occurred, a "near by" version has. A method of searching for "near by" ~~etc~~ strings of chars, ...
 One way to deal w. problem is to quantify signal so as to minimize ~~etc~~ bc of signal.
 We can do this repeatedly for different "error" levels (i.e. rms error) for i. LPC code.

10

There is of course an optimum mix of ~~etc~~ (LPC code + code for error) i.e. MDL

~~etc~~ It is may be way to get "most accurate" LPC code — which might be good ~~etc~~

to Use PPM

8

PS Move to ≈ 3 Major problems in getting GP to be "practical"

- 1) Move efficient search over poplms for code (PPM = ~~etc~~ like "good try") Phase 1
Phase 2
 - 2) What to use as initial population is an O2 prob. (in advanced TM — perhaps in phase 2?)
 - 3) How to recognize that a given problem is "similar" to an older "solved" problem.
 - 4) To translate ~~etc~~ formal "well defined problem" into an initial population ~~etc~~ problem sol lang ← This could be purely from "Microscopic Content" → it's not visible.
- ~~etc~~ a suitable ~~etc~~ (i.e. perhaps a suitable mut/cross algms).

20

2,3,4 are closely related (may be identical).

Anyway, after TM recognizes a new problem z to decompose z_1, z_2, z_3 , resp. ~~etc~~ to past prob, ~~etc~~ prob_{1,2,3} — we can look at i. populations of these

30

3 problem problems at various stages of completion (cell but test... = no. betw 0 & 1 $\approx \alpha$).
 So maybe ~~etc~~ mix i. ≈ 3 populations w. $\alpha_1 = \alpha_2 = \alpha_3$. Actually, we have 2 params for each ~~etc~~
 of $E_i \approx 3$ populations: α_i & its covariance α_i f closeness of population w/ α_i
 We must somehow mix i. ≈ 3 populations ≈ 6 corresponding params to get ~~etc~~ an initial population for i. "new" problem.

- 1) initial ~~etc~~ of i. input problem's "into" \rightarrow A lang for expressing soln
- 2) an initial population \approx possibly a mut/cross or PPM or some other way to generate code.

More Generally: how to take more prob decn into an appropriate PSM for it.

6/14/04
4 PM

Very General Approach to "Phase 1": (153.38-40)

LINEAR Predictive Coding

Spec

00: 153.40: Given PPM, a TSCQ is a massive CJS for Out: we have for hill climbing problem of reducing to CJS. The kinds of possible moves into climb are like 153.38-40. Try and
Way to train f. frame (among other things).

Actually, T. Good about "Hill" is t. \leq CJS; for: entire sequence. Say we keep final problem in t. TSCQ constant. We could change: scope of probs w/o modify PPM or we ~~use~~ method of induction much different from PPM (say GP).

If I use PPM, I may want to ~~use~~ "BOOT" TM with a set of "free" solns. to a large set of probs. This is because PPM may not be able to do much w/o a large SSZ.

If so perhaps it will always learn very slowly: needing many examples of a conc. before it can usefully use it in prod. Corresponding to J's OOB results, if I wanted TM to learn recursive functs, I'd put lots of examples "Boot".

Re. 06 ff: If PPM is slow, - well it's probably still much faster than normal GA.

T. Q is: would it be as "Creative" as GA, or would it be too "Elitist"?

Also I suspect that local Extrema using PPM for GP probs will not be a problem... but we shall see. 22

13: 129.40

~~BN~~ I had idea of applying PPM to the BW XFM since
b. xfm corpus has some info (except for PC = 1/2 length of xfm), any compression would be useful concept for additional $\frac{1}{2}$ factor! Well, I believe that successive use of WB xfm does not add extra $\frac{1}{2}$ factors! They merely permute t. corpus: so we only have to use $\frac{1}{2}$ factor once. If this is true, it ~~is~~ makes it much more likely that t. ~~is~~ subseq. WB xfm will give some extra compression. (Try to find previous ref.) Sec 154-160

If t. first BW xfm shifts corpus by α_1 , it becomes BW xfm shifted by α_2 , then to get original corpus, we need only shift $(\alpha_1 + \alpha_2) \bmod N$ - which is of PC, $\frac{1}{2}$.
I'm still not really sure of this, huh!

22: 112

~~BN~~ PPM didn't do well w. "Short Corp": Perhaps for a long time (many probs in TSCQ) The per problem will continue to \uparrow . - But till it has a large amount of corpus; then the per problem will become more "constant". I also considered (06) of starting TM w. a bunch of prob/soln. pairs so PPM would begin w. reasonably usefull set of per lines, so t. "flattening out of 22-23 would occur sooner into TSCQ.

But, I think t. main idea in TSCQ writing, is to really know reasonable hints for each prob. Ideally, these hints should be (read by TM, ~~or~~ if not directly, used as few "hints" as poss. to get TM to acquire t. hints w. acceptable cc. We necessarily need TM to be able to discover hints of all kinds.

30:

BN

Linear Predictive Coding is "similar" to PPM in that both xfm original signal into a signal of narrower band.

Also, both use "recurrent context". I guess LPC is in 2 parts (1) The set of slowly changing coeffs (2) t. error in prediction. For lossy compression, t. error signal may be omitted.

LPC as is would probably be poor for speaker ID. (Hrr, say we use 4 coeffs: This is 2 bits in 4 space. These "pts in 4 space" would be characteristic of a speaker & could be so effectively coded via PPM. So, each of t. 4 coeffs is xmt'd digitally w. its own precision level.

It may be worth to use $\frac{1}{2}$ PPM on the set of "characters". Each point in 4 space is a scanned in t. language. To denote profile of these "chars" in 4 space could be one characteristic of a speaker. - But sequences of characters (PMM) would also be useful.

4TM

TSQ 12

Charles Bloom has written about ways to improve PPM... I have implemented many of them. I have 2 relevant papers by them.

"2 kinds of PD's":

PD1 to define Game for QATM PD2 "guide to reach"

explained in 3TM 424.00

3TM 394.22 Was early version of problem

TSQ v.g.

3TM 425.00

.. 419.28

.. 403.00

1) 402.12 on various kinds of "reaches".

3TM 394.22 on "colours"

Spec
00: 152.40: Making Conc. nets for solns w. "logical reasoning": Logical reasoning "maps" / "forms" & even problems into a problem in logic, then solves it as a logical problem.

So Given problem: To try to map it into a logical problem is one possy. (We want to use context to make Reiz more likely (if context is appropriate). If Reiz is successful, we continue w. that soln. If not, it doesn't appear in the conc. net. (only successful/usable operations occur in conc. net).

- 1) also Note 158.17!
- 2) Conc net has to be "denser" more similar-like learning. CJS better. Conc has to be
- 3) We need a faster CPU & 2/3 wait until factors are available.
- 4) refer to 2) We need a "better" TSQ.
- 5) Some time ago, I write about "TSQ Repair" - find it. It was a no! I scanned 3TM 386-452 conditions "TSQ repair"

12: 152.40 I've written lots about TSQ writing: One v.g. idea was to write "ruff's dirty" TSQ is by writing a worded list of prob/solns in actual PPM code.

Then look at CJS's. Try to repair by adding problems, modifying TSQ, improving context dependence of tokens. When I finally do get to TSQ writing, perhaps read's review last few yrs of TM writings & try to summarize impl. ideas.

One complaint about induction TSQ's is that they are not readily progressed to induction problems. For, induction seems to be easy to find reasonable probs (e.g. try the data, etc) for it, & it would teach me more about TSQ writing of conc. net construction.

[SN] On finding Schemes to solve QATM! One easy way to get a pd: have O(Qi) give 1 (or a few) outputs. Then we have these outputs as "clusters" of E (letters, w/ a "distance function" from cluster center that determines U in PC. We may have some non-normalized. The induction function itself could be a clustering alg. in which, in the past, the alg. has been successful in finding sets of cluster centers & assoc. distance functions to describe solns of particular QA problems. - i.e. input = Qi, output = set of centers, & assoc. dist. func. The distance func. can > 1 but $\sum_{all} A_i$ to have PC > 0; - which is very desirable! (But other assignments need to be given PC = 0, was a profound criticism of Bloom).

[SN] on "2 kinds of PD": PD1 ~~is~~ f. Game for QATM. PD2 is ~~an~~ a function of how many QA's were done! Its ~~the~~ an attempt to get a PD assoc. w/ corpus. Qi, Q2 ... On. Essentially a SUMAC. In addition to that Corpus, it QA's & any external "contextual" info should be included. It's TM's PD for $\frac{1}{n!}$ instead of all of it. I'd bet it's ever been done. With corpus Qi, ... On: PPM could give a PD on $\frac{1}{n!}$ On. T. Qi's: would it have usable CJS? Could we repair it? TSQ? to give a scriptable CJS? Could we "repair" PPM to give usable CJS? That TSQ may be beyond Phase 1. Use an entirely different TSQ! See. of a "TSQ repair" (Spec 157.00)

4TM

TSQ: .30 Jump



0: (150.90) ! To start w. struct ~~code~~ ; w. N \geq set of N problems; ~~showing~~ if O solves m out of N , it get $pc = \frac{m}{N}$ which is a somehow useful hill (Pro it considers ≥ 1 probab to be of $=$ value). In next step, it may try to manipulate ~~reorganize~~ probs where O didn't work & find corrective steps. Or find corrective errors on O 's prob work for all (or more) problems high below. This is perhaps same of N to what I consider in the latest version of IDSA report: early chapter on QA / map.

Hm, if we use
-i. Current QA
Gave; then
any outside problem
gives ~~error~~
 $pc = 0$ or
(in case $= -\infty$.
If we assign
 $pc > \phi$ to unresolv
probs, we can continue
to Hill climb on
other problems.

Put less "intentional sheep"! I want to go for TSQ writing part as scores best. - since this seems like the most seriously diff part of TM design. T. Early \rightarrow Qs can be practically anything. They are mainly designed to teach me how to write TSQ's.

OK: Back to 139:00 outline of ~~review~~ review.
I may want Ray typed up to refer to in future. Like IDSA report:
So to start w. descn of Phase 1 of TM. \rightarrow also Phase 2: How Ray differ! ~~Give~~ Give
(narrow) refs in IDSA Report.

T. present most direct goal is to use \approx PPM methods to assign pc 's of tokens in Phase 1, L such. for QA problems (both d'funcs & S'funcs).

150.10 for S'Inductive! At first use Bern inductively (freq concept).
So induction becomes problem of finding easily defined sub seqs. (not really arbitrary sections of corpus)
Section is "Alphabets" of tokens for Phase sub seqs. so one can do Bern. induction,
I.E. we want to be able to "partition" the corpus into more easily codable/predictable parts.

150.26 \rightarrow The idea of set of ~~con~~ (PPM type) concepts is of interest. It increases size of the data for
predn. One way to get such Ngrams: Consider the set of all (previous) contexts justward
followed by the token, "c". Put them in (backward) lex order. We may be able to find Ngrams
as "clusters" of prefixes.

30: **TSQ:** It would seem that writing TSQ's would, by far, be the main problem TM.
(This doesn't necessarily make it a diff problem, hm.). One writes the TSQ & its
concept net. From this, it should be poss. to estimate CJS's of solns to problems.
Using more sophisticated Phase 1, L term "Guiding Pd" will perhaps quickly ϕ
 pc 's of tokens that desc. solns. - But after the TSQ is written we can then look
for contexts (150.10 - 26; 152.20 ff) ~~some~~ \approx examples of contexts that
would be helpful for the TSQ solns. expected by Trorer. The \rightarrow 3 should fall
on how to document/verify PPM to ϕ pc 's of solns.
Spec.
 \rightarrow 153.00; .12

ATJM

139.00 is start of review! Some imp. ideas to inclusion "redundancy" 3TM 420.20 starts w/ a note: redundancy "problem"

02:139.06 R T. Q of to what extent PPM is as good as (or better than) the use of definitions. (Ref: 139.06R)

At that point, I felt that definitions summarized a lot of the use of PPM to "rewrite" a corpus, saved much time in looking for higher order regularities

It may be that PPM is as good as or better than defs, if one only considers "one layer": But defs ~~make~~ of ngrams (= "words") make defs of Ngrams possible... which is an imp't "Hyerorder Concept".

3TM 403.24 discusses "2 PD's":

I don't remember what they were - but I think they were very imp't at that time!

2 passy: 1) In Louchter Phrasal induction:

T: guiding Pt is usually a fairly simple 1d, like AZ or OOPS... very Bernoulli-ish:

BZip2 is a more sophisticated passy,

2) could be a "true" (incomplete) PD implied by some vnc. or equivalent.

Another thing I discussed much in the past was "What is Context" I think the final conclusion was, (but maybe not):

It amounts to "How is the pd of the next symbol (or the pd of the next (or distributed) symbol), dependent upon the context of the corpus (or the prediction problem) PPM uses a simple freq. model of probab. - w. some unneeded refinements for "symbols that have never occur'd before (escapes). But, in general, the Q of "What is Context" is the same as the general prediction problem

How can we reduce PPM's concept of "Context"? - using ALP, say.

AZ-14 is a sort of way to do this. Using definitions, Prem Ngrams, Expansion into Context.

I don't have a grasp of how to look at the problem in a good, general way.

PPM takes Corpus (usually) to \vec{L} . Can I think of some, better, ways to predict on \vec{L} ? I've written

Some on this last.

6.12.04 Originally, the idea of "Context" being imp't arose when I was using simply the freq. of tokens to predict the next token. "Context" suggested ways to get better p's for tokens.

But in general, ALP gives the best ways to get p's of tokens is its most general way to get p's.

"Recent Ngrams" is the most restricted idea of context... this is used in PPM

I think that this is modified to consider the size of each context & how often it gave a better pred. (i.e. how sharp its implied P.D. was - also was it deterministic (i.e. only one continuation ever observed).

For some ways of using contexts, they observed that using longer contexts \neq entropy/symbol. This must be because they were not using the right way for

evaluating a set of long contexts. I looked one aspect of this in my analysis of

"One Shot Frag" in AZ-14x (see also letter to Wolff). -> Another tack is 5.2.20

[5N] Most recently, I've been considering having TM start of on MTM problems, then go to what seems to be: NMTM. However, maybe better to start w. Stochastic (= MTM) at the beginning, because ~~one is a stochastic~~ "out kill".

In many open prbs, getting on a reasonably steep part of the hill, is often a very difficult problem.

We could use a MTM seq but a start for evaluation - so you wouldn't have to get all the problems into to get a good score (actually any score at all!). In starting

w. a true MTM-based score, I had considered the problem of moving out to the structure score. I don't know if I ever found a good way to do this!

would starting w. structure score "imply" be better? - or would I have some difficulty transitioning to NMTM problems.

→ (Spec 52.00 52.16)

474

o (5000) (148.1) cont.: I have been considering competition among cands: Could they not cooperate (also/instead)?
 1. advantage of Cooperation is that Puz furnishes a mini TSG. Could we obtain such a Puz w. Cooperation? (Also Note 148.30)

So: Q's: ① What are good algms for competition or cooperation in a Community of Cands?
 ② Could such a community "live" in the internet as "worms"? If so, how would they deal w. less ^{benign} (benign) worms/viruses/spyware/trojans?

As is, they ^{if suitably designed} "Tiera" cands could, perhaps, live in any computers as "viruses" & have offspring, but it would send out over the net. Here the goal of the system would be that of "Real Life" - i.e. no goal at all - just reproduction as a "sort of" goal. Children would not be sent out very frequently (since Puz would be noticed by web admin) - so much time would be spent desiring ~~each~~ ^{each} child to send out on the net.
 Ray specifically avoided any Machine Code in "Ancestors", to avoid .off ... by using code involving an interpreter Engine. But no other computer would have.

So Ideas I'm working on now:

- 1) I want to write summary of my recent work on implementing "BZZ" to do GP's Phase I of TM
- 2) I want to write more on future of AI for talk.
- 3) T. idea of what kind of social activity to have in a Life community designed to solve ~~the~~ external probs / TSG's. How to implement a prob solving "culture" (147.28ff) → 168.00
- 4) Perhaps how to get Puz done on web re "benign worms" or just voluntary time on "Unused Machine Cycles".

$$e^{1.2} = x$$

$$1.2 \ln e^{1.2} = x$$

$$1.2 / 1.2 = .10 - .05 / 1.2$$

Real 3) Perhaps give more detail look at T. Ray's work (also partly M. Levy's "A Life")

24 Perhaps look at Ray's writings (see ~~the~~ his website in my "Bookmarks" on NP500). 5/24/04

25: 147.30: Another direction of the community Ming Li (et al) has a psm "Pattern Hunter II. PS C:\PS\..."
 This detects, counts homologous ^{two} betw. DNA ~~seqs~~ or prot seqs. These are "similar sequences".
 Below. Paramⁿ measured by match, mismatch & gap scores. They say they can do this very fast.

Q's: just what do they look for & to what extent would Puz be useful in inducing patterns in the pps that denote "useful functions"? T. Q's are re to Puz for its utility (or induction) of "Regular Expressions" ← (which even, I think, related to fractal state seqs)

↳ This psm may be able to detect similar sub seqs in same order in different cands.

ATM

ALIFE (47.256 - 148.40)

Spec 147.40: ATM A crad. has to be able to solve a problem himself, before he can present it to another crad.

We arrange "problem rules" so that this is a useful challenge. e.g. One crad just mult 3 (x-2)'s together (using random 2's) & ask for a solution to resultant polynomial. — Somehow, the Challenger & Challengee have to be on same level.

Crad can get "money" for solving top level problems in "problem pool" of 147.30.

They can also get money from one another by winning problem challenges.

A crad. has to make his problem look easy, but be hard. A challengee has only a certain amt of time to solve a problem.

There is also big Q of how do we generate children? Maybe we just look at total amt of money earned per unit time for each known crad, & we generate ~~more~~ L search-wise,

Crad based on Ric Cooper's (like B22), with weights.

We discard crads from population but are not in for ↑ percentage of earning rate.

We may use ordinary GP method of cross mut for 0.02/10 — they may be better.

Another Q is just how much money to give for each soln — so the amount won is a useful sense, a measure of how good that crad. is in problem solving.

Suppose crads challenge one another. If they know their "worth" (≡ rate of wealth gain) of all crads, they will want to pick a poor crad, since a poor crad is easier to beat.

The poor crad will not want to take a challenge, since he will probably lose. → 149.00

[SN] Computer "Worms": Use of Computers on "Net" to reproduce worms: This could be a great barrier in "ALife": The worms would be rather benign to hosts & use very little computer capacity of host. To be resistant against Anti-Virus SW. They would have to have varying structure — difficult to differentiate from other, common SW. in the Computers.

But the viability is not the most interesting part. How to get them to evolve to better survive in user env, is main problem. The evolution so as to be more harder to detect/eliminate

would be part of it all over evolution. First list various ways worm could be more "curious" / diff. to detect / on desirable operator's. It may be a good detector/detector of other "hostile" viruses. — May look for ways that usability of computer.

→ How Hackers may modify worm so it only does "useful" things & has no agenda of its own — could pass off for "benign"? in env. makes they find themselves? find themselves? (probably could!)

Also NB early (est?) papers on viruses assumed that they do not exist before a properly constructed

Viruses is a legit computer prog. is a "understandable" task. A "friendly" life form could meet & be "friendly" into some in a computer! If they fight one another, they will cooperate, or agree to both (room host because fighting is expensive. For each cooperative has, means saving resources of host — (which they don't want to do) (Absolute 149.00!)

REV: 02 on time to do 2-3 Trup pm net works in PB35 that is very close to what is needed in PBCC; so I don't have to make many changes into the new environment.

So, I should be able to make 2-3 Trup pm net works in PB35 that is very close to what is needed in PBCC; so I don't have to make many changes into the new environment.

General Conclusions about speeds of PB35, PBCC, Pent 500, Pent 4 2500:
PBCC is about 1.3 times as fast as PB35 (145:10); (145:05 vs. 145:10)
on PBCC, Pent 2500 is about $\frac{2.6}{1.7} = 1.85$ (145:07, 145:10)
faster than P500

Not much difference: Perhaps check this more carefully. My impression was more like 2.5 to 30. I expect P2500 to be faster on 32 bit PBCC.

We can use large amounts of RAM in PBCC w.o. slowdown, however.

For time needed for various corpus sizes (lower Bud & outcomes) see 144.18-28

T. max corpus size for PB35 to use normal memory is 7.85k

Using Virtual Memory 2.85 M (but speed ÷ 60)

Using Pent 500 in PB35 & small corpus, it takes 10us to insert 1 symbol for corpus length 10.

Time/symbol $\propto \ln 10$; so corpus length of 10^n takes $10^n \ln 10$ sec per insertion.

A "regular expression" is a compact way to describe a language.

SN

In Unix/Linux

There is a program that searches for regular expressions in files. Recent Gnutella is one kind of "regular expression" that can be searched for.

reg. expressions may be a kind of P5000. (= P500-number)

"Regular Expression":

However many other types of expression are searchable: we should make a project of finding kinds of regular expressions that are

(a) findable rapidly

(b) useful for prediction. See folder "Trees": Near identical paper on "Regular Expressions" - Ray see used 3/15/04 - term. paper as well. 149.25

SN

My impression has been that ANN is good for predictions involving non-linear continuous functions. This Guy in Bioinformatics says best guess for predicting structure from DNA sequences, ~~has~~ has been a ANN pm. - which seems to be a discontinuous problem! - But I would have to know much more about what the pm does!

SN

ALIFE: "Avida" at CalTech, (see website/Bookman) They like this pm (free available)

To create organisms (Eggs) from small no. of bits. Idea is to use ALIFE to solve an evolutionary problem.

There are several 10 devices P5000 "Hideo". Each has own set of RAM. This is distributed. There.

A way I might do it: There is a pool of probs of differing difficulty. Each has a label & rewarding pm

Remove if solved. We solve how strange to making mutations. But there is a pressure to

"Get max reward": which means "solve hardest probs". We may allow creatures

to eat one another, to get parts to make "strong" (= reward getting) children.

It would seem like there could be no point of "contribution", since one normally takes to genotypes of successful cands & mixes them w. genotypes of other successful cands. Contribution wouldn't yield anything better!

T. poss. a drawback of an ALIFE problem solver: T. cands can present problems (of competition) to one another. This is an automatically updated TSQ. A problem is to get problems involved in surviving in the society to be related to the problems presented to the system. Perhaps the social inheritance consists of cands presenting problems to one another.

4TM

2³² 2³⁰ = 10⁹ so PG = 2³².

running Big Mem. Bus in ~~the~~ PBCC.
Dim A (268000000) As Byte. 268M
T1 = Timer
For x = 8000000 to 250000000 250M

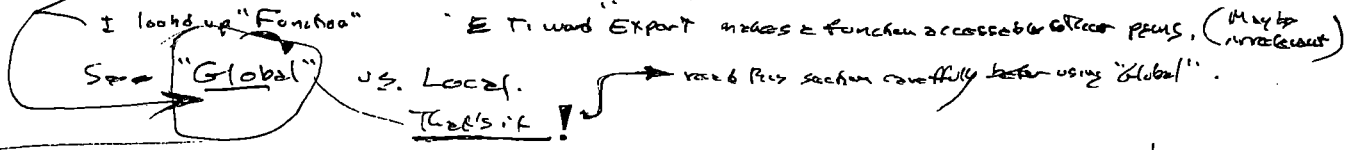
... etc.

It took 105.131 sec. ! ~ 100 times slower than Big Mem. Bus.
T. Diff. betw. Big Mem. & Big mem is rather small!
GH! I did 25M loops rather than 25M - so 10 times as long - not 100 times.
Try Big Mem w. 25,000,000 only ; Also note: Disc does 2100.
10.1065 sec!

try Big Mem: 2.68M only! The try just stand down Big Mem. Bus! (.00)
Use 268M in A(); Big Mem. Bus only 26.8M in A(0)
Say 200M in A. for Big Mem, only 25M loops -
Took 2.42 sec ; so using > available RAM, it used to disc a lot.
So remember: Don't use so much RAM! 200M is probably fine, ~~the~~.

I may get 512 Ram, hrs.

(17: 145.40) on PBCC: "DATA" and "READ" Maybe some what different from PB35?
No such word as "SHARED" is used. This may be serious diffy
I have "VARPTR" function. (but not VarPtr 32' - I guess its ok. - it returns
a 32 bit address. - see VarPtr VARPTR



READ\$ -> Read\$ is only ASCII is readable. See it Read\$ is a way to get numerical data in!
Or use I can write in contents of E. Fast Node. Or how is a file is containing data, Read it up locks.

For setting up the initial contents of the first Node, user is exp. of assignment statements:
put it in a string so it can be read when needed. In PBCC; I could just use MAT
assignment of MAT A() = B() ; where B() is trivially w. pre initial Node values.

T. System may have a Read var: try reout!
(string)

Read(7) is a string of 7th Data object:

This can be converted to a number via Inst like CVBYT, CVD, CVWVD, etc

So no real problem.

Also is STR\$ and VAL (i M K I \$)

Actually Val converts string express. to no. This is what I need to put numbers in
DATA 3 normally.

I think try it to VAL can have decimal pb. like "38.753".
See "VAL" keyword. With 33, 923, 13 over columns,
I'm not sure about columns,
lower columns.

PB35 has Val in many of its assoc. functs: So I can do changes within PB35
I test them, to be more to PBCC. How PB35 does perhaps "Read"; only has Read\$

Speeds of various entire pgms \approx Part 3, 500 vs. Part 4, 2800!

1) To do loop of $A(i) = A(i) + 1$ for 25k times.

The do it. 25k times 4 times each, so 25k operations.

In PB35 this took 3.4 sec $\approx 11.05''$ on HP 500

so $\frac{3.4}{25} \text{ ms} = 136 \text{ ns/op}$ or 68 clocks per op.

2) To do 25 M operations in one loop using $\approx 2.8 \text{ M}$ by array

In PBCC took about 1.4 sec ≈ 24 on Part 2800

so 56 ns per op. or $56 \times 2.8 = 156.8 \text{ clocks/op}$

$$\frac{\text{part 2800 speed}}{\text{part 500 speed}} = \frac{2.6}{1.4} = 1.857$$

3) to do each of .06 on PBCC HP 500 took 2.6 sec ≈ 2500 .

or 65 ns per op.

so PBCC is $\approx \frac{3.4}{2.6} = 1.3$ times as fast as PB35

So it looks like PCCC can do larger memory w.o. slowing down - so we could get

factor of $\frac{3.4}{1.4} = 2.4$ So using Part 2800 in PCCC will increase speed of structures

So conclusion of 144.18 - 28: speedup by $60 \times 2.4 = 144$ over 144.18 - 28

1 M corpus would take $\frac{60}{2.4} = 25 \text{ sec.}$ for 1 M symbol corpus.

Maybe use BAD! $145.18 \times 90 \rightarrow 146.17 \text{ ff}$

Unfortunately, porting of a pgm from PB35 to PCCC is extremely non-trivial (often syntax different and difficult to find). The printout of error(s) into compiler is interesting:

It can give a number like (10,7) which means 7 ~~is~~ word in 10th line.

There is something in PB35 about references to Arrays via STRINS being by "value" rather than "?". I don't know what this means - Part 4, 2-3T pgm does use > STRIN.

2/08 defined macros, so it does work OK.

First Transpare SN 14 (Part 4, 2-3T pgm)

Part 4 SN 14

Just try transcribing the Part 2-3T pgm. Make a copy of what to various SN PGMs are:

On top of ~~SN 14~~ SN 18 v. 6.5

> SN 14 may not have to reinitialize all arrays to zero - so it works once, but can't be

repeated. (w.o. zeroing). [It may be easy to zero an array using "MAT" command in PCCC.

The more obvious difference PCCC: the first line is FUNCTION PBMAIN() (see " " END function.

SW to M 32 bit Dos;
Sort < Myfile.txt >
Sort.txt

In the "Document on line", the first line is always "Function PBMAIN() AS LONG"

My experience is that "AS LONG" is necessary. Try deleting other parts of PGM.

If there is any problem, I think it's easy to add 4 lines:

See: keyword "WINMAINFunction" for alternative "AS LONG".

10 If INKEYS = "" then goto 10; otherwise if screen will disappear soon.

PGM ends.

The technique is to write a PGM, check if it compiles OK. (Compilation may have useful feedback e.g. 18-19) After compilation, it's possible to run the display & goto to "CC dir"; do refresh & double click on the "error" file - usually second from top.

$\rightarrow 146.17 \text{ sec}$

Speed of My 2-3 Trac ppm

The x60 factor seems familiar to me! ~~Just~~ Just what I can do w. that factor, is unclear!

Now in virtual mode see nearby H1 can be.

try $k_y M_x = 5$	$H_y = 10 k$	$k \dots = 5 \text{ sec}$	-30 as output!	$6 \rightarrow 342$
		new $+527$	$7 \rightarrow 1476$	$6 \rightarrow 350$
			$7 \rightarrow 1463$	$7 \rightarrow 1718$

$H_y = 100k$	$7 \rightarrow 1698$	$H_y = 1M$	$7 \rightarrow 1703$	$H_y = 100M$	$7 \rightarrow \text{exp. memory error}$
--------------	----------------------	------------	----------------------	--------------	--

$H_y = 5M$	memory error	$H_y = 2M$	OK	$7 \rightarrow 1465$	$H_y = 3M$	no good	$H_y = 2.5M$	$7 \rightarrow 1700$	OK.
------------	--------------	------------	----	----------------------	------------	---------	--------------	----------------------	-----

$H_y = 2.7M$	OK.	$H_y = 2.8M$	OK	$H_y = 2.9M$	error	$H_y = 2.85M$	
--------------	-----	--------------	----	--------------	-------	---------------	--

So $H_y = 2.85M$ OK. $H_y = 2.9M$ bad. P_2 's close count. $20 \times 2.85 = 57$ Megabytes!

This is \gg the $16.727/M$ by P_2 for (-11) said I had!

For $k_y m_y = 10$ spot 126256. (≈ 121106 or $143,402$) so it may be wrong.

AM It was w 10ms (symbol) ~~in normal~~ using normal (rom) so 16ms / 24 Mbol in virtual memory.

so .6" for k symbols $600" = 10'$ for M symbols 6 minutes for 600k symbols.

2.8 minutes for 2.8M symbols — so I can use key to compare w. 1M corpus.

at 600k corpus.

So I can do ~~stuff~~ stuff in **Colony Corpus**

Conclusion: In normal PB35 memory 7.85k corpus is max size (w/ 43.40)

Using Virtual Memory, 2.85M corpus is possible, but speed is $\approx 1/60$ of normal memory.

Normally it takes 10ms per symbol for 2 corpus of key in 10, 2nd time w. N (N).

So $N=100$	$\approx 20ms$ /symbol.	$N=100k$	$\approx 50ms$ /symbol
1000	$\approx 30ms$	$1M$	$\approx 60s$ (symbol).

So a 1M corpus would take $\sqrt{.6ms/symbol \times 60 \times 3600/symbol} \rightarrow 60 sec = 1 Min.$ $\rightarrow 25 sec$

~~Using Virtual Memory~~: Using virtual memory it would take 1 hr. $\rightarrow 25 sec$

Using P_2 at 2.86 (2.5 times as fast) $\frac{60}{2.5} = 24$ minutes for 1M symbol corpus. $\rightarrow 25 sec$

.6 x 24 = 14.4 minutes for 600k corpus. $\rightarrow 25 sec$

These figures will be mult by maybe 2 for rest of calculation.

50 BUG! Error in computation of Memory needed for 2-3-T for large Corpus!

T. present ppm. Uses words for addresses. This is 16 bits $2^{16} = 64k$ only. For corpus of $> 64k$.

We need 3000 bytes/address. This will \approx halve the size of compressible corpus:

I had 4 byte and 2 words, ~~8~~ 8 word errors so $4+16=20$; now I need $4+32=36$ bytes/character.

57M by $\frac{100}{36} = 1.58$ characters in max corpus

Use of Double words may slow down PB35 further!

Perhaps try using PBCC (console compiler — it runs at 32 bits!)

Time some runs in PB35 vs. PBCC using P2.8. See if there is much difference.

In PBCC 250M by available so corpus of $\frac{250}{26} \approx 7M$ symbols.

150K 500 hrs
300/hr
10/hr:
50 hrs.
w/2 days.

10 : **NB** I was thinking of using **Bzip2** output for estimation of Δ bit cost; but result is only ± 4 bits:
factor of 16, which is not good enough! Looks like Δ will have to (at least) rem. Quicksort or
Heapsort. As is, I understand how 2-3 Tree, Φ sort work. I want to understand
Heap sort also. It may involve **Tree Structures**.

of! 140.90 **142.28 = .90** is a nice "My level" approach to TSQ construction: **GOOD TO KEEP IN MIND!**

35 **Try to see how various func. types, prob types, soln. types can detect this scheme.** (12)

11 **SN** Under what conditions can we use ALP to get product of known expected error
yet not use "Cross Validation"? Perhaps it ~~estimates~~ is entirely a priori:
we hadn't seen data when ppm. was written: **Exactly what does this Mean?**
After the ppm is run on new (previously unseen) data, we have a certain Δ cost
for bit of original data. This is **expected value of future (next symbol)**; bc / symbol. (151.12)

If seq. is assumed to
be "stationary"
Then assign costs
bc / symbol to 2
subsequence probs
a % "unbiased" estimate
of bc / symbol.

12: .05 \rightarrow One way is to write tsq "in English" at My level. Then put in problems and
of TSQ "Problems" (i.e. not yet) pro corpus.
Both in actual code. Both probs & solns. should be expressed in
"factored" form as much as possibl. using solns, cons, defns, that I think are
"Generally Useful". Then each of these concepts ("factors") is either to be
learned by TM viz a tsq. or inserted into it. This insertion **can take form**
of a large no. of (prob/soln) pairs that use ^{these} other concepts, > being
made part of the **pro corpus**, just ppm codes before TM starts (rng).

This "pro corpus" gives ppm lots of important "contexts" & their likely continuations.
A desc of a soln. could include "Use hour #18". TM will have had to correlate f. Various
hours w. various problem desc. Hour #18 would be followed by various choices used in its implementation.
of f. hour. For IM to try to desc hours, she'd - TSQ's of a list of probs + solns,
but could have used ~~particular~~ a particular hour \rightarrow so it could both invent f. hour, as well
as "rules" for its implementation.

- 1. The initial TSQ's are to teach TM elementary skills! Mainly lang techniques.
- 2. Mainly to teach Me how to design TSQ's to get TM to acquire skills of various kinds.
- 3. Ultimately to get skills needed for Phase 2

**I could take to present 2-3 prob ppm, see how large corpus I can do in
power basic Ram; then see how slow it gets for corpus using PB35 w.
My impression is that it's like using a HDD for f. data, using a RAM disc. \rightarrow per byte
acceptable for internal in many "projects".
See ~ 81.00 ft for studies of ~~tsq~~ speed of prob 2-3 ppm. T. ppm. very used**

was EN 18 B. Bas : also SN 18 C. Bas : I think f. array use 20 bytes per symbol of corpus.

PB35 allows maybe 200k for arrays! so corpus = 10k symbols.

SN 18 B. Bas	H1=1000; DIM DN(CH) as word.	So: DYNomic Memory?
H1=10k (out of memory)	6k ok	7.9k out of mem
2k ok	7k ok	7.85k out of memory
4k ok	7.5k ok	So 7.8k is ok.
8k out of memory	7.8k ok	7.85k is out of memory

So I'm surprised to do a 7.85k corpus

In SN 18 C. Bas all arrays are virtual
I want to know how many 100 sec or
more! For K/M X=10
try K/M X=5 output = 15.958 I think it
doesn't work for K/M X=10
For 24 M X=10 I got 12.1106...
2115 W/O V/M 1021
V. 1, ...

IMPT → TSQ's Reference UMCS .28 ff to .40 at least
2.13.04, 12

PPM (objection) 0.02
criticism

10 : p 233 of Num Rec: Says for Quicksort, worst case is an array that's been already sorted!
I don't see this!

02 SN Possible Criticism of PPM! Cont contains that only occurs ^{only} 2 or 3 times in a very long corpus, and unlikely to be ^{get + pred.} predicted! (From the best work on $A \geq 100$) see at the end of
If a second layer of definition is used: T. freq. of symbols should be based on the most likely parse - i.e. not all occurrences of α (where α has been defined) will be instances of α . A certain no. of them (e.g. by offset) will be due to "coincidence".

On Q sort: 1) we can (perhaps) significantly speed up comparisons by starting comparisons of 2 strings on the var symbol. $k=0$ for top & bottom "chunks of base"; for a base chunk one always knows 1 string that is all strings in chunk, so compare those 2 strings to get k .
2) Horowitz's s. get t. expected value of c for Quicksort soln by solving a difference eq. Probably one can get difference eq. for d & get var. of d.t. - which would be very useful!

At first glance, is not difference eqs for expected values & for var., seem identical: so $\sigma^2 =$ expected value! But check this! My impression was that v. d.t.'s

looked like M rather than W

No! Corresponding to 14.12: $S(N) = \frac{2}{N-1} \sum_{i=1}^{N-1} S(i)$. Note absence of (N^2) term (Horowitz may have gotten $N+1$ for that term, but its irrelevant here).

$$NS_{N+1} - (N-1)S_N = 2S_N$$

$$NS_{N+1} - NS_N + S_N - 2S_N = 0 \implies NS_{N+1} = (N+1)S_N$$

$$\frac{S_{N+1}}{N+1} = \frac{S_N}{N}$$

$$S_N = kN$$

$$\frac{\sigma}{E} \approx \frac{1}{\sqrt{N \ln N}} = \frac{\sigma}{E} \approx \frac{\sqrt{N}}{N \ln N}$$

18 TSQ's: A Major problem is to describe (not just / red. / say) so that it means I want TM to use can be easily expressed in Set Language & / or T. lang is "extensible" so the needed concepts are easily definable. In any computational universe (say, I must be able to define (or point) concepts so that 28 is satisfied. It amounts to putting certain my kinds of a priori info into the Red. Machine (e.g. see 38)

28 ff may be the main principle of TSQ, Red machine descrip. T. rest of: "Learnability" of T. TSQ should be implemented via (Bain seq) / (PPM) / (reg).

So I write the Seqs of prob (i.e. possible solns), & get see what the p.c.s PPM gets for my solns. [I don't need to guess it for this - BZ it's itself would be a good info!]
38 One way to implement (35-32) would be to give TM certain (prob/solns) free as part of its (a priori corpus) see 143 of

4TM

One Q is whether I can do w/o 2-3 lines psmi (i.e. just use Hoare sort like BZ1P2).
See if I can find it.

Quick sort (not ~~with~~ ~~bad~~); see ~~with~~ ~~CPP~~ ~~in~~ ~~C:/PS/BW~~ ~~xxx~~ ~~-Nelson~~ ~~code~~
H(x) look at ~~xxxx~~ ~~H7ALL~~ ~~due~~ ~~is~~ ~~is~~ ~~but~~ ~~CPP~~

Seems easy to psm: see Horowitz, S24m1 P.121 for quick sort - also, NUMPROC. P235
I don't yet see why it's needed! It is used int. simple, recursive form of Quick sort.

Rec Quick sort: If partition pt. is chosen at random each time, $i \in \{1, \dots, N\}$ is expected time to sort a list of size N , then $E(N) = \frac{1}{N-1} \sum_{i=1}^{N-1} (E(i) + E(N-i)) = \frac{2}{N-1} \sum_{i=1}^{N-1} E(i)$.

One solution $E(N) = \frac{2}{N-1} \sum_{i=1}^{N-1} i = \frac{N(N-1)}{2} \Rightarrow \frac{2}{N-1} \cdot \frac{N(N-1)}{2} = N$. $1+2 = \frac{2(3)}{2}$

Eq. -08 is wrong. After a random point is chosen, we have to do $\sim N$ comparisons which takes time $\sim N$ add in .08: Assume exchange takes \ll comparison time E_u .

$E(N) = N + \frac{2}{N-1} \sum_{i=1}^{N-1} E(i)$
 $E(N) \approx N + \frac{2}{N-1} \left(\frac{1}{2} N^2 - 1 \right) = N + \frac{2N^2}{N-1} \left(\frac{1}{2} - \frac{1}{N^2} \right) = N + 2(N-1) \left(\frac{1}{2} - \frac{1}{N^2} \right) = N + (N-1) \ln(N-1) + 2$
Hand S P12 & Eq. 3.4 get slightly different eq. $E(N) \approx N + \frac{1}{N} \sum_{i=1}^{N-1} \frac{4-i}{2}$
Guess $E(N) = 2N \ln N$.
K(a) Time of comparisons has depends on ~~ordering of data~~. That's applying to Horowitz's S's proof.

$\approx 0.1N$: consider $\sum X_i \ln X_i$
 $\Delta x^2 \ln x = 2x \ln x + x$
 $\sum x \ln x = \frac{x^2}{2} \ln x - \frac{x^2}{2}$
 $\sum x \ln x = \frac{x^2}{2} \ln x - \frac{x^2}{2}$

so $E(N) \approx N + \frac{2}{N-1} \left(\frac{1}{2} N^2 - 1 \right) = N + \frac{2N^2}{N-1} \left(\frac{1}{2} - \frac{1}{N^2} \right) = N + 2(N-1) \left(\frac{1}{2} - \frac{1}{N^2} \right) = N + (N-1) \ln(N-1) + 2$

so what's $\sum i \ln i$? $\frac{x^2}{2} \ln x - x^2 =$
 $\frac{x^2}{2} \ln x - \frac{(x-1) \ln(x-1)}{2} - x^2 + (x-1)^2 = -2x + 1$

1	3	3
N	4	8
1	5	12
1	8	24
1	12	48
1	15	75
1	20	120
1	25	175
1	30	240
1	35	315
1	40	400
1	45	500
1	50	625

$X_{n+1} = N+1 + \frac{1}{N} \sum_{i=1}^N N_i$
 $X_{n+1} = N+1 + \frac{1}{N} \sum_{i=1}^N N_i$
 $\sum_{i=1}^N N_i = \sum_{i=1}^N N_i + X_{n+1}$
 $\sum_{i=1}^N N_i = \sum_{i=1}^N N_i$

S=1
For N=1 to 100
 $X = N+1 + 2S/N$

if ~~xxxx~~ point $N+1, X, X/(N+1)/\log(N+1)$

N	X_{n+1}	$X/(N+1)/\ln(N+1)$
2	4	2.88
3	8	4.6
4	12.6	5.23
5	17.2	5.6
6	22.4	5.6
7	28.3	2.15
100	937.47	2.035
200		2.03609
500		2.025
1000		2.0225
10000		2.01677

Solver says N $X_{n+1} \approx 2N \ln N$.

This is "Expected value". No Δ^2 is given!
It's "proof pt." is chosen random, it would seem due to initial ordering of data would not affect result!

100 : 350
200 : 360
500 : 270
1000 : 225
10000 : 10.777

100 50
50 30
20 10
10 5
5 3
3 2
2 1

250 - 20x
225 - 10x
167.7 - x
165.2 - 0
2.01652 + 1/N

N.B. Forgo. D. difference Equs can be used to get Δ^2 at each p.d. of cc of soln.

! (SN): For a T.S. of 0's & 1's only: How close is a curve $P(x)$ to max PC curve & pc is Maximizing approx?

I.e. $S(x) = 0, 1, 1/2, \dots$ how sources:

to find $f(x) \Rightarrow \sum_t (S(x) - f(x))^2$ is min. This is same as $\sum_t (S(x) - f(x))^2 + f(x)^2$.

to find $P(x) \Rightarrow \prod_{x=1}^n P(x) \cdot (1-P(x))^{1-S(x)}$ is max.

$\max = \prod_{x=1}^n [\sum_{S(x)=1}^1 2 f(x)] - \sum_{all} f(x)^2$

$= \prod_{x=1}^n P(x) \cdot \prod_{x=0}^n (1-P(x))$

→ Hvr, (03) can also be viewed as a max PC model w. Gaussian D.F. → So check that out!!

HA! So they are both very similar. 03 uses Gaussian D.F. for error. 06 uses a linear d.f.

→ There is some funky parity w. G d.f. since words not mutually static or not.

One "Lower" way to solve 06 is to solve 03 (against approx.) & least squares method + data pts in accord w. 06 & using Gaussian error is simpler soln. Two approx. methods for data pts & use 03 soln. R^2 converges to correct soln. to 06.

Maybe $\ln(X(1-X))$ is more exact. And (07) The (06) is not exact.

SN Instead of PPM, I could use that old \sum or method, in which I only make 1 deduction using the received content. I do this for all of the contents that have occurred before ($SSZ \geq 2$) & use the previous prediction. While this may take much longer, it may be more accurate. It may be able to help get regular PPM more accurate, by falling back to it.

predicts of various context strings. I do have exact formula for PPM's using factorial algebra. Hvr, my approximations may not be correct. Check using Maple. I would like to compare it to PPM.

Actually, it doesn't really implement "deductions" (to create new items from old) any more than PPM does. I am hoping that PPM will be powerful to get over thru Phase II. Another way to get thru Phase I is by GP... possibly using PPM for LSrch — which is a bit like using a "second order" Phase I — The induction method used is universal, rather than limited modeling of PPM. In such a system, the lowest level (TSQ must have problem)

that are (line extrapolating a seq. of functions. I.e. in a QATM, we have a sequence of points that are related to $(A_n)_{n=1..n}$. A func. for each value of n : $f_n(x)$ — We expect this as an "unordered" seq. ... Reo it is ordered. As "unordered" seq. it is a set of approxs to a func. that will work for all "all QATM's".

We want to (try to) rapidly extrapolate a PD on $f_n(x)$: PPM is able to do this rapidly, but much better performance is probably possible. Actually Phase 2 sort of does this.

We assume that Phase 1 has found to solve problems that are Phase 2 meta problems.

Phase 1 could just as well apply such methods to its own ~~own~~ own problems, of various kinds. Thus at 1. beginning, Phase 1 will simply solve problems — to train Q_1 or well as TM. I will learn how to write TSQ's & TM will learn to solve Phase TSQ's.

REV

- At least 2 imp't items
- 1) Implementation of PPM & use of the L such
 - 2) TSB writing (is how a character by @)

201.138.40

Ok: Some things to include in the "10 of 10"!

1) Explanation of why we decided initially to do 2-3 trees ppm!

- a) $N \log N$ time.
- b) It may be available $\approx 2 \times 2$ done p.p.m.

for a small corpus

- a) It could do insertion
- b) " " " deletion

Apparent faults of 2-3 trees:

- 1) My p.p.m. uses at least 20 bytes per symbol of corpus, is 36 bytes per symbol for larger corpus. (144.29)

06

1.5) Discuss recent work on how 2-3 trees may be done: 138.13

2) Perhaps initial discussion of "PPM": why it looks good!

Initially I did it was in Costa Rica.
Nov 2003: 3TM 42.13 mentions B22 as
pass: 50 < B22, 3TM: 405.3) str. 405.12
3TM 41.50 ... 5 bit. rep. ... 150.02

Q: just what were the nature of the "improvements", in terms of sequential produ. of \tilde{L} problems of computing normal forms & expressing kernel

2.5) The kernel approach: how to update to kernel cheaply: how to ~~get~~ get individual tokens cheaply

3) Discussion of 3 kinds of L such ~~T < I~~, Π , (random)

Why Π & Mc looked much better at first (factor of k or $2k$ where k is length of s.d.m in tokens)

Conds under which this factor does not occur.

So discuss conds under which $T < I$ seems best.

4) The "frequency of updating" problem and T. possibl. user of Heapsort (?)

rather than 2-3 trees. (Initially, I may want to use 2-3 trees)

Go thru recent work & see if forey do indeed cover imp. ideas.

14.30
15.20 of TSB Report

Also go thru recent readings of ideas on **TSQ's** - partly modified by recent

ideas about PPM as a cheap ~~low level predictor~~. low level predictor.

Discussion of adequacy of PPM in this app'n. 136.02, 132.20 is imp't idea (PPM does ^{not} have to be Universal)

Q: Improvements in predicting \tilde{L} sequences. Can all improvements in predicting \tilde{L} be affected by ?? Perhaps ~~not~~ because \tilde{L} has just about all info about original ~~corpus~~ corpus.

It takes very little extra info to go from \tilde{L} to corpus.

Could we apply PPM to \tilde{L} w. any signif. compressa? It would be easy to apply PPM to \tilde{L} sequences!

It would be really easy if I could get the B22 code & ~~then~~ find its parent ~~code~~ \tilde{L} .

I have no idea as to whether this would work!

Maybe read original BW paper: 58 had lots on mechanics of sorting.

It may be that .23 is sensitive to just what ^{improvement} variety of PPM one is using. For some, recursive may do nothing - for others, it may do a lot, others it may have negative effect.

No! on second thought, all varieties of "PPM" do the same in that they convert corpus to \tilde{L} .

T. differences on how they try to predict \tilde{L} - so for several recursive at .23, we only end up using the "different version part" of PPM once!

We can try various methods of produ. on \tilde{L} : linear's non-linear (as mentioned before). Even try Lempel-Ziv!

33

Actually, I'm not terribly optimistic about recursive PPM, but it's so easy to try, that it certainly should be tried! So when I have the \tilde{L} generation p.p.m., I will want to try these tricks. (.23 - .33) \rightarrow HVR, see (54.13)!

about who reads / and? T. Neplatech mechanic can be given problems of a dirty
in "solve", "interpret", "simplify", "pick", etc

SN Back to an earlier Q on the "adequacy" of PPM for L-subs: As an extreme case,
say we used full Unl. def for L-subs over success in problems in TSCQ, using some of
previous problems as corpus. This would recognize all poss. tags in corpus & use them
to prep up L-subs. By using a less optimum P.D. to Guideo (erch), we end up w. solns of
smaller pc : i. larger such times

T. m. in. iden. of TM, = PM₂ was that no could somehow make a problem of im. proving t.
Guiding P.D., ~~problem for st. im~~ a part of: novel TSCQ. The phases → phases 2
is one way to try to implement Piv.

SN

on found for 2-3 trees & its large need for RAM
If we put stuff in lex order cheaply, as in BZCP, say, ~~could~~ could use the L &
corpus for prediction of individual tokens wo. having to insert any pnyms into Γ . wo
would have to know where to insert & p. tree (or context strings, hvr. It suffices were
in actual lex order (by address of), we can ~~use~~ do a binary decision
search taking ~ logN decisions.

T doesn't know sort = Radix Sort.
This uses some RAM, but not
near as much as
2-3 Trees.

say we have 2 tables of N entries long. T_1 & T_2 . $DM T_1(N), T_2(N)$

$T_1(x)$ gives the Lexical order of x . R_1 -shift of t corpus.
 $T_2(z)$ is inverse of $T_1(x)$: If x shifts t , shift no. of t is R_1 (lexical context).

To find where z contains y (its order no. will usually be a float \neq integer):
we compare y w. shift no. $T_2(\frac{N}{2})$. If its lower, y have comparison no. $T_2(\frac{N}{2} + \frac{N}{4})$
... " " " " " " " " " " " " $T_2(\frac{N}{4} + \frac{N}{8})$

So we find y 's position after $\log_2 N$ comparisons.


Given Table $T_1(x)$ we compare table $T_2(z)$ in N steps.
for $z = 1$ to N $T_2(T_1(z)) = z$: next

periodic updating of Γ

The process could be speeded up somewhat by noting how many common
symbols occur in z 's parity. This would give idea of whether
more should jump, say to $A + \frac{N}{16}$ or, if it is closer, to $A + \frac{N}{8}$ or if it were

Consider Lempel Ziv: compression $\approx x \frac{1}{2}$; v.s. Bzip2: $x \frac{1}{3}$. For ASCII: $x \frac{1}{3}$ is 4 bits/symbol.
 $\frac{1}{3}$ vs. 2.67 bits/symbol. ← Actually its more like 2.2 bits/symbol. Any text is
 LZ may also be better so difference is maybe $(4 - 2 \frac{2}{3}) \times \frac{2.2}{2.67} \approx \frac{4}{3} \times 1.1$ bits/symbol.
 For a 10 symbol pgm! Plus is $x \frac{1}{3}$ ~~2.67~~ $2^{10} = 2048$! So it does pay to use ~~it~~
 Bzip2 rather than LZ. ~~There~~ no punishment on Bzip2 maybe may not be very
 significant, hrr. — I'll have to try this out.

In the category corpus to compression factor for LZSP codes wasn't very good, but ~~my~~
 "Lisp"

"Lisp" will be for (I hope) a fairly redundant corpus!


It looks like I will have to pen. PPM for generating corpus; any way — Re: C mite get

so useful. Feed back from PPM (Bzip2) pgm.

So Best: write ^{down} ~~revised~~ of how I will use PPM to do LZSP. — why I will use TAZT,
 Then write it in PBasic. See if ~~it~~ ~~can~~ accommodate large text files

is so useful — It ~~is~~ ^{is} ~~not~~ ^{is} possible to get PB to accept large files, but it slows things down
 tremendously! Hrr, I may be able to tolerate this: i. only large files I'm interested in

are English texts; If it takes a day or ~~two~~ ^{week} to do 6i file... still ok.

For testing Supus. I probably will have to do files of $< 1k$ symbols, ~~with~~
 But to test PPM on GP, I may need ~~some~~ ^{some} files of $> 1k$.

NLP: $\sigma \rightarrow \sigma$
 easy to show for $k=1$: It
 not sure how to do it
 $k=2$

SN Why was $\sigma \rightarrow \sigma^2 \frac{N+k}{N-k}$ not a logit soln. to "How Many Gifts?" problem?

J. "Proof" was that v. factor $\frac{N}{N-k}$ gave direction from model using correct cuts. σ^2 if present
 cuts were correct. That v. ~~additional~~ ^{additional} $\frac{N+k}{N}$ factor was because cuts found were not correct.

— The Bzip2 eqn was true for non-linear stuff, too, but "N was large and — so ~~it~~
 SO system was "locally linear" (would large N do this? — or would small σ be needed?)

Re: .29: This seems closely related to PPM: we have various "Match lengths"
 corresponding to k : we could use "backshifts" (\approx MDL) on wtd sum (Univ. D.F.)

Approaches to "M.C.": ①. 29.5 ②. Gatspirt from previous corpus ③. use marks
 like PPM (which maybe ②)

TSQ's: One simple way to start: write a bunch of sample problems, & write their
 soln. pgms. ~~Not~~ ^{Not} necessarily written w. ~~each~~ ^{each} in mind!

This is the start of T. TSQ. It gives an initial P.d. for L-section subseq. problems
 But may be more carefully designed (scheduled). We can use PPM for our ~~TSQ~~, & its continue

It would seem that this is similar to Rex PLIT "Guys" pgm that lessons to pgm. "by example".
 In ~~my~~ ^{my} ~~case~~, all previous correct pgms (both ^{for solns} given by instructor or solved by TM) are used as

"examples". See how his pgm worked. It may have used old ideas. He has ~~a book~~ ^{a book} of
 papers by various people on PPM.

Re: 136. (5 of on 129 1/2). It is certainly not clear as to how just how we ~~would~~ ^{would} want TM to
 respond to problems: 129 1/2 ... How does it know what is wanted?

A kind of T would be to learn to respond properly to questions asked of Maple 2/3/4/5/6/7/8/9/10/11/12.
 So an iter. item would be to compare to "problems" of 129 1/2 w. PPM of Maple. Why is 129 1/2

so unclear some Maple "well definedness" of Maple's problems.
 items: 1) Evaluate a function either numerically or have done by literal's give values to param's,
 2) specify what characteristics one wants in "simple" —

3) solve 4) prove 5) disprove 6) minimize (wrt param set) — usually continuous, but

distal = Bochl
 medial = Inval
 distal = Bochl
 (= not Angel)
 #10 took 5 min
 v. no. distal.
 Distal 2
 12:30 Thurs Day
 Re: Bvo/Seb.
 Monday
 7 AM Sci.
 The present on
 PPM — But it's
 pretty 16 bit.

0
 1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35

w. 2 params: one would desc. behavior near $x=0$; t. other could desc. behavior at $x=\infty$.

We could try 3 decays (6 \rightarrow 4 params) to see if we do any better.

Another posy is 1 decay. (w.o. $\frac{df(x)}{dx} = 0$) so only 1 adjustable param.

I may be able to solve for t. 2 params analytically & then use it as a first approx. for a 4 param case.

Hrr. I really need a good review of recent Menkes' work on PPM for data, act. & how to use it in Lsuch.

List of ideas to think about (2 STACK):

1) Difference betw. OOPS lang & Lispish lang. in Lsuch: Is it a fundamental difference?

Idea of doing part of PPM & having part of output. In imp't Q is whether savings in pms occurring in Lsuch are more likely to be useful in industry than things occurring in OOPS P&Ms. Perhaps try them out - to get ideas.

2) TSO's: 129 1/2 is very interesting! It is a list of mathematical facts that I expect TM would be able to learn. A fact, for TM is defined by how it is used to help solve problems.

Hrr. I don't think I ever discovered a reasonable way for TM to learn "facts" of 129 1/2! For each fact, I want to say just how it affects TM's (behavior / problem solving routines).
 \rightarrow .15 ff may be out of t. BIG Bottlenecks in TM TSO writing!

SAAB was largely concerned w. ~~the~~ diffys of this sort. I put out a QA form of problems ~~which~~ would cover an enormous range of problem types - so writing TSO's would be much easier...
- but still Not Easy!

I still want to write up a detailed outline on how to use "PPM" for Lsuch - but I think

.15 is t. BIG PROBLEM. Having t. SW. Pgm for PPM would seem to make it easier to write TSO's - but only of t. kind where I didn't know in detail, how to write.

Waste to be done. - Bob Note (31) $20 \frac{2.5}{2} = 1.25 \left(\frac{3}{4}\right)^{20} \approx e^{-5} (1.25)^{20} = 867$

My impression is that it's best to write TSO's slowly by hand (English), then work on down to

P&Ms. One top down way is to start w. Big Problem soln & factor it down to t. primitives:

If some parts are too hard to factor! (leave them as primitives).

31 .27 \rightarrow These TSO's can be evaluated for incremental Δ of each problem soln. - for trouble shooting

We can put PC of 2 particular symbols. \rightarrow This can be done using BZIP 25 or BIT Cost

w.o. my P&M supply - still bit cost is \approx so PC factor of $(\sqrt{2})^t$.

(Probably M or B) good point for trouble shooting! Several diff. solns. for same problem

[SN] in TSO's, we may want to put multiple solns. into t. corpus, so that TM could learn from them.

[SN] If we \downarrow bc per symbol from 2.5 to 2 : a factor of 1.25 from

20 symbols \rightarrow PC $\rightarrow (86.7)^t$; symbols $\rightarrow (9.3)^t$.

00:13430: ^{Writing} Also, writing these TSO's will give me a better idea of what I want/need in a guiding P.D. (analogous to PPM). One big diffe betw. present state & 1990 (w/ same) is that by using v. R.H. TSO, I can write TSO's a lot easier. A TSO having much more of Definitions, may make it easy to write them. In fact learning of definitions can replace "Telling" students things

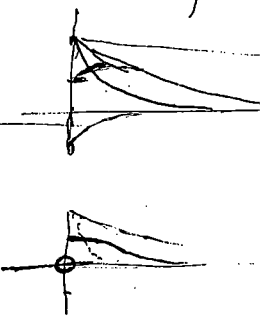
The "Lap rule" Born sup. seems to be close to a kind of suggestion for new Caus that I expect in lang. There certainly are other impl. ways, but basically, I think that v. Born's first main fundamental, a law of nature, ^(lang. occurs) ~~is~~ ^{by which} (P.D.'s are usually obtained.) PPM is a rather fast way to implement something very much like Lap's "rule".

- List of #0055
- 40
 - 20 in HPS00
 - 17.5
 - 7.1
 - 2.3 - probability
 - 2.3 - spec
 - 2.35
 - 2.1

10:13419: The idea of new P.D.'s or tokens as "data objects" or "material objects" to be combined with other P.D.'s & w. various other data objects, to create new P.D.'s. So we have an algebra of P.D.'s to create new P.D.'s from old (this tests what much of classical probability theory is a book!) i.e. combining old P.D.'s to get new ones. No good ways to get P.D.'s from R.U./Data.

In OOPS, v. idea of associating a P.D. w. a desc. of use of Tokens in a particular problem sol'n. seems very interesting: — But next, v. small size, in what loss of a statistical (d.f.), Real 2 logical construct. i.e. ~~to~~ ^{to} Set of Tokens used in a PPM is an interesting, perhaps useful, "data" type. Learning to estimate, manipulate these P.D.'s could be some kind of useful Math. problem. Assoc. w. it is learning how to generate a combiner & manipulate S-functs. — which it will need to know, any way.

2007/By
20
1000/By



2.0
2.1
(5N) A 2-param family of decay functs w. nice properties.

$f(x) = a_1 e^{-bx} + c e^{-dx}$. If we want $\frac{df(x)}{dx} = 0$ at $x=0$ we still have

3 params to adjust. for $f'(0) = 0 \Rightarrow -ab - cd = 0 \Rightarrow ab = -cd$. Another constraint: we may want $f(x) = 1$ at $x=0$ — so only 2 params. We adjust f : 3 params so that as a prediction $f(x)$ has impact per PC for an int

compute. Th. score for each $f(x)$: is $\ln(f(x))$ if $y(x)=1$; is $\ln(1-f(x))$ if $y(x)=0$. So we want to maximize $\sum \ln(\cdot)$ or Maxz product $\prod f(x)^{y(x)} (1-f(x))^{1-y(x)}$. This may be easy to compute: (Go we may want to use pseudo loglik pt.)

$f(x) \approx \frac{a}{b} + \frac{c}{d} = \alpha$; $a \cdot ab = -cd \Rightarrow a = -\frac{cd}{b}$

$a \cdot d + bc = \alpha \cdot b \cdot d$	$a \cdot b = -cd$
---	-------------------

$-\frac{cd}{b} + \frac{c}{d} = \alpha \Rightarrow \frac{c}{c} \left(\frac{d}{b^2} + \frac{1}{d} \right) = \frac{\alpha}{cd} \Rightarrow \left(-\frac{1}{b^2} + \frac{1}{d^2} \right) = \frac{\alpha}{cd} = \frac{\alpha}{ab}$

No! " $\frac{1}{b^2}, \frac{1}{d^2} \approx \frac{1}{c}$ are linearly related so ~~any~~ ^{any} param b, c, d will determine OR, or not easily.

For computational we'd like to be able to pick the values for 2 of 4 variables & then easily find the other 2. say we know $a \cdot b = -cd \Rightarrow cd = -\beta$ $a \cdot d + bc = b \cdot \alpha$
 $= =$ somehow $a \cdot b = \beta$.

$(a - \alpha b) \cdot d + bc = 0$ also $c \cdot d = -\beta$ so we know $\frac{c}{d} = \frac{\beta}{d^2}$ and $c \cdot d = -\beta$ $c^2 = \frac{\beta \cdot c}{d}$
 \hookrightarrow gives $\frac{c}{d} = v$ $c \cdot d = v \cdot (-\beta)$ $d^2 = \frac{\beta}{v}$

4 PM

Talked Royal Coast.
4.5 June Friday

One way to compare PPM. Compare pc of a soln: i.e. cost finding it via PPM v.s. via GP.

To compare Soln. times we have to consider long & during such... which is what GP always does. This sounds Diff! ... at a more "Q1" level, Comparison of "Goodness" of GP & PPM w. a fixed population. For GP, no long during a generation, i.e. no long during a Generation for PPM as well.

One (perhaps) advantage of LSrch (using PPM) is that short codes are automatically "selected for".

In normal GP, one has to modify the "Fitness Funct" to include criteria that short codes are better. Lsrch ^{using} same PD, does this automatically. — GP deselects (using fitness funct) i.e. a code after trial because of excess length. Lsrch deselects before trial, which seems much more efficient.

One Big Q is PPM v.s. "OOPS".

We have to consider various (i.e. A.H.) mechanisms that have been & can continue to be introduced into OOPS but could significantly ↑ its power. The extent to which it can do this in PPM is unclear. (Also note 136.10)

T. Useful Mechanisms for induction in OOPS:

- ① ? usual Lap's distribn. for t. next token (Rem Lag).
- ② Ability to determine how tokens
- ③ Ability to (combine, use, modify) presentations of t. son tokens used in previously successful codes — to be used for present problem. → This sounds like a v.g. idea → 136.10

It may be that PPM does do ③ (i.e.) to some extent — perhaps better than GP!

It does this by suitable choice of "escape" probs (or event). My "Korax" does this...

How effectively, is not clear.

- ① Is pretty clear as well as less much better by PPM
- ② In Costa Rica, I got idea that PPM would simulate affected of data.

Now, how I'm not so sure! — Because definitions can create dealing names & various production "language".

Most recently I had idea that it would be useful to try PPM (as my version of) on English text (to see how it compared w. regular PPM & its refinements).

But more important: I wanted to try it on GP problems to see if it gave significant speedup.

→ So: What is expected value of PPM v PPM for Lsrch? What do we expect to get from it?

1) It will be based on of the phases (i.e.) to demonstrate useful TSO writing, to enable TSO writing & testing. → 135.00

Well: Two results: Don't seem to fit together. (This is Re: Murray's Computer)

I put Murray's original 4.36 HDD into my case; replacing my 4.2 HDD w. vacuum pressure the 4.0 GB HDD.

On start up: I got to Logo Logo screen & before 2 cycles routine (prob on HDD) — Never got further. This is same behaviour that we get using Murray's 4.36 HDD and Murray's ESD. So, it works same but it's i. SW on Murray's 4.36 HDD.

2) I formatted the 7.5 GB HDD & put in Murray's Hardware ESD hardware. I installed w 93 F.E. w. no trouble, but on reboot: 2 total mess (maybe I didn't format HDD properly?) suggesting that ESD HW is at fault.

Re: Murray's Computer
ESD

for Murray's
Computer
5/7/04

Sum
Priority

4TH

Note that in all of the foregoing discussion, the fact that PPM could not detect regys like 131.09-13 is quite irrelevant.

A discussion of a soln. to 132.15 ff is in the NIPS report: The section on the R (recognition functions) - in the first or second section of the report (pages 8 of 4/11/04 document)

So the point is, I really should be worried (now) about 131.09-13 - that is, approaches of 132.27 ff is 133.02 and move to the point.

A more General Q: Certain mutation/crossover algorithms have greater expressive power than ~~them~~ in some respects. How to characterize this: Possibilities to the power of PPM.

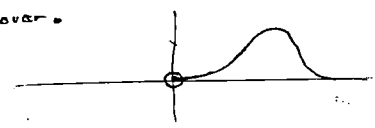
Consider 2 parents ~~with~~ fitness functions. Consider all α crossover. This gives a uniform d.f. over a finite set of kids. Considering all α pairs of parents, we get a D.F. out: "next Generation"

How does this D.F. differ from those obtained by PPM?

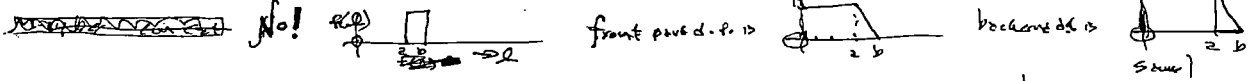
The GP distrib: ~~say~~ w. a finite population, we only get a d.f. over a very finite offspring population - kids are restricted to being found of 2 parents in the population, crossed over.

In a PPM d.f. all possible strings have prob > 0 .

The GP d.f. is formed by the concat of 2 d.f.s. ~~Choose~~ choose a parent at random, ~~choose~~ choose a crossover pt. ~~is~~ the top end of the crossover is one D.F. the bottom end of the crossover is the other D.F.



Perhaps (it looks like) if $f(x)$ is the distrib. of lengths of regys, then after crossover the D.F. of lengths is something like $f(x) * f(x)$ (ie. autoconvolution)



I think in general the form is back end d.f.'s are v. similar.
$$1 = \int_0^1 dx - \int_0^1 f(x) dx \quad \left| \quad \alpha = \int_0^\infty f(x) dx = 1 \right.$$

So = ~~some~~
$$\int_0^\infty f(x) dx = 1$$

The d.f. for concat. is conv. of $h(x)$ w. $h(-x)$.

~~conv~~
$$f_{conv}(x) = \int_0^x h(x) \cdot h(x-x) dx$$
 which isn't exactly a convolution.

Well actually, since $h(x) = 0$ for $x < 0$, $f_{conv}(x) = \int_0^x h(x) \cdot h(x-x) dx$ - which is

a kind of convolution. It's either a convolution or a conv. of $h(x)$ w. $h(-x)$.

Different authors may use different definitions.

So: My impression is that whenever concat 2 d.f.s like this, the d.f. of the resultant ~~is~~ get wider, by say, $\sqrt{2}$ factor. (Unless of course, we have a "selection" for shorter regys).

There is also Q of: even if the GP & PPM give the same D.F., is PPM more

"efficient" than GP? - My impression is that PPM is rather "efficient"

Also its use of Lsearch is efficient.

4TM

Intelligence PQ2

Intelligence Strick 2.0

1620P 53:04

Tom Word 781 646 3703
Joff " 781 355 9934
Mori
message
message

It may be that w. 131.35, I'll not need a "Phase 2"! Tho still, Phase 1 doesn't understand fully what "often" means. — So something has to be added to 131.35!

Hvr: — T. first "layer" is already Universal. — What does it. second layer add to it?

Well, it would seem that finding ways that PPM couldn't find, would be a v.g. thing!

NA. "Universality" means discovery in $T \leq \infty$; But it can be very long! An already Universal induction system can be improved by speeding it up.

Another tack: Derb. just what I want TM to learn in 131.09-13 : write psm that could do it. Then make TCS for TM to learn that psm. JS

SN in f. assoc. ex libid 8/3/00 — A study for a TSM Algebra: "What I want TM to know"

"equals plus equals = equivalent" I want it to know that plus, but e. implications of it. — i.e. I want TM to "understand it" to some extent.

In f. case 131.09-13, T. universal system must eventually be able to do that! How can we speed up the discovery? We should be able to write a TSM to enable/speed up that discovery

I think the original problem was this: that TM learns to do problem types $\sum_{i=1}^k$. T. has a table by trainer using f. \sum symbols. This sort of thing would seem to be not too difficult. Next we want TM to learn to recognize f. a prob. types w.o. f. Trainers including f. \sum 's. Let's look at the details of both types of problem solving.

SN Perhaps f. soln. is that PPM doesn't have to be universal: best there will always be imp. kinds of ways that it cannot (see/dec). Hvr, the main way discovery process is (search psm in a Univ. lang. — i.e. this is universal, T. finds details of 131.09-13 should be (must be) solvable as a main problem type re 07,08, 15, 19... So I think 131.09-13 is not a log r. Criticism of PPM's use in guiding search on a universal set of instructions.

15th is of interest/import. hvr. I'd like to have a cheap "soln" to it before I go on. (Tho this is not obvious.)

Going Back to 15: A couple of ways to solve that part!

1) We give TM (Prob₁, a₁) (a₁ is main) either it solves it w.o. or it (Prob₂, a₂): It has to find a function that solves both probs. — We then give (Prob₂, a₂) it has to find common soln. for all $\sum \dots$ etc.

2) We give TM (prob, a): It finds soln. We then give TM (prob₂, a₂), but we are satisfied w. a "less than perfect soln." e.g. It may find F_1 that solves Prob₁ & F_2 that solves prob₂ & ($F_1 \neq F_2$).

However, when its total funct "now" for any problem (old or new) is now f. function prob₂ \sum (F_1). — This certainly better than no soln at all, but if Prob₂ has a function that solves both, exactly, that would be better. At first, it would seem that TM has no way of knowing that a function exists that solves both Prob_{1,2}. (Hvr, we may start TM out w. its assuming that all probs we give it have a common function soln.)

It should be easy for TM to learn to correlate \sum a₂ w. \sum f₂ — so it could get 100% correct answers.

Bill 7532
Bill: 7532 NO
617: 492

4.24.04
4 TM

00: (30.40 : 130.37-38 is imp. Learning to map probs to indices easy (fast) to base solns. Its mechanism that leads to inference needed to check solns.

Th. problem to be solved is finding a best fitting struct - parallel to level (random method of finding params work?)

05 → How a major weakness of PPM is to be in ability to try params (or substructures) in different orders (?)
Check to what extent this is true. (The universal lang can do some stuff & this may deal w. lang problem (no perhaps not in best way).

09 In general, if α & β should be followed by β and α can have many (various) values; PPM has
10 learn each case separately!
In general, I certainly don't expect PPM to be able to recognize all kinds of reps but

12 I'm disturbed that it can't deal w. bits particularly short reps type. - Just how far can it go
13 t. reps of interest? That's reps is internal to a seq. yet its needed response is

after t sequence. PPM only recognizes reps that end at end of a loop → seq 133 as for conclusion of discussion
Could we get a recognition type that we want by > 1 "layer" of prodn?
Laws of physics all disallow "action at a distance" - so all causality properties via a chain of adjacent disturbances - yet physics is able to deal w. very complicated causal seq relations.

20 To what extent is GP (crossings of function trees) able to deal w. very complicated causal seq relations.
eg. we want to be able to categorize a set of seqs by know reps before. (Backwards: Contact)
or there is a certain characteristic of first 10 tokens obtained by a certain function (strings → trees)
PPM couldnt discover much of this. If vr, if we tried to find functions out first 10 tokens

23 PPM could help find a prop made ones.
Hvr, any reps that can be done by PPM can be done "As a type level" by a univ lang. - just how
to implement this in a good, general way, is unclear. Perhaps UMC looks at first level codes.
tries to find regularities! Or, UMC looks at t & tries to find reps - in addition to "kernel" & is
normal regression repts. Perhaps how UMC look at t error PPM uses! After running a loop long time
looking for reps on PPM errors, can a very long corpus, t. UMC interpret a "universal correction"
for PPM that significantly fits accuracy.

30 On the other hand, I expect that grammatical-like reps (nouns etc) would not be found by PPM
& reps considerably. would be able to deal w. city of reps. (say 12-13).
A finite state grammar with help (NFA). A possy: t. set of contents that predicts
either a particular token or a set of "backward contexts" (usually starting w. same token)
This set is a "state" candidate.

35 Essentially, the problem of 09, 12-13 is a categorization problem.
→ A perhaps V.G. way to deal w. t & higher problem of capacity of PPM. The problem of
improving pd obtained by "PPM" becomes one of "low level problems" of t system. By
improving pd I mean to include cc as well as pc aspects of t. "pd": i.e. it is a "found pd"
in which both pc & cc are considered: if pc is found of cc as well as cc & corpus.
Do I want to (hard to) do this on phases?

I write table to part 4. I ~~start~~ out of existing BZZ egms - & how try to produce it -

Compare w. BZZ, PPM, act.

As far as TSQ's: I could write a tsq w. possl. solns. to problems, & see how many bits were needed for each additional soln. T. values would be within a factor of $\sqrt{2}$ (if it's too small to X.I. & too large) - Good cert! @

I would like to see if it could get solns better than mine! Particularly if I use a font t. seq., w. lots of problem solns - so that it might find parts of solns that I wouldn't think of using!

[SN] on TSQ Design: I will have to solve a certain class of probs that have a commonish soln.

Then learn to solve Then

and so on for several classes.

Then it has to learn how what class a problem is in & apply a (pre-learned) program soln.

If would seem that this kind of thing should be readily (run by PPM. "The class of problem decided" (whoa (run) would be a / prefix to ~~some~~. 2nd attempt to solve a problem. T. fudge discuss suggests NO!.. This is normally done in Phase 2

On second thought, I don't see how this could work! (Deciding what class a problem is, is a "meta problem":

The "Meta soln." (Meta function): Looks at f. problem & outputs a Rd on possl. solns: This is probably could be in

f. form of a narrow distribution on "contexts" (these contexts being ones that are followed by possl. prob. solns. of context)

One way to a program this "meta problem" soln: The corpus consists of (problem decdn soln) pairs.

Unclear as to how PPM would be effective in this case. I think T. Charac would depend on to context

& f. end of the problem decn. It ~~is not possible~~ This would usually contain info needed by f. "meta decdn" of what problem class we have. It's f. kind of problem that Phase 2 is designed for.

One way would be to try to find a function that maps from problem decdn to a pth on \hat{L} (distribution of pths on \hat{L}). A corpus for ~~such~~ (ing such a func, could be the problems \geq ~~some~~ pth w. pos on \hat{L} . This is a semi discrete problem, since small errors in the output are not imp. (ie. if ϵ of \hat{L} with ϵ kernel, then errors $< \epsilon$ are acceptable)

At f. beginning of TM's TSQ, a per null context will eventually give a ~~pd~~ over useful context that are solns to problems. The Meta function will map to a "sub distrib" of Pths & pd. "Narrow it down", so to speak

to "solns" in .15 correspond to pths on \hat{L} that are ~~followed~~ followed by null context.

Remember many possl. functions for problem space: Context of them are useful to solve f. problems:

Context are useful to solve f. meta problem.

In a preliminary training phase, a problem-by-problem could be placed (by forner) as a f. problem.

If we then use Problem Solvers as corpus, then inserting a new problem m, would end in the proper "index" of problem type & condition a good soln. key, when f. problem (as context) was inserted into \hat{L} , to find how a soln may be generated. After PPM has been done for many problems & problem types: we want to be able to insert a problem w. o. to the data.

If would not work if the meta came before problem decn... This is a major criticism of PPM method

I think we turn into f. diff. of .15-18! T. point is $\alpha \wedge$ is not as close to $\alpha \wedge$ in PPM.

Well, say we solve a large set of f. first order problems (w. or w.o. index... - no indexing would speed up solns). After have solved these then, we use that same \hat{L} , to develop, T. meta egm now maps from problems into "soln. space", or simply from problem decn to indexes. While f. input subsets for f. Meta Solns would be like those used to solve f. problems, the output functions, would be much different, unless we have given TM problems that are similar to it.

Meta problem (at least in output form). {Remember Ob/op algbra: we have similarity of probs to meta probs on the obs but not to ops.

8/13/00 - Dup. of 24 earlier page.

This will be an ordered list of problems, tasks, definitions, ...
--- toward the construction of an initial TSQ for TM.

167.24-.34

l, m, n are 32 bit random numbers

u, v, x, y, z are variables

":" separates examples. ", " separates data within an example

cond means "conditions for this problem"

imp means "what is implied by these conditions"

n=n e.g. 3=3, 7=7

cond x=n imp n=x

"[" and "]" are "metasymbols"

(I guess I want [l+m] to be my numerical value of l+m.

l+m=[l+m] e.g. 4+5=9 --- learning Addition

[l+m]=l+m e.g. 9=4+5 --- Equality Commutes

cond x=y imp y=x Equality Commutes

l+m=[m+l] Addition Commutes

x+y=y+x Addition Commutes

x+m=m+x Addition Commutes

l+(m+n)=(l+m)+n Addition is Associative

cond x=m imp x+n=m+n --- if equals are added to equals the sums are equal.

cond x=y, u=v imp x+u=y+v --- as in previous example

l-m=[l-m] learning Subtraction

m-m=0 meaning of Zero

x-x=0

cond x=m imp x-n=m-n -- if equals subtracted from equals, remainders are equal

Note: TM may learn $(\frac{m}{x}) \cdot x = m$; for many values of m, x; but for 1 value of x, $(\frac{m}{x})$ is

not true: so its true w. probab.!

T. freq. are not ^{good} Q.A. problems: e.g. $\underbrace{cond\ x=u}_{Q} \text{ imp } \underbrace{m=x}_{A}$

Since cond x=u implies lots of things, like x=x, m=x, 1=1, etc.

If seems better to give problems in which various alg principles are "discovered", i.e. because that TM has no defined

mechanisms that make certain (operations/prob.solv. method) very likely. TM acts "as if" it knew that

w/lt. was commutative. Also, we want TM to realize that all commutative funcs are in imp. ways

One way to deal w. this: First give TM so work prob as if it knew about commutativity of several

l+m=m+l funcs. Then later get it to learn how v. commonness of l+m commutative rule

is used in ~~the~~ various levels of math problems. (It may be that this aspect of commutativity is an

essentially difficult discovery that it took the Math community much time to discover ... in which

case we may need to give TM ~~the~~ "Hints".

So 128.23 - to express the problem of computing the kernel as one way to do linear prod. — that auto correlation is a soln. of Torgnyz Matry with given M coeffs of kernel in $N \times N$ time (N is original corpus length). Hvr, other means of N linear prod. (128.26 ff.)

It is not unreasonable to do a kernel w. several thousand pts. ("logs"). Hvr, we do want to avoid ~~the compute coeffs~~ (in kernel) to be $\ll N$ because of overlapping data

To start trying this out, all we need to do is put a corpus in a lexical order! Perhaps we have to know it.

Taken following each entry I think it is the sequence of "following entries" that we want to auto corr of.

Note that we want to auto corr. for each to kern, so we have R different corpi.

SM could we use cross-correlation between R time series? Well look at some cross corrs: see at they are "zero". "Zero" in sense of useless for prod. (23)

SM

10 stocks: each has 2 tokens! Make corpus of tokens ~~of~~ of stock that are max % of each day.

Then use a "hedge bet" on PC's for each stock, every day (the bet will be a vector of up to 10 different stocks — like a hedge fund). Actually, one may not be able to "hedge bet" because of unavailability of leverage.

Picking 10 stocks to use, is not easy! One way: pick in same industry! (likely to be correlated).

Another: use currencies: correlated somewhat. Another use indices DJI, S&P 500, etc.

Use various riders, demands (2 indicators)

Another trick! Each stock has 2 tokens: (1) when it's high on that day. — so do time series structure, (2) when it's low on that day.

But each day we have a pair of tokens in say high/low order.

IN SM, we have fixed set of tokens so we don't have to do "Exceptions" every day.

161.00

In time series curve fitting for kernel, we still have the problem of "How many Gaps". T. must account "soln" (which sounds reasonable), is to make gaps based on the past. (T. must account past poss. for as close to present problem as poss.)

Hvr, in the present case, the spread of kernel coeffs is rather limited: they could be \ll or $\gg 1$. We could probably narrow it down further by looking at past kernels of content u with present one.

If we use ~~assume~~ coeffs are \downarrow w. distance, we have a relatively narrow distribution. Also the ratios of previous adjacent coeffs could be ~~in~~ in a rather narrow range. (Actually the ratios of successive coeffs that are actually used can be > 1 because the previous k it could be wrong)

If we assume the coeffs \downarrow w. distance, the ~~code~~ code cost of the successive coeffs \downarrow a lot.

After eyeballing the set of coeffs for successive window large no. of coeffs considered, I can probably suggest some good terms for the kernel. — that rather than successive coeffs vary slowly

in a way that I can optimize. so $e^{-ax} - bx^2 + c$ maybe good (only 3 param!) — only 2 if we want a second kernel.

The question of a prep for non-linear prod. like $X_{n-k} \cdot X_n$ together, product X_n . If they are distinct plus just special int to product X_n .

These may be to only ~~symmetrical~~ symmetrical second (or third) order products.

I'd like to be able to express the kernel as $\sum a_i e^{-i \cdot x}$ because it's easy to compute the values, so try it w. 2, 4, 6, 8 values of $-i$ & optimize w. Marquand's method.

$$\begin{aligned} \text{or } & \frac{d}{ax^2 + bx^2 + c} \\ \text{or } & \frac{d}{\cosh bx} \\ & = \frac{d}{e^{bx} + e^{-bx}} \end{aligned}$$

A good way to estimate kernel: for each token sequence, $(O_{or}())$, do an "auto correl".

It may fit to an exponential so take ~~the~~ $\frac{1}{\lambda}$, since it's ~~the~~ auto correlation of string function.

say $X_i \Rightarrow$ a randomish seq. of 0's & 1's (at least for λ).

$$\text{we want } \lambda, \alpha \Rightarrow \sum_{j=i-\alpha}^i X_j - \frac{\alpha}{\lambda} \sum_{j=i-\alpha}^i X_j e^{-\frac{j-i}{\lambda}} = \text{min.}$$

I think $\alpha \Rightarrow$ something like λ . I think we want

The idea is, $\forall X_i = 1$ for all i then we want λ to be 1. At any rate, for linear predn. the mean of predn. must = mean of predicted.

If X_i are random 1's w density d ($0 \leq d \leq 1$).

Then for large λ and α normalized kernel, we will get

α predn of d for ϵ next X_i , which is to correct average.

The larger λ is, the more longer ϵ $\leq \alpha \leq \lambda$ & ϵ closer we (usually) get to prediction, d . The goodness of prediction relatively insensitive to λ , but \Rightarrow best for $\lambda = \infty$. Here we assume ϵ 1's are of a uniform density.

It may converge "bursts" (which is causal interest) from $\lambda < \infty$ may be better.

one can't ~~not~~ ~~short~~ ~~but~~ but one can't ~~kernel~~ α should be no longer normalized.

say X_i comes in Gaussian shaped bursts α to burst area of density d .

Converge (at random), one every $\frac{1}{d}$ (at least). Simplest $d \ll 1$ - so infrequent bursts. I guess λ would be about burst width, α maybe normalized kernel.

Is it possible that there is a constant term in λ predn.

normalized kernel to be normalized

$$\sum_{i=1}^{\infty} k_i = 1$$

$k_0 \in \phi$ ← first prediction of $t=0$.

$$\sum_{j=-1}^{\infty} \alpha \frac{e^{-\frac{j}{\lambda}}}{\lambda} = 1$$

$$= \sum_{j=1}^{\infty} \frac{e^{-\frac{j}{\lambda}}}{\lambda} = \frac{1}{\lambda}$$

$$= \sum_{j=0}^{\infty} \frac{e^{-\frac{j}{\lambda}}}{\lambda} = \frac{1}{\lambda} - 1$$

$$\frac{1}{1 - e^{-\frac{1}{\lambda}}} + 1 = \frac{2 - e^{-\frac{1}{\lambda}}}{1 - e^{-\frac{1}{\lambda}}}$$

Then, essentially, it's a problem predn but one is allowed future info as well as past.

I could try non-linear (cross products) terms, If I want to do linear predn,

then cross correl ~~into~~ (\equiv autocorrel) has all needed info.

Cubic terms might be sig! $\sum X_{i-1} \cdot X_i \cdot X_{i+1}$: If $X_{i-1} = X_{i+1}$ then

Parasitic/very likely that $X_i = 1$. Also $X_{i-h} \cdot X_i \cdot X_{i+h}$ terms.

This could be a very interesting study in "Prediction" (Vibrose's). Was also have to state results so that λ ~~minimum~~ α when α of α term. is a function of term indices divided by α of ϵ original corpus.

The X_i is only 0 or 1, our ~~pre~~ predn are ~~between~~ 0 & 1 (real) α , are essentially pc's.

Treating the kernel discovery as a linear predn. problem, using auto correl info; It's a Toeplitz (P47: Nuss, Rice, Matrix this is to be normalized \Rightarrow takes N^2 operations rather than N . P 433 has said for common, less general form fast E need - It could take $\ll N^2$... but I have to study this. Apparently takes $M \times N$ operations - M being no. of bits & N no. of data pts in corpus. Actually, it takes about $M \times N$ operations to get + auto corr function! P 430-433 of Nuss, Rice. (in Power spectrum Estimation by Max Ent Cell placement Method) seems directly relevant, since it seems to use X_i auto correl. of t & t equal.

4IM

Somewhat funny about λ kernel d.f. | Say we have a distribution of λ at various latitudes of λ kernel.
 For a finite distance D . We slide the distribution past itself. For each λ kernel displacement,
 we count — No! It is not what we want!

Simulate a continuous d.f. ...

Say we have a seq. of unit λ pulses / counts in bursts usually, but also some random
 background & weaker bursts & broader bursts. To predict $(0, 0, 1)$ or λ point, using
 to knowledge of all other pts: (1) What would be a good kernel? Say we use Max likelihood — we also have to
 predict zeros.

→ Actually, the real problem is to assign ~~kernel~~ to each token & to get max pc ~~for~~ for product
 of pc's assigned to all tokens. One way to do this is to try to do a kernel for 1 token at a time;

to get max product of "Yes" & "No" products by λ kernel. While back I discussed the construction
 of an approx to a kernel (~ 17.32 ft) Minus of $\ln(\frac{count}{EUV})$ but this assumed I had a "non-parametric"
 model of λ kernel to "smooth" (E approximate) — so it doesn't look so useful now...

But the idea of expressing E & λ kernel (as a sum of exponentials (Laplace form))

Still sounds good — but λ error criterion will be difficult — simply max product pc. \equiv Max likelihood.

I can either do each token individually ^{or} do them all at once. Calculate whether it would
 be much faster to do them individually. — One way: Assuming all tokens have same

kernel — one could do one token ^{type} by itself, then use that kernel as an initial
 approx. for doing them all together.

I can just use a 1 sided kernel ~~to~~ to get λ kernel params. then use both
 sides for actual param. Perhaps try each side individually, to make sure they are about a same.

Eyeball λ data to estimate good starting params. The open program is pretty much indie
 of λ core. I can vary both λ coeffs & λ exponents to use in λ Laplace form.

— so start maybe w/ 2 exponentials (4 params) then 3 expts then 4, etc.

A major problem is to make sure λ kernel is ≥ 0 for all values.

It's prob. that Maple has a PDE standard ~~of~~ Non-linear Opten param.

— But putting λ data to it would be non-trivial.

Now. Recipes describes "Marquardt method" pp 523-528; ← Try Google, Also see Book on
of Math Modeling
670k.

IN PS folder: dumpster - "Marquardt" 3/25/04 is m djvu (1987)

Maybe get Alex to print it. — Github does DJV files.

Around 3/25/04 "EM algorithm" There are 6 files 3/25/04 on an email which
 is usually used for opten of Max likelihood... but it's misssydate

EM. PS gives good picture on pp 1, 2, 3. The method itself is "obvious", but proof of always
 converging, is not. From + lines 38: It's not clear how to apply EM to λ ~~prob~~ problems of
 getting a good kernel model.

T. Various EM papers may tell how to apply the idea to very many non-parametric ~~problems~~
 possibly suggesting how to do " λ kernel" (which may be a common problem).

4TM

So good pc's / solns of probs that are continuous of v. TSP!

SA

In Induction/TSP acquisition! We have ① a universal lang. to express program solns.

② an initial pd./on/ statements in that lang. This initial pd./is not a universal d.f., but it enables 2 certain ans. of lang.

For P₀: 3 degrees of sophistication! ① symbols of L₀ each have ^{known} fixed pc's so pc of a cond. is something like expr lang. (not cond. term). ② like 1 but pc's change, accumulating pc via L₀ replace... assume known seq. for taking of lang. ③ PPM - like d.f. on Tokens.

4/23/04

Q1: For a universal lang., is it easy to have "derivations" or exp.?

To some extent P₀ imposes itself on responses to TSP.

Q2: At what pt. can we ~~modify~~ give a system to prob. from of improving P₀ to L₀?

→ If deriv. are poss., it is poss. for L₀ to be effectively modified to any lang. (Coul. count)

Q2: Demands in L₀, P₀ is TSP.

In "Phase I" of ALPHA, we give ^{QA} ~~prob.~~ problems! T. system wants to find a single funct. to solve to whole set: At first Deterministic, then probabilistic. In either case, past solns (w. subsets of "T. present" TSP) are used to give a PD. for search on t. respect "Corpus" (old ~~exp.~~ soln to new augmented corpus). The Corpus is a set of solns to "t. same problem" - or to "aspects of t. same problem". It seems to be a ~~GA~~ GA problem, i.e. a PPM guided search.

SA) R₀: GP, GA! There is something unrepresentable about the "Monte Carlo" search.

Taking ~~exp.~~ ~~prob.~~ time, N (where N is no. of diff. cond.) N can be ∞, yet if system seems to work!

- How? Well, one way: we pick cond. z' w. P_z (where P_z is pc of cond. being "correct/soln." (var. ~~prob.~~ testable) We decide on time limit, say 10⁵ = d_z.

We do M trials, each lasting T secs, max. We then report after T ≤ 2T, etc.

If t. soln. has p_z > 10⁻⁶ and it takes time = cc_z to generate & test it, it will take total of w cc_z = 10⁶ time to find it. - which is much worse than L search, if p_z >> 10⁻⁶.

L search prob (z) $\frac{cc_z}{P_z}$ time for search.

In QA Phase I: just how did I go about representing probabilistic models?

I Phase I wrote a lot about P_z: perhaps including summary.

For Phase 2 I mainly needed P.D.'s on PSM_z problem pairs - i.e. P_z's was a 3 param distribution mean, var. & ~~prob.~~ moment; so 0, 1 & 2 moments.

Before going on to v. Details of Phase 1 & 2, write a good usable summary of just how T2T 2 // is and so work!

I don't quite remember just how the kernel was computed! Somewhat: 103.25 ft: ~~was~~ 99.00 off early value... 2 ft. unclear ~~to~~ ~~not~~ ~~bad!~~

107.20 off isn't 100% clear, but I think we can find a kernel for each token, indep. of P_z's: They are probably all some shape, so we can average them. (Be aware on whether t. shapes are indeed etc. same).

FTM

00: 12440: perhaps I can update in T2T when T interval gets sufficiently $> \Delta$ (Time needed to update)

So: I guess that I have to. ppm in end data) to write out a more detailed & cr of all of Lsrn for T2T Lsrn: is for Π is random to some extent.

It might be worth to do T2T Lsrn. I've forgotten the details of the proof of optimality! This does have a slight advantage over T < 2T Lsrn.

I really can't update more frequently than after each prob. soln. for T2T Inv prob. Since process is noisy to update with!

Hvr., for OZ problems, perhaps I can update before final soln. Since I can have already obtained some conc. at suppt. w. t. soln. This is essentially very \sim to GA. Since I ultimately (for phase 2) want to solve OZ prob, this is GA is an impt direction to go in.

The eq of 124.33 $m = \sqrt{\frac{2.5}{\delta}}$ assumes updating & soln. time. In OZ prob,

this idea has to be modified: one way: compare + pc's of by G. could be before & after an update. I. e. after an update, find best pc's of various by G. could that were obtained before update. Since time to soln is $\propto \frac{1}{pc}$ this is directly translatable

into "S". Actually, updating is designed to \uparrow pc's of programs by G. could. If should do this w/o overriding, \rightarrow (17)

Work out the operation of Lsrn for OZ problems. Certainly (10) seems to be an impt part of it. Perhaps (14) can be used to estimate S more accurately!

17: (14) \rightarrow T. way it works: After every update, we recompute S (because it takes little time). After an update of m. ~~new~~ words, we compute the new pc's of each of the m words w.r. f. corpus & compare these pc's to the values obtained before this update. The total difference in $\frac{1}{pc}$ will be $m(m-1)\delta$ (note factor of 2 in m^2 w.r. f. the $\delta \frac{m(m-1)}{2}$ of 124.31.) This "S" computation gives us real feedback

then we may want to delete a word from corpus both forward & backward PC were both corpus

how well the system is "working"! \leftarrow Very impt! S is a very incisive remark about efficacy of TSD! This "S" was very small in the ~~English~~ English Corpus used in original BZW 1994 paper.

The S & w. corpus length is scaled to \rightarrow for corpus of \sim 150M chars. Make S estimates for 100M char corpus.

Can we use a (very fast) BZZ program to estimate S in TM! \odot - we can do this because to estimate S, we don't need to pc's of individual tokens: only the products of those pc's - which BZZ gives quickly.

From 1994 BZW paper: Table 2, p14:

Corpus 14 \rightarrow 4K symbols. $\frac{.52}{34} = .17$ m bits/symbol (?). $\frac{435}{34} = 3.93$ $\Delta = .52$ But $\frac{.52}{(34)^2} = .06$ M per symbol.

Interpretation of BZW data is not so simple!

BZW Data		
14	4.35	.52
4K	3.93	.44
6K	3.39	.41
435	2.98	.32
250K	2.65	.32
1M	2.43	.32

Anyway, what I can do is write a set of programs & solns, & see how to ppm given pc's to new probs in an acceptable way. - Doesn't pick up the sub-functions that I expect to see. This, of course, takes much less time than letting TM randomly search for solns. - Also, it's much easier to ppm! I could get a useful approx. from BZZ!

This base is of interest! Write some ~~basic~~ series of programs solns & see if BZZ can

4-TM

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

∴

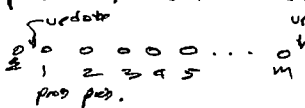
$$3 + 2 + 1 = 6$$

$$\frac{3+1}{2}$$

$$\frac{m^2 + m}{2}$$

00: 123.90

Away to think about it: Consider a seq of probs. that have been done w. updates after each problem. This is "system 0" ($\in S_0$). How much worse/better is it to update every m problems? (Time needed to update is same constant in all cases).



call S_m system that does this.

No: $\frac{(m-1)(m-1)}{2} \cdot \delta$

The update on one problem causes subseq problems to be worked faster for future S .

At first problem after update, both systems (S_0, S_m) are same.

After second "	"	S_m is behind by δ
third	"	$\delta + 2\delta$
fourth	"	$\delta + 2\delta + 3\delta$
m th	"	$\delta(1+2+3+\dots+m-1) = \delta \frac{(m-1)(m-2)}{2}$

show by $\frac{m(m-1)}{2}$

impt!
 If you don't understand why of this "Sp 1w":
 Oh! The idea that the problem is deficient by not having $(k-1)$ updates of previous problems

Say Δ is time needed to update. when m is such that $\Delta = \delta \frac{(m-1)(m-2)}{2}$

System is as good as $S_m \rightarrow$ as fast as S_0 .

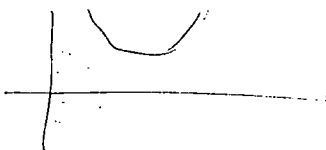
At end of m probs no update so total cost per problems is

$\frac{(m-1)(m-2)}{2} \delta + \Delta$

To find $m \geq 0$ this is min.

$$\frac{m-1}{m} (m-2) \frac{\delta}{2} + \frac{\Delta}{m} = \min$$

$$\frac{m^2 - 3m + 2}{m} \frac{\delta}{2} + \frac{\Delta}{m}$$



$$x = m-1 \quad \frac{x}{x+1} \quad \frac{\delta}{2} + \frac{\Delta}{x+1} = \min$$

$$\frac{(m-3 + \frac{2}{m}) \delta}{2} + \frac{\Delta}{m}$$

$$1 - \frac{1}{m} + m \cdot \frac{\delta}{2} - \delta + \frac{\Delta}{m} = \min$$

$$\frac{d}{dm} \left(\frac{m-1}{m} (m-2) \frac{\delta}{2} + \frac{\Delta}{m} \right) = \left[\frac{(m-2) - 1 + \frac{2}{m}}{m + \frac{2}{m}} \right] \frac{\delta}{2} + \frac{\Delta}{m} = \min$$

$$\frac{d}{dm} \left(\frac{m\delta}{2} + \frac{\delta}{m} + \frac{\Delta}{m} \right) = \frac{\delta}{2} - \frac{b+\Delta}{m^2} = 0 \quad m^2 = 2 \left(\frac{\delta+\Delta}{\delta} \right)$$

~~$m = \sqrt{2 \left(\frac{\delta+\Delta}{\delta} \right)}$~~

can know Δ but δ will be quite uncertain. I can determine δ empirically, by working a seq. of problems; then know where is "point 2".

from pt. 2 we work a bunch of problems. — we see how long it takes. (Go back to "pt 2")

update δ previous m_0 problems, then work to solve batch of problems & see how much δ in time went.

26

Or, just after an update, we read several problems that occurred before that update. We can pick an initial m , then do 26 periodically to estimate δ . Since we only read, say, 10% or 20% of k problems for 26, this will not be much of a penalty for determination of δ .

0

31

23

$$\delta \frac{m(m-1)}{2} + \Delta \Big/ m = \min \rightarrow \frac{\delta}{2} \frac{m-1}{m} + \frac{\Delta}{m} : \frac{d}{dm} = \frac{\delta}{2} - \frac{\Delta}{m^2} = 0 \quad m = \sqrt{\frac{2\Delta}{\delta}}$$

so $m = \sqrt{\frac{2\Delta}{\delta}}$

see 125.00 for translation of this eq. for 02 probs.
 see 138.13 for a cheap way to do updates.

Anyway from 123.90, 13, 84; $\Delta \in \mathbb{N} \cdot R \cdot \eta_e$. This could be 200 NL: NL is no. of school m. to solve a set of problems. T. time to solve a problem could be 10^9 tokens.

Here "solver" is not proper term for 02 probs. 100 problem w. 10 tokens each is NL = 100;

So $\Delta = 200k$: which is \ll to 10^9 to solve a problem. So may be update every problem?

At first, it should take $\ll 10^6$ to solve a problem so may not be feasible.

In T2T such I was remaining of updating more frequently than just "solve problems". \rightarrow (Spoo) 125.00

4/21/08
47M

Mandas
Dachner
Harold
Ray

10:
03

Trouble. (22.37): If a cond is \in (any R D_c is \in kernel width W . — It takes $\sim W \cdot R_c$ operations for "correct" \vec{L} profile with addition (or deletion) of kernel.

On 2 other hands, a complete update of L 's pc's requires $\sim R \cdot N_L \cdot n_2$ operations.

Maybe only $N_L n_2$: We move along the L corpus 1 by 1; each time we do a block, we update only 1 of R kernelizations of L corpus.

T. update algo is $Y_{n+1} = (1-\alpha) Y_n + \alpha X_n$. $K_n = \text{const}$. (perhaps a continuous parameter \therefore a shift.)

So $(1-\alpha) Y_n = Y_n - \text{shift } Y_n$... So $Y_{n+1} = Y_n - \text{shift } Y_n + \alpha X_n$ (or a constant)

Shift Y_n ... seems to be a fact in SE — 200 is clock in 406 — maybe less in Pentium.

So we do $Y_{n+1}^j = Y_n^j - \text{shift } Y_n^j$ for all R components. We need n shifts.

add a small constant to one of the R components: (Presumably $n \leftarrow n+1$)

So probably time may be $\sim R \cdot n$ shift. It is units by done in (1. in Part II III or IV.

So it is $N_L \cdot R \cdot n_2$: But this maybe \in no. of clocks or no. of clocks $\times 8$ (or $\times 16$ for pentium) ... Say $R=10$ (min); $n_2 = ??$ 20? @: $200 \cdot N_L$.

Say we are keeping only the top 10% of ranks — still, N_L can be large. It corresponds to GA "population size" is $N_c = 1000$; $R \cdot N_L \cdot n_2 \approx 200k$ clocks for one update of \vec{L} 's pc's.

N.B. ... small constant can be \in to **Some of the cond** at this pt., w.o.

appreciable \uparrow in cc of L update!

SN "Inner loops": He has a 10 clock algo for generating random num. which he says is Very Good, fast.

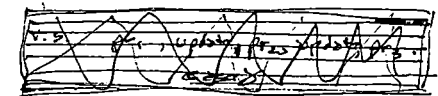
So: See how much time update takes \in compare it to amount of Benefit we get from it.

BAR expressed in common unity of "time lost or saved". This is a classic problem about

Scheduling Problems v.s. Meta problems! — But more realistic = better parameter knowledge etc.

So say update takes time T . If we update first, subsequent problems will be solved \in a fraction of time \in sooner (T \rightarrow T(1- δ)). \uparrow is reduction of \in some of \in problems in \in update.

So P_1 , update, P_2 , P_3 v.s. P_1 , P_2 , update, P_3 .



In δ cases, compare times for P_2 & P_3 times for P_3 : α v.s. β In α P_2 is solved faster, but in β \uparrow update includes P_2 in addition, so P_3 will be solved faster.

In δ , \in 3 problems are each solved as fast as possible, but we have cost of an extra update.

I need \in model telling how much an update & soln. time as a function of what's updating.

Consider only first-order models: So \in P_1 , P_2 , update and P_1 update, P_2 update. Both solve P_3 in \in same time. This disregards that P_1 in δ , \in soln to P_1 is better than that of β , so \in second update of δ uses better soln. than \in update of β . In spite of this \in betterment of \in P_3 soln. by δ , it costs \in whole extra update.

So, assume an update at a particular point summarizes \in update into update.

ATM

10: : in advance by convoluting the kernel w. \vec{L} , R times
 -01 Say N_L is length of \vec{L} 's corpus. Then it takes $T_{time} \cdot N_L \cdot R$ to convolve f_i kernel w.
 -02 \vec{L} , R times.

On the other hand using T2T ≈ 12.28 we take about N^{cands} (no. of cands
 passed) computations of p.c.'s of Nodes. Each node may take time $N_L \cdot C$ to compute.
 I don't know what C 's value is, at all! - It is < 1 , but how much less... I just don't know or here.

say Good ideas! So total time \downarrow cands. $N_L \cdot C = C \cdot N^{cands} \cdot N_L$
 So compare $N_L \cdot R$ to $N_L \cdot C \cdot N^{cands}$ certainly $C \cdot N^{cands} > R$.
 So it would seem best to calculate all of the p.c.'s in advance via $2.00 - .02$.
 A better comparison is $(.00) \cdot C \approx N_L \cdot R \cdot C$ Vs $N^{cands} \cdot N_L$ on R.v.s. N^{cands} .
 R is still $< N^{cands}$. \therefore So $.00 - .02$ looks best.

The analysis of 12.28 of ~~T2T~~ still is relevant & makes T2T look up.
 We start out by computing \vec{L} , then we compute the kernel then we get the Laplace xfm of the
 kernel (via $117.00 - 119.03$). Then we get R p.c. values for each point along \vec{L} by
 convolving the ~~kernel~~ approximate kernel w. \vec{L} .
 We can then do T2T or (I don't know) L such.

T2T seems only a factor of 2 slower than (I or Random's yet requires Much less RAM.
 Random's has speedup ~~in~~ that PC's used not be calculated, & a perhaps fast (quick
 and dirty) pseudo-random No. Gen can be used. But, it still seems to use a lot of RAM.
 If we had RAM available. (It needs $\sim N^{cands} \cdot D$ or RAM), would it be
 faster than T2T? Well even if we had ∞ many variables, the time $N_L \cdot R \cdot C$ factor
 is PERHAPS much smaller than N^{cands} (if no. of Nodes, but we need to compute).

In .01 instead of multiplying by C , perhaps mult by $(N \cdot C)$? More likely: Mult \approx
 $\approx N_L \cdot R$ of .01 by N_L which is no. of ~~bits~~ wts we use in approximating f_i kernel by wtd exponentials.
 Is $N^{cands} \approx (p.c.s)^{-1}$? So upto 10^{10} ? On the other hand N_L will probably be $< 10^6$
 So it looks like $.00 - .03$ should take relatively little time - Negligible.

T factor is 0.4, if N_L is $< N^{cands}$; but for a single T2T Run to soln. ; but
 If we want to modify the P.D. in view of problem solns, It seems like a different story.

32 { On the other hand, a single set of cands w. assoc Goves is a known kernel, affects it.
 various PC's in a simple way; Can I implement it in ~~an~~ an inexpensive way?
 Removing a Cand from the L distn. has a similar effect - calculate it the same way.
 34 Can we use the Laplace xfm. approx to cheaply implement the addition/deletion?
 I think we just consider to ~~start~~ state seq. ΔL given by the new cand.
 After it gives the Lap. xfm. kernel is added to the old \vec{L} ? How much time? ($< C$)

4 TM

Bill's birthday 10 Tribes WARLINE Tuesday 11/1/04 6:30

20

Speech Recognition: To discover how humans recognize phonemes? See what kinds of errors they make (say w. 1 of noise: Gaussian random phoneme noise) Also do error study for synthetically constructed phonemes vs. ones obtained from humans.

03: 120.40!

T2T: I need (100) pc's for each of 6 Tokens during cond construction

How much processor needed (XNUMX)

One (long) way: Get \tilde{z} & smooth w. kernel, R times cost pc vectors.

A very approx. way: at each pt on L, we move along until we find a noise to hours. -- (do both trials, 28 on acc). I have to do this until I get to pt. at

Which we have pt of least likely token. -- which is an awful lot of testing! T. no. of tests depends on corpus size. Also with that kernel. (Normalized kernel & fit constant width)

10

We can estimate fr. amt. of time spent on such activity. Essentially we have a Bern

seq. of choices & fr. probability fr. of each choice by pt. & its pc.

$\sum_{i=1}^R \frac{p_i}{P_i} = C \cdot R$. $\frac{C}{P_i}$ is time spent on i 's most likely choice

\Rightarrow T. constant, C will perhaps be \propto corpus size. (\approx unnormalized width of kernel).

Here, \tilde{z} is corpus that generates \tilde{z} , not a population of words.

A possibly very useful cc. reducer! That organizes all of the pc's on a single run along L.

Perhaps this amounts to saying each node comp. (R potential pc's/computed) takes out average, t. same amt. of time. T. no. of nodes is about equal to total no. of words created (at a particular T c.B.).

So we create a cond. for each node developed. For words w. \tilde{z} tokens dec by beam, T2T's 2 big savings.

"System of states" amounts at most to \tilde{z} nodes work.

This technique of Dory T2T may give us a factor of D instead of w.o.f. enormous amount of memory needed by U- or Random search. (The outsource Prof. it

may be possible to use this same trick to ↓ amt of RAM needed for both L & Random search)

As a result, \tilde{z} may be better by T2T by a factor of $\sqrt{2}$, but otherwise,

use about same amt. of RAM. No, it's hard to store much more than T2T.

22

HVR: (NB) In T2T, one must store all of the states/nodes that have not yet been

completely developed. In depth-first search, one explores a tree in a "stack of left" manner,

so not many nodes are "in process" at any one time. -- hence \rightarrow depth of i branch

being worked on. See the diagram pgm for L search for INV problem in my

1985 "Optimum Sequential Search"

I may want to write up this analysis of "Optimum T2T" that it's almost best

as it, and, ~~is~~ w.o. much RAM needed, seems like a "BIG Deal"!

10

A perhaps impl. Q: Do the search algos apply to List type omcs, in which one normally finishes a devn of a cond. before testing it?

T. for Question that brot on this latest development (is surprised) wds: how much time it took to calculate pc's "on the fly", rather than compute them all

Refs to QNX
Do QNX on Google.
QNX Real time Platform.
Search Phoenix Developer Conferences
Google: QNX
QNX Momentics
NC download G.2.1
Non commercial.

4/19

20: : An Alternative approach to optimization (v.s. Phase) is GA guided by a PPM. We can keep population at ~~constant size~~ (larger constant size by deleting cands, when better cands are found).
 On second Part of seeing that I have to do L search (of one of 3 types) to do GA, any way!
 So GA is not an Alternative to (L search) PPM: it's just another. (training & Analysis, & Comparison)
 & Part of L search Development.

Presumably, after a sequence of problems of sufficient size, they become solved, PPM guiding Pd is changed. ~~What's the~~ Just how PPM should be done, is unclear...
 One could ^{continuously} add cands to the basis of L corpus, but computing PC's becomes harder
 i.e. one can't smooth the L distribution w. t. kernel - which normally enables very fast operation of the system.

Maybe 2 phases of operation: At first, before we have any good cands, we use a constant L distribution. When we accumulate out better cands, we add them to the L corpus & recompute L distribution. In general, we may want to do this - but we only add a batch of cands to the L corpus (i.e. we compute L), whenever the set of new cands really looks rather good, i.e. has lots of members above the 90 percentile of score.

Important, one does do L search "to help GA", one uses L search instead of GA.
 Here, to be competitive, the L distribution has to be continuously or periodically updated (i.e. Learning during L search).

If we use a TSQL, we can update L. def. between problems or ~~per a few~~ every several problems.

3. 119.37: Another poss. advantage of Random L search is that the L d.f. can be continuously updated. It's not necessary to smooth L w. t. kernel. We assume, here, that the kernel's constant - or that it is Updated irregularly.

A factor that the prop. doesn't explicitly consider is, in Optm problems, the cands aren't "solns" or "non-solns" - they get G values ($\in \mathbb{R}$ values). One way of dealing w. this is to include in the corpus of L, only the top 10% ^{or whatever} percentile of the ~~best~~ G values, ~~with~~ ~~threshold~~ ~~to~~ ~~include~~ ~~them~~.
 Threshold-to-inclusion \uparrow as we "climb the hill". This would be nice too w. (23) (random L search)

In T2T L search for OZ problems; we do a "T" run ~~until~~ until we get some output; then we modify L for ~~next~~ $T < 2T$.

To start out, our initial L could be constructed w. corpus of Hand-solved problems!

I had originally had a heavy TM solve a bunch of problems using some fixed a prop (I had a Z(17), but my Z(17) doesn't change the pc's w. new solns, rapidly and do be much better than constant. After that prob has been solved, after TM has solved each prob PPM way, we use PPM for building pd.

Drew up outline for arch of 3. L search systems ^{but} giving back to how each part works. Process how each is used: what kind of TSQL it needs, what kind of GA prop it can solve.

4/18/04

4 TDs

118.90: Another way to get min $(\ln x - \ln y)^2$ is use non-linear optimization using Newton's second order method I invented! It is probably identical to "Marginal" method. I think it converges rapidly.

So anyway, say I have a good way to get PC's of Tokens. I consider it to realize 1) // Lsrch, 2) Random Lsrch 3) TZT Lsrch. ← Call it TZT Lsrch

For TZT Lsrch, [all I need is 2-R params that "partly" describe kernel: i.e. pair 2-R pts of equal percentile intervals.] ← NO! I need to PC's of each of R Tokens (or better still, their logs... that I can get approx using that trick from "floody pt". Actually keeping track of "where ones" in TZT is not so easy! Drawing would be: draw every fixed ΔT intervals, and check to see if $\ln x_{PC} > \ln x_{\Delta T}$

I might want to use pseudo floody pt: I use 2 registers for each no. The characteristics is multiplied to get characteristic of product. ~~the numbers are added~~. The numbers are added. We look at char of product. If it has a leading zero, we shift it's decimal mantissa by 1. (2) → it looks like a very fast way to keep track of products of PC's! but I'm not sure integer math is any faster than floating pt. math! Its floating pt. addition must slow.

When I'm doing $n \leftarrow n+1$, I just keep track of carries. As soon as one occurs, I inc. ~~the~~ mantissa and shift characteristic. So I always know m bits, & I always know approx (which good unit) its log₂.

5th to get log₂ approx: The characteristic is always between .5 & 1. To change bits to \ln log₂, we need to subtract a certain amt, then multiply by a constant. Since accuracy (not by) we can't multiply or divide by small 1 or 2 shifts is additions/subtractions.

Also, we might be able to do approx. business between nos. in bits & floating pt. We do a few "shifts & subtracts". On basis of sizes of results, we can get approx. & push back. T. program has method "in line" to ↓ time cost.

T. stuff from 117.00 to here is all v.g. & eventually, parts useful. ^{How earlier discussions to 98.00 & earlier are also impl.} T. cc of compn. is as imp. as to accuracy of PC being calculated... So I will eventually have to do that bit by bit analysis. But now, I want some f.b. on whether you can help with TSD's

to make Phase I feasible — so I want to open it, & if TZT is trying to fix problems. Also perhaps, try it on GA? Real has been demo used using standard GA methods I want to know how many trials are needed for soln. of GA probs U.S. normal means.

Going over this is first! T. advantage of // & Random Lsrch is $\frac{1}{N}$ -factor of N, ~~time~~ in Time. N is no of tokens in soln. T. Disadvantage is used for much Memory. However, there is also a factor of N reduction in Memory! But T. resort may be not once can do it w. not so much RAM, maybe GG. Advantage of Random over // is that maybe less computation ^{to get} choices of Tokens. → 120.23

Advantage of TZT is (probably) much less RAM & perhaps much less keeping track of things (just samples for PC-first trace search).

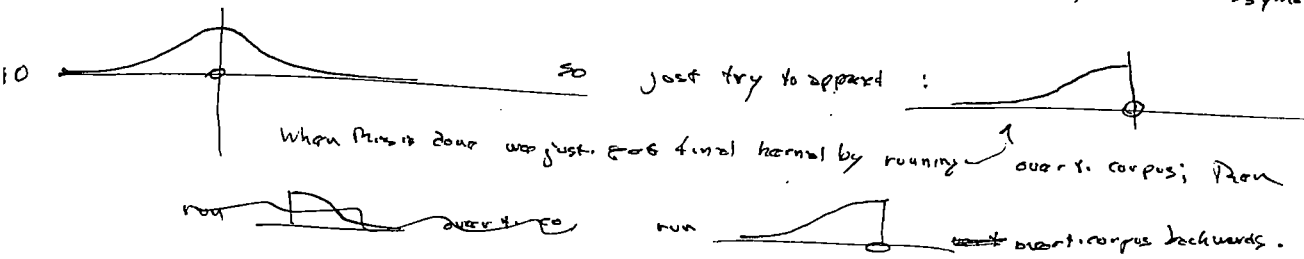
4TM

0:117.40: Probly a good idea to segment to 2-3 tree nodes w. recursive, then given wave

Some of waves, transfer it ~~to~~ completely or ~~most~~ mostly, to Assembly node,

02:117.05! The WTS, (Cj of 4. kernel) can be found by curve fitting, using a ~~least~~ "error² of ln(y)" ~~is~~ ~~the~~ ~~same~~ as GRC - but look into it GRC carefully! The reason we must minimize error in ln(kernel values) is that these are closely related to pc's & its product of pc's that we want to maximize - or to sum of ln's that we want to maximize.

If may be easier to divide the approx. kernel into 2 parts, because it is symmetric!



I think I had a way of ~~doing~~ doing non-ms. error minimize by successive approach

The thing that ~~is~~ varied in successive trials was wt. given to each. ~~the~~ approach pt.

So we did by ms error, but weight shifts so that it becomes equiv. to ms of error in ln.

Say x is a data pt. & y is a pt we want to approximate x . we want $(\ln \frac{y}{x})^2$ to be min.

We assign a wt, w , so $w \cdot (x-y)^2 = (\ln \frac{y}{x})^2$.

First we try to set ~~the~~ $(x-y)^2$ to be as small, using $w_0 = 1$. Next time,

21 $w_{n+1} = w_n \cdot \left(\frac{\ln y_n - \ln x}{x - y_n} \right)^2$ ~~we continue doing this from the new set of y_n 's;~~

we can compute w_n in accord w. 21.

If this doesn't converge: 2 poss. things: 1) do slower convergence, so

$w = \left(\frac{\ln y - \ln x}{y - x} \right)^{2 \cdot \delta}$ in which $\delta < 1$. as δ approaches 0, w approaches 1

2) Some how, do it backward: i.e. in solving $x = f(x)$ (for known $f(x)$)

we can try $x_n = f(x_{n-1})$ successive approx. If it diverges, then its likely

30 that $x_n = f^{-1}(x_{n-1})$ will converge. How we would do this w. 21

31 in 21 since y_n is a function of w_n , we can write $w_{n+1} = f(w_n)$.

Here, the inverse of f is quite complicated.

Another diff: It may be that using 21, w converges at some x pts but diverges

at others. Inverting eqn 21 would then simplify

Since the $(\ln x - \ln y)^2$ error criterion seems to emphasize behavior for small x ,

for initial w , perhaps make $w \propto \frac{1}{x}$ (Note in previous case x is always > 0)

we don't ever want $y < 0$ (so $\ln y$ is complex)

∴ At present time! It would be poss. to convolve L w. kernel in time $\approx T$ or $T \ln T$.

Time $\approx T$ units by doing by simulating kernel by its Laplace xfm, then do $x_i = \alpha_j y_i + (1-\alpha_j) y_{i-1}$

More exactly, ~~kernel~~ \rightarrow ~~smooth~~ \rightarrow X , $x_i = \alpha_j L_i + L_{i-1}(1-\alpha_j)$

Then kernel smoothed L_i is $\sum C_j X_i^j$. Have C_j is the filter Laplace xfm of kernel, K .

We may be able to use same kernel for a slightly changing \tilde{L} , which would be easier if we keep T constant, by deleting 2 corpus members every time we add 6 ones (corpus members).

So if they work out we get a bunch of smoothed L_i : each is a vector (for each) having a no. of components \approx the no. of Tokens.

This would facilitate $T \approx 2T$ $L_{i+1} \approx 1/2 L_{i+1}$ each bit output over 10M machine cycles!

The random L_{i+1} needs a fast random no. generator, set random nos before 0.2, \approx we make boundaries, and a fast series of L_i to random nos. (maybe 2^{-3} tree) Inner loops has a fast random no. generator.

On second bit - 12 doesn't seem relevant! Use .00-11 to get R components of a PC vector (base 6 or R distinct Tokens). If $S = \sum_{i=1}^R$ of R components, we pick random no.

on $(0, S)$ to determine what interval it falls in. One way $S_{2^k} \geq 1$ is random no $0 \leq Z \leq S$, we do $\sum_{i=1}^R$ comparisons: $\geq Z$ we put min $r \geq 1$ in that, it takes on average $\frac{R}{2}$ additions - which may not be bad! Trouble is, we have to do $\frac{R}{2}$ compares to get S !

So we do it $\frac{1}{2}$ times! On the other hand, at each S point we may compute a "S" anyway, to use for normalization. But usually R additions will take little time.

We have to load a register repeatedly in new components.

use X_{2^i}, X_{2^j} : $\left[\begin{array}{l} \text{add } X_2 \leftarrow X_2 [X_2 + X_4] \\ \text{inc } X_4 \\ \text{if } X_4 < R \text{ goto } \rightarrow \end{array} \right.$ This may not take so much time.

I may not need v.g. random nos. - try some very fast (not necessarily good) ones. See if using a better random no. gen. changes rate of (avg. min).

32 \rightarrow Recursing to ends of \tilde{L} : Bound w only used $\frac{1}{2}$ of kernel (?) - but anyway, we can reflect these ends of \tilde{L} to continue them to full kernel xfm!

e.g. $L_{-2}, L_{-1}, L_1, L_2, L_3, \dots$ (and T_1, T_2, T_3 is continuous $T_4, T_3, T_2, T_1, T_2, T_3, \dots$)

So ~~the~~ Kern problems.

ATM

SPOC

20 : 115.90: Another viewpoint for GPD! That we use the "same GPD" for the present problem & for the general (problem of Improving induction Alg. In an early version of I. I. D. S. I. A. reports, (probably ^{still} Abstract) I had this GPD for ^{regular} problems: "Hierarchical problems" like deriving a better Alg. for ^{regular} problems. I spoke of the "GPD" to do this - Rinky of the "GPD" ~~in an~~ abstract way: But if we consider a PD to be a program that effectively obtains P.D.'s;

20 Then it's a more practical, usable way to look at the problem. Reading ^{my} most recent I. I. D. S. I. A. Report, is suggestive: Even the 1975 TB/isi paper.

The idea is that ~~at the~~ lowest level, "main line" problems of TM should be similar to the main problem of improving the Alg. that solves the low level problem. So each soln of a 1st L.R. (low level) problem can go, (to some extent) into the "corpus" ~~of~~ (or the solved corpus) of the h.l. (high level, over the problem) of improving the method to solve low level problems. the "L" system suggest one way to do this.

One thing I hadn't thought about enuf, is improvement by Logical Analysis.

The way I expected to do this is to first teach TM how to solve formal logical problems, then teach it the correspondence between other types of real problems & the formalisms of deductive logic - sort of applying deductive logic to those & u. probs.

A possible approach would be to teach deduction as a kind of induction - other words nothing special about this & it may be possible that it will be wrong.

20 (15.8) (15.8) The idea of a **PD** having a CC Param seems very simple & the idea of "improving a pd. by adding more CC" of 114.10 if it's one way of ~~or~~ (partially?) ordering such d.f.s. The defn. of the Univ. d.f. as the limit of a seq. of P.D.'s obtained from "FOR's" of successively logical CB, - is another way ~~of~~ order of them. So we can have a vast family of different PD types, with different systems, different output forms & different relationships of CC to output PC. We probably need algms to translate from one PD type of S/O, to another.

I could have maintained L d.f. w. G values as a "Utility" & how TM devised different ways to use this utility for problem solving. It's something to System has to consider:

Like a bunch of floating pt. insfs. that we give ~~to~~ TM to use as it likes, or some analog devices, or a Dictionary or Encyc.

Anyway, I've considered 2 extreme paths to TM:
 \rightarrow heuristics, grammar...

I) Phase 1, Phase II, Phase III Grammar of P.D.'s

II) GPD having pd used to solve low level probs; Same PD is used to change ^{improve} Algms used in ^{L.R.} [problem solving of all kinds].

while II) is probably better in a Procedural sense, it is probably best to do I) because I get faster ~~to~~ A.B., e.g. it's better to build successively better objects rather than try to make a final object & write.

ETM

10:11:40: How can we ↓ Reg ~~formal~~ CC? Well one way is to perhaps drop idea of a single corpus, ^{problem}
& use a TSCQ to educate a TM so it can have a better Guiding pd for LSch.

At first, (as in Phase 1) we use a simple induction system to summarize ~~the~~ regps found in ~~the~~
TSCQ in the past. There is probably a limited amount of cleverness in this Guiding Pd., but it is used
by the system with "auto notice" - be configured so that the Guiding Pd. is as effective as possible...
remember the basic prop. being used of Universal, & - sort of a "simulate" model any conceivable instruction
set.

I've noted that the Lisp lang seems more suitable for PPM guiding & PD. Rank forth. - If this is
indeed, so, can forth simulate Lisp in a way that makes PPM a good guide for the resultant lang?
Well forth can simulate Lisp by defining a set of ~~basic~~ Lispish instructions. One Q is -
What kind of TSCQ would keep this about? Essentially, we want forth to define words
→ consist of those words tend to be used in cases & patterns be worthy of definition.
This definition of new tokens really modifies PPM a lot!

13 Pro in AZ 141, which is a Lispish lang, I also define new tokens. (116-20)

~~SN~~ SN I am concerned that in many induction environments, I
need a fast approx. PC values for the "Guiding PD". One issue of induction without
speed has been emphasized: Speech recognition; Analog input - digital output.

- SN In a TM system as contemplated I need:
- 1) Alg. to generate new cand.
 - 2) Fast as poss. way to test cand. (Gorc)
 - 3) way to remember params of each cand, & its GORC.
 - 4) method of mapping 3) into 2)

poss. use of GP to discover good Alg. for Cand. Generation.
In GP/GA we do have to store/remember population. Here in the "L" or PPM method we
don't have to, so it is a distribution ^{plus garbage} has all into a population, but we have used Forum.
So while it does ~~not~~ probably require more memory than GP, it doesn't require much more
~~memory~~ memory (that write by a crasier to do direct mapping w. GA/GP population??) - or more diff.!

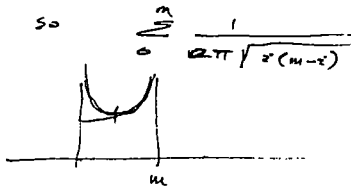
It would be nice if we could do a by comp. capacity run of GP to get a good
cand. generating Alg. - Or a good pd generator to generate cand.
If this were poss., it might be best to use ^{a system} the IPL (Inf. Proc. Ctr) ~~or some~~ & use its output
for the next run of cand. generators, etc. It is prob much more efficient to generate improvements
in the system supervised ("Boots trng") just doing it all in one big run.

The GORC must (perhaps) use a default set of problems from the problem of generating good models.
→ I have to work this out much more carefully!

HA! What .30 amounts to (if one uses ~~the~~ $\frac{1}{10}$ instead of $\frac{1}{10}$ in .32, it's updating
cand. generator a few during the Listen! - I.e. learn while speaking.
~~will~~ will is group.

ATM

20



$$\sum_{z=0}^m \frac{1}{2\pi\sqrt{z(m-z)}} = \sum_{z=0}^m \frac{1}{\sqrt{z(m-z)}} \approx \frac{2}{\pi}$$

$$\int_0^m \frac{dz}{\sqrt{z(m-z)}} = \sum_{z=0}^m \frac{dz}{\sqrt{z(m-z)}} \approx \frac{2}{\pi}$$

So $\frac{2}{\pi} \approx 0.637$

05

So 113.87 becomes $\frac{e^m}{2}$. (No, I have to go back to 113.12 (or maybe 113.11!) to see just what it means!

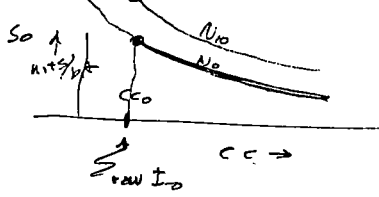
$$\sum_{z=0}^m z^z \approx \frac{e^m m!}{2} \approx \frac{m^m}{2\pi m}$$

I guess the original (error) started w. 109.06

SN

A measure of Utility or Merit, of an induction system! Suppose an induction system I_0 is used in a direct manner, it obtains at an entropy of N_0 bits/symbol at a cc_0 cost.

By using I_0 to direct (i.e. L such) we can decrease N to N_1 bits/symbol, but it will take more cc .
By \uparrow to cc , we obtain further \downarrow of N bits/symbol.



$$N \text{ units/bit} = N(I_0, cc)$$

Now it may be that the form of the curve $N(I_0, cc)$ is universal (i.e. is the same for all induction systems! i.e. we just have to know one particular $N(I_0, cc)$ curve & we can immediately obtain all of the others.

So if I_{10} , say had entropy N value of N_{10} w. a minimum cc of cc_{10} , then we just get

$$N_{10} > cc_{10} \text{ on } N \text{ curve: } N_{10}(I_{10}, cc_{10}) = N_{10} ; \text{ then}$$

$$N_{10}(I_{10}, cc) = N_0(I_0, cc) \cdot \frac{N_{10}(I_{10}, cc_{10})}{N_0(I_0, cc_{10})}$$

The figure of Merit of a Ind system with N units/symbol at cc level.

The way we are able to get better N results from an induction system I_0

we use I_0 to construct the guiding PD, for an L such to find better Ind. systems

27

28

w. better N /symbol. Of course the choice of models, etc. for the deductive L such is critical, & is important in determining how better entropy \downarrow w. cc .

But given 27-28, and an initial Ind system I_0 , we can get a new Ind system I_1 w. a larger cc .

Presumably, we can repeat this operation indefinitely, getting I_0, I_1, I_2, \dots each w. more cc .

Would we get as good, or better results by just using "one stage of simplification" i.e. cc_0

higher cc for L such? (Perhaps not! i.e. successive I_0, I_1, I_2 gives better & better models.

I, "just use bigger cc " / Method uses same cons every time. The successively better I_2 is more like Phase 1, Phase 2, Phase 3, etc.

Oh the other hand, using a previously discovered PD to guide L such is invaluable in improving cc!

So for this "much better PD" one says Haufault.

Speaking loosely, we know that using a non induction system to guide L such, we can

approach to some the discovered cc, but only close --- but w. much expenditure of cc .

(12-90
10:11:24

We use prefix code A^2BCD AAA "2" is sep between nodes to be "2"
 1. first Δ mem "end of a data" is start of α over. The second Δ means "end of prefix"
 We assume $\alpha = a$ (y) α is known but the first repeated symbol in α is start of α .
 What if $\alpha = a$ corpus "madus" is that symbol α now follows a^2bcd of each α .

A^2BCD AAA ... Prefix α regular corpus follows.
 Code α : AAA of α corpus costs for a (a^2bcd a multiple of α if $\alpha = a$). Δ is α means α end of α prefix.
 so pc of a is $\frac{1}{5}$; next symbol route a^2bcd so pc of a^2 is $\frac{1}{6}$; next a complete a^3bcd so
 (or pc is $\frac{2}{7}$)

1. "regular" corpus: N symbols, " a " occurs in times. e^{-x^2} , e^{-Ax}
 A naive code gives α repeated a a pc of $\frac{m}{m+N}$
 $\frac{m}{m+N}$ is $\frac{m}{m+N}$ symbols a corpus

A less naive pc would be given by e^{-x^2} "factorial formula".

where we do i α 's Row no i $m-i$ a 's: so pc of α is $\frac{m-i}{m}$ i times

so we have something like $P_i = \frac{m!}{i!(m-i)!}$

so $\sum_{i=0}^m \frac{e^{-m} e^{m \cdot \frac{i}{m}}}{2^i (m-i)!} = e^{-m} \sum_{i=0}^m \frac{e^{i} (e^{-1})^i}{2^i (m-i)!}$

So, it looks like I get out of partition of P_{m-i} as i numbers!

then $i=0$; $i=m$ terms: $\frac{1}{0! m!} \approx \frac{e^m}{\sqrt{2\pi m}}$ so P_0 and P_m are $\frac{e^m}{\sqrt{2\pi m}}$
 so $S = e^m \left(\frac{2}{\sqrt{2\pi m}} + \sum_{i=1}^{m-1} \frac{e^{-i} (e^{-1})^i}{2^i (m-i)!} \right)$

$\int_0^1 \frac{dx}{\sqrt{x(1-x)}}$ $\int_0^m \frac{dx}{\sqrt{x(m-x)}}$ $\int_0^2 \frac{dx}{\sqrt{x^2-1/2}}$
 $\frac{1}{\sqrt{x(1-x)}} = \frac{1}{\sqrt{x(1-x)}}$ $\frac{1}{\sqrt{x(m-x)}}$ $\frac{1}{\sqrt{x^2-1/2}}$
 $\int \frac{dx}{\sqrt{x(1-x)}} = \ln \left| \frac{1-x}{\sqrt{1-x^2}} \right| + C$
 $\int \frac{dx}{\sqrt{x(m-x)}} = \ln \left| \frac{x}{\sqrt{x^2-mx}} \right| + C$
 $\int \frac{dx}{\sqrt{x^2-1/2}} = \ln \left| x + \sqrt{x^2-1/2} \right| + C$

$\frac{dx}{\sqrt{x^2-1/2}} = \ln \left(x + \sqrt{x^2-1/2} \right)$
 $\int \frac{dx}{\sqrt{x^2-1/2}} = \ln \left(x + \sqrt{x^2-1/2} \right)$
 $\frac{1}{\sqrt{x^2-1/2}} = \frac{1}{x + \sqrt{x^2-1/2}}$
 for $x=2$, $\ln(2 + \sqrt{2^2-1/2}) = \ln(2 + \sqrt{3.5})$

$\int_0^{1/2} \frac{dx}{\sqrt{1/2-x^2}} = \sin^{-1} \frac{x}{1/\sqrt{2}} = \sin^{-1} \frac{x\sqrt{2}}{1}$
 $\int_{-1/2}^{1/2} \frac{dx}{\sqrt{1/2-x^2}} = \pi$

$$Z \equiv \frac{m!}{\sqrt{2\pi}} \beta^{-m} \sum_{j=0}^{m-1} \beta^{mj} \frac{1}{j! \sqrt{m-j}} = \frac{m!}{\sqrt{2\pi}} \beta^{-m} \sum_{j=0}^{m-1} \beta^{mj} \frac{1}{j! \sqrt{m-j}}$$

$$Z = Z(m, \beta) = \begin{array}{ccc} m=2 & 2.0 & \leftarrow \beta = \frac{m}{2} \\ 1.3 & 7.05 & \\ \\ 1.022 & 2.9 & \leftarrow \beta = .8M \\ \\ .6 & 1.8 & \leftarrow \beta = .9M \\ \\ .52 & 1.28 & \leftarrow \beta = M \end{array}$$

$\beta = N P_{\alpha} \quad \text{so } P_{\alpha} = \frac{\beta}{N}$
 $\delta = \frac{m}{N - km} \approx \frac{m}{N}$
 $m > \beta \quad m \text{ all cases of interest.}$

1) Check to make sure, 2) Do eqns in Maple in Σ form. 3) Do eqns exactly, w.o. approx — see if it makes any difference — also see if it gives the "reasonable" results!

Seems that $\delta = \frac{m}{N - km}$ maybe slightly wrong: should be $\frac{m}{N - km + m} = \frac{m}{N - (k-1)m}$

$$\delta = \frac{m}{N - km + m} \quad 1 - \delta = \frac{N - km}{N - km + m} \quad (1 - \delta)^{N - km}$$

"Brushing δ kx under rug" say .23 is 0.14.

$$\left(\frac{m}{N - km + m}\right)^m P_{\dots} \quad NP_{\dots} \equiv \beta \cdot \left(\frac{m}{\beta}\right)^m P_{\dots}$$

Shouldn't some $m!$ term ϕ sum? I guess this is legit code: it consists of making ϕ data of x , but not using it at all.

$$P_{\dots} \left(\sum_{i=0}^m \left(\frac{i}{\beta}\right)^i \frac{m!}{i! (m-i)!} \right) + f$$

$$\left(\frac{i}{\beta}\right)^i \approx \frac{i!}{\sqrt{2\pi i}} \quad i \neq 0$$

$$P_{\dots} \left(\sum_{i=0}^m \beta^{-i} \left(\frac{i}{\beta}\right)^i \frac{m!}{i! (m-i)!} \right)$$

So do sum from 1 to m , but add $i=0$ term. $\lim_{x \rightarrow 0} x^x = 1 \quad x \ln x \rightarrow 0 \rightarrow x \rightarrow 0$.

$$P_{\dots} \left(1 + \sum_{i=1}^m \beta^{-i} \frac{m!}{\sqrt{2\pi i} (m-i)!} \right)$$

which is $\frac{1}{\beta}$. Factor of summand occurs \approx when $\beta = m - i$ for each $i \rightarrow 2+i$, since summand \uparrow by factor of $(m-i)$.

$$1 + m! \sum_{i=1}^m \beta^{-i} \frac{1}{\sqrt{2\pi i} (m-i)!}$$

Factor of summand occurs \approx when $\beta = m - i$ for each $i \rightarrow 2+i$, since summand \uparrow by factor of $(m-i)$. $\approx \uparrow$ by β factor of β . (Also \downarrow by factor $\sqrt{\frac{2+i}{i}}$)
The would be expression is not attractive when \approx mult by e^{-x}

$$1 + \frac{m!}{\sqrt{2\pi}} \sum_{j=0}^{m-1} \beta^{j-m} \frac{1}{j! \sqrt{m-j}}$$

Is it all reasonable that $Z(m, \beta)$ should be ≈ 1 for all m ?
 Growth if $m > \beta$; \approx for all m ?
 $Z(m, m) \approx 1$ clearly \uparrow w. m .

$$Z(m, \beta) = 1 + \frac{m!}{\beta^m \sqrt{2\pi}} \sum_{j=0}^{m-1} \frac{\beta^j}{j! \sqrt{m-j}}$$

Could this slow \uparrow of $Z(m, m)$ w. m be \approx for approx answer?
 I would think that $Z(m, m)$ should be ≈ 1 for all m , because it makes extra copies of i identically same PC's, so $(1 - \delta)^{N - km}$ factor, which is < 1 , must add up to ≈ 1 via $\frac{m!}{m! (m-i)!}$.

This factor is exactly $(1 - \delta)^{N - km}$ so is $\sum_{i=0}^m (1 - \delta)^{N - km} \frac{m!}{i! (m-i)!} \approx 1$?

Seems like the $(1 - \delta)^m$ factors should exactly cancel out \approx

Try is using exact summations and exact δ values, not \approx

Try coding it fully! Burn seq: no define / one of i squawks: substitution of δ derived α δ

Many ways should give $\approx \phi$ change in PC (I think). If not then I'm losing some PC somewhere.

4TM

0: (109.40)!

T. Stirling's approx is only good for $m > 1$, say, $\frac{2^2}{e^2} \cdot \sqrt{2\pi \cdot 2} = 1.919$

But 0! does not work w. Stirling!

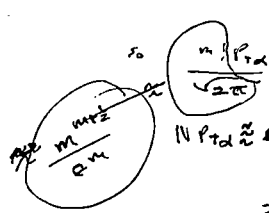
$1! = \frac{1}{e} \sqrt{2\pi} = .922$... not bad!

In 109.25: $P_{Tot} = \sum_{m'=1}^m \left(\frac{m}{N P_{Tot}}\right)^{m'} \frac{m!}{m'!(m-m)!} = \frac{m! P_{Tot}}{m'!} \sum_{m'=1}^m \left(\frac{1}{N P_{Tot}}\right)^{m'} \cdot \frac{(m!)^{m'}}{e^{m'}} \cdot \frac{1}{m'!} \cdot \frac{1}{(m-m)!}$

Since there are no neg terms, $\sum_{m'=1}^m$ (never zero)

$\left(\frac{m'}{e}\right) \frac{1}{m'!} \approx \frac{1}{\sqrt{m' \cdot 2\pi}}$ $\left(\frac{1}{e}\right)^{m'} \cdot \sqrt{2\pi} \approx 2!$

0.9



$\sum_{m'=1}^m \left(\frac{1}{N P_{Tot}}\right)^{m'} \frac{1}{(m-m)! \cdot \sqrt{2\pi}}$

$z \approx \sum_{m'=1}^m \beta^{-m'} \frac{1}{(m-m)! \sqrt{2\pi}}$

$j = m - m'$ $m' = m - j$

$z \approx \beta^{-m} \sum_{j=0}^{m-1} \beta^j \frac{1}{j! \sqrt{m-j}}$

$\beta^m \sum_{j=0}^m \beta^j \frac{1}{j! \sqrt{m-j}}$

So by $\sum_{j=0}^{m-1} \beta^j \frac{1}{j! \sqrt{m-j}} \approx \sum_{j=0}^{m-1} \beta^j \frac{1}{j! \sqrt{m-j}}$ \Rightarrow a function of m .

$m \leq 1$ perhaps doesn't occur $z < 1, m \geq 1$

$z(2,2) = 2^{-2} \sum_{j=0}^1 2^j \frac{1}{j! \sqrt{2-j}} = 2^{-2} \left(\frac{1}{\sqrt{2}} + 2\right) \approx .675$

$z(3,3) =$

$m^m (m+1/2) \exp(-m + 1/2 \ln m) \approx 2.506628$
 2.506628

to evaluate $z(m, \beta)$

$A_j = \frac{\beta^j}{j!} : A_j = \frac{\beta}{j} \cdot A_{j-1}$

$A=1 : s = m^{-1} (-.5)$ for $j=0$

For $j=1$ to $m-1$

$A = A_{j-1} \beta / j : s = s + A_{j-1} \beta \cdot (-.5) (m-j)^{-1} (-.5)$

Next

Print $M, B, s \cdot B^1(-m)$

! B35! 4TM III, BAS
 $1! \approx \frac{1}{e} \sqrt{2\pi} = .92213$
 $\times \sqrt{2} = .99898$
 $2! = \frac{4 \sqrt{2 \cdot 2\pi}}{e^2} = \frac{8 \sqrt{4\pi}}{e^2} = 1.919009$
 $\times \left(1 + \frac{1}{24}\right) = 1.99896$

M can be quite large: K is fine that $m \leq 1$ ppm

Even w. $\beta = \frac{1}{2} M$, and large M w. β very small z for β (big M !)

woops! I left out the $m!$ factors.

$m^m \approx \exp(-m) \approx (2 \times 3.14159 \times M)^{-1} (-.5)$

putting $\beta = m/2$: for $m=2$ to 20 $z = 3.2 \times 10^{-17}$

I have to multiply by $P_{Tot} \approx \frac{1}{\sqrt{2\pi}} \approx 2.5 \dots$

P_{Tot} can be quite small if m is large. It is + original cost of deriving α - m's \dots

If only the P_{Tot} α is not used: Part one could give α (of α \dots)

divided by $\sqrt{2\pi}$ we get 1.3 to 705

for $\beta = .8 m$ P_{Tot} goes from (.022 to 2.9

for $\beta = .9 m$ P_{Tot} =

$m=2$	$m=6$	$m=20$
.6	1	1.8
$m=10$	$m=20$	
	1	1.28

for $\beta = m$ P_{Tot} = .52

It would seem P_{Tot} for $\beta = m$, \dots the deriv. of α should not help for any values of m !

Also, it would seem that for any $\beta < m$, P_{Tot} \dots $\frac{\beta}{A} = \alpha \approx c < 1$ Rem

for large m and M , z should be truly large.

$\beta \approx N \cdot P_{Tot} \therefore P_{Tot} = \frac{\beta}{N} < \frac{m}{N}$

also: look at $m=2, 3, 4$ etc (common cases) - w. large k , hvr. k will be m.f.

Re: Mocha CC for colns: For PPM say output is 2^5 , say w. u.p. results -
even w. PPM w. large k allowed, m or t large ones would vary infrequently
"exact"

So perhaps my original factorial version of P_{T_x} -

~~At~~ If its easier to calculate using "mind" version, using a $(1-g)^m$ factor
to compensate for the defn. of α .

While 108.40 might be ok for obtaining P_{T_x} , ~~it~~ it seems very expensive. On the other hand,
using α case costs of k factors to 2^5 in factorial (the Prod of 2^5 etc) would seem to be catastrophically

poorer (because of commonly accuracy "correct product") than 108.40! So 2^{141} may be better
than 2^5 but probably much worse than PPM! (15)

T. Expense of 108.40 isn't altogether ~~expensive~~ outrageous! We will be keeping track of 21 combinatorial
PC's ordinary pointint corpus. This might be done using a kind of "floating pt.", in which 2^5 is
fixed pt. multiplication, but I keep track of shifts needed to keep mantissa ≈ 1 .

One could "do" 2^{141} , but use 108.40 for evaluating P_{T_x} . Perhaps the system would
start out pretty much like the old Z41, but gradually improve as the predictions improved P_{T_x} was
computed w. more accuracy. This is perhaps the best way to do "Z41" any way!
We may start out using the factorial formulas

I do need to get that 2^{141} coding straightened out! I will probably need it when I do Grammer disarray
to segment "Phase 1". I think I'll need it for (13) even! Perhaps by preloading Sol 6!

4.14.04 I can "sort of" see why I felt that OOPS' language wasn't as good as Lisp-like AZ141 -
IA AZ (41), the sub-sequences always (or usually) have maximums as job-functions.

In OOPS's forthish lang, this does not seem to be true. One could, perhaps, in OOPS
adopt a very forthish programming style, in which ~~the~~ (most) all parts were shown series of
dictionary entries - in which cases, the desired context structure utilized by PPM
might work o.k. In fact PPM might impose a "style of programming" on the system in a way so that
PPM would be useful (a sort of self-compensating hypothermia) - in which case
one could use (perhaps) a (most) any universal formalism that could express a desired solns
as "short codes". Hvr, w. some universal lang, it may take a much longer corpus;
before the system begins having useful PPM-type regularities. -> On the other hand

P_{T_x} does not seem to be implemented by a suitable initial TSD.

SN "To first order" PPM does every thing that does do, but cheaper -
Hvr, on "second order" (advanced lang) PPM may be unable to do the equivalent of
defining α constants to obtain the implied "grammar".

30

20:10840 :

If α is assigned prob δ , all Tokens will have their prob mult by $(1-\delta)$.
 α contains k tokens. α occurs m times in (augmented by α) corpus.

The joint prob of k pc's of α symbols, k tokens of α in corpus is P_{α}

There are, under instr of α include α itself "1"; N tokens in corpus.
 Factor of change of pc of corpus brought by use of α , m times is:

$$\left(\frac{\delta}{P_{\alpha}}\right)^m \cdot P_{\alpha} \cdot (1-\delta)^{N-km} \approx e^{-\delta(N-km)} = e^{-\delta(N-km)}$$

$$\left(\frac{\delta}{P_{\alpha}}\right)^m \cdot P_{\alpha} \cdot \frac{e^{-\delta km}}{e^{-\delta N}} = \left(\frac{\delta \cdot e^{\delta k}}{P_{\alpha}}\right)^m \cdot \frac{P_{\alpha}}{e^{\delta N}}$$

$$\delta^A \cdot (1-\delta)^B = \text{max}$$

$$\frac{A}{\delta} \ln \delta + B \ln(1-\delta)$$

$$\frac{A}{\delta^2} = \frac{B}{1-\delta}$$

$$B\delta = A - \delta A$$

$$\delta = \frac{A}{A+B}$$

So perhaps assign value to δ so that best fit of pc is max

T. values dependent on δ is $\frac{\delta^m}{e^{(N-km)\delta}} = \frac{\delta^m}{A^{\delta}}$

$m \ln \delta - \delta \ln A = \text{max}$ $\frac{m}{\ln A} \cdot \ln \delta - \delta = \text{max}$

$\frac{m}{\ln A} \cdot \frac{1}{\delta} - 1 = 0$ $\delta = \frac{m}{N-km}$ $\ln A = N-km$ $\delta = \frac{m}{N-km}$

$e^{-\delta(N-km)} = e^{-m}$

So $\delta = \left(\frac{m}{(N-km)P_{\alpha}}\right)^m e^{-m} = \left(\frac{m}{e^{(N-km)P_{\alpha}}}\right)^m P_{\alpha}$

So we get factor of δ again!
 I assume $km \ll N$ so $\delta \approx \frac{m}{N}$
 Prob it is correct out!!!

$\left(\frac{\delta}{e P_{\alpha}}\right)^m \cdot P_{\alpha} = \frac{m}{N-km}$

$\left(\frac{m}{N P_{\alpha}}\right)^m P_{\alpha} \left(\frac{m}{e}\right)^m \approx \frac{m!}{\sqrt{2\pi m}}$

Next, we have to consider all possible strings of m tokens.

$P_{\alpha} \sum_{m=1}^M \left(\frac{m!}{(N P_{\alpha})^m}\right) \cdot \frac{m!}{m!(m-m)!}$

$\approx P_{\alpha} \sum_{m=1}^M \left(\frac{m!}{N^m P_{\alpha}^m}\right) \cdot \frac{e^{m'}}{\sqrt{2\pi m'}} \cdot \frac{m!}{(m-m)!}$

$= \frac{P_{\alpha} m!}{\sqrt{2\pi}} \sum_{m=1}^M \left(N P_{\alpha}\right)^{-m} \frac{1}{\sqrt{m!} (m-m)!}$

for $\sum_{j=1}^m \beta^{-j} \frac{1}{\sqrt{j!} (m-j)!} = \sum_{z=0}^{m-1} \beta^{z-m} \frac{1}{\sqrt{m-z!} z!}$

$= \text{some function of } M \text{ \& } \beta$

I could plot it for various β, m . — But β will be $< M$.

So do this for various large m , $\beta = .9, .8, .7, \dots$

Go over these eqns carefully before plugging! Getting rid of $\frac{1}{e}$ factor is a great relief!

Don't move directly w. $(1-\delta)^{N-km}$ using second order formula
 $\delta = .23 \rightarrow .25$ can be used for x ! $\rightarrow \left(\frac{SP_{\alpha}}{100.00}\right)$

00:107.40: Q: Is t. kernel width of t? Insertion into L at various pts, will occur w. different t values
 Hvr, at each t value, t. total length of L is $\approx t$. The width of kernel will always be a certain fraction of
 t. $\text{freq} \approx \text{len} L$ ($\approx t$). If we have a fixed L (or a L of fixed length - obtained
 by deleting a card whenever one inserts one), then we have a kernel of fixed width.

[SN] Adding 2 new tokens is certainly not at the beginning — yet it ~~does seem to screw up~~
 the needed shape of t. kernel! Another disturbance is later: when a new token is
deleted (as in OOPS). In t. ~~early~~ first case (early ~~time~~ life) we have to give equal pc
 to several possi. symbols — can't be done w. a "pro corpus", because symbols in pro corpus
 are at different distances... (?) no! — those are distances in corpus space (t. direction)
 we are concerned w. "L" spaces, ... E_{L^2} directions

For t. design of sequences of instructions; we can have t. "pro corpus" in the form of k distinct
 cards: $\Delta T_1; \Delta T_2; \dots \Delta T_k$. T_i is the token (or C): $\Delta T_1, \Delta T_2, \dots$ (Δ is null context)

We could probably add in definitions as well, in this way: a bit Δk , hvr.

Error of t. kernel is probably second order effect, but t. relevant error is not log of t.
 kernel, since we are interested in maximizing t. pc of t. foly — t. product of pc's of tokens
 = $\exp(\sum \ln(\text{tokens}))$.

t. idea of this "kernel" (has one apparent flaw), T. linear ordering of t. tokens ~~is~~ in kernel

pc a lot. When one sorts t. corpus "locally" $\rightarrow L$, P_i is automatically ordered. (3) is another way to deal w. it

One way to avoid this — The only measure of similarity of 2 contexts is t. length of their common
suffix. This seems to be odd in conventional PPM is PPM*. The relative wts of ~~some~~ phrases

common suffixes is (perhaps) shared by all suffixes. An alternative way would be to view each } is this
 $\approx \sum (1/P_i)$

suffix as a PC, — a form seq. of a pro corpus (or w. $\frac{1}{2}$ probab). So t. past success of t. suffix is
 used to help get its wt. \rightarrow (33)

Q: To what extent is this "escape" probability close to that obtained by using a "pro corpus"?

In both cases t. pc of t. ^{new} token is a best $\frac{1}{t}$ (T being no. of tokens, P_{i+1} for L or $\frac{1}{t}$)

(Apparently) t. city may do this, \rightarrow to derive k ~~tokens~~ (there are k bit first tokens) & id'nt

So the answer: for each context, i consider t. best way to code t. corpus using ~~the~~ best definition

In t. PPM discussion of ~1995 (written, Chao, Mittal) they use a "Tri" structure that somehow
 seems to avoid this error of low order — but it may be more expensive in ~~time~~ time and memory

Expanding this a bit: ~~to~~ \rightarrow code using t. α ~~pro corpus~~ \rightarrow a particular suffix ending in "2"
 t. suffix is α of length k.

First define α in t. pro corpus. pc will be \approx product (pc's ~~with~~ with corpus) of α tokens.

Also, ~~the~~ pc of symbol domain "definition"

In Pro corpus we can compute the pc of each sequence of tokens replaced by each α .

This is done by looking at t. cumulative pc of t. corpus at t. beginning and ends of each α insertion

APM



10.10.06: Hvr, I really have no good idea as to how long it takes to generate a corpus card. Currently, it takes ~ 2000 ind ~~to~~ clocks to put ~~1~~ token into \hat{D} . Also, ~~the~~ $(a-10)$ & by of Ram ~~into~~ by Enuf for all that I need to do.

In $T \approx 2T$ (which seems simplest to pym) + Bottleneck seems to be the comp. off. pc's factors.

Try to see how much accuracy I need in various when pc's.

As is, I don't have a clear idea as to how much heavy I need, how fast, how good a TSQ is needed to get to ~~Phase 2~~ Phase 2.

For ~~the~~ Token re calcus: I may be able to ~~sim~~ simulate kernel by Laplace xcu, Reu convolution using $x_n = \sum \alpha Y_{n+1} (1-\alpha) x_{n-1}$ - using many α 's in $(-)$ is positive or negative. This could be ~~linear~~ in $SS2$. As is, I have no idea as to what the kernel looks like.

I could just accumulate an empirical kernel, as 103.05 : (hvr) I had idea that it was a kernel kernel!

One kernel for each of tokens types!

(Thinking out loud): When one token asks for the pc of a context: we look at longest historical matches in corpus. - How far along "do we have to go back we find a particular token? - ~~the way to that~~ to L sequence of tokens can be treated as a sequence of "Token vectors"; (A token vector ~~is~~ is a basis vector in Token space $\approx k$ dim. space historical tokens) ~~a seq. of vectors~~ So we convolve L. w. f. (smooth kernel) ~~is~~ we get ~~sequences~~ sequences of (smoothed) vectors.

To store these smoothed vectors would require $k \cdot N$ bytes or words. (N is no. of tokens in corpus).

Hvr, because of the ~~smoothness~~ smoothness of kernel, it shouldn't be very to store so much info (smoothness like only 2-4 points per second for a "bandwidth signal"). If the kernel is of width W , it may be easy to ~~store~~ store every $W/2$ point of a suitably filtered form of L : perhaps a "boxcar" filter of

width W or $W/2$ would be good.

1) ~~the~~ say kernel is of width w . - store every w or $w/2$ points of ~~the~~ filter of L .

2) we don't need $\frac{N \cdot k}{w}$ points at all: Because of filter compression as we go along, - the kernel "sort of" puts w. dev. as we move along & ~~step~~ step. It may be that to use of pts we need to store is indep of N (!). If the kernel is ~~of~~ of width w then maybe we need only (w) vector ~~pts~~ ~~pts~~.

So its ~~imp~~ imp to know just what fraction of width of kernel is of the corpus, N . - will vary in how accurate ppm is for that kind of corpus. width will be ~~or~~ or bits/char compression perhaps.

3) Actually, the entropy of the corpus depends not so much on the width of the kernel, but

on how rapidly the corpus shifts from emphasis on one token, to emphasis on another.

4) Consider the following ~~form~~ form: $(Token_1 \text{ m times})(Token_2 \text{ m times})(Token_3 \text{ m times}) \dots (Token_k \text{ m times})$. ~~usually by the most likely one~~

If the width of kernel is $\ll m$, substrings will be short. The correct token prodn will ~~be~~ be the ~~most~~ most likely one. Or a times, it will be the ~~second~~ second most likely one. with only 2 symbols, the entropy per symbol is < 1 bit/symbol.

(and actually < 1 bit/symbol). It would seem that BW's method (or BZ type) with, in the present

cases, do better than 1-kernel - i.e. it would pick the closest correct context - which would be correct, except at boundaries between symbols (i.e. transitions between symbols). BW would give very low PC's to

tokens at transitions, because they ~~are~~ are not the most recent, & haven't occurred in a long time.

With $k \approx L$ of form .29, we can compare which kernel would be liked! I think it would be very narrow! say 3 sybds wide (apparently indep of position in corpus!!) $\approx ?!$

Another Q: ~~Does~~ Does a corpus exist \Rightarrow ~~its~~ its L is ~~sub~~ sub. in .29?

29

4TM

18:10β
-19

23ENE 2 24K Books

0106.40 : get this local concentration of cc. The rate of which this occurs will depend much on how fast our algm for β & pc's of contacts of β & γ cards. This algm should depend on how near to β of γ card is to β . best β 's thus per — so ~~as we climb hill~~ as we climb hill, pc's cards not longer in the "top percentages" will actually \downarrow — β & γ pc's assoc w. their contacts will \downarrow .

O.K. : so write up all 3 Lenn Models in a little detail. Devise suitable notation — to describe them & to facilitate programming.

Q: Is there any secondary method of \bar{L} in its applic, assoc. w. fact that it's β & γ will be useful? Short — \leftarrow no on 20 chars? "d" will occur often in the "corpus" in \bar{L} . I think there is a relevant record decrease in deciding if β needed \geq symbols β & γ for start & stop. See 97.14 — 98.14 At 98.12-19, I think I had "Boundary Control".

[SN] In testing to present system or problems others how work on: Consider that our system will work \leftarrow if it has good, reasonable contacts: These can be acquired by ^{suitable} initialization or by a suitable Pre-TSQ. It has to acquire functions, etc, that can easily be put together using a PPM. Presumably, after it has solved a suitable large batch of problems, it will have to have such a set of functions, pre-TSQ. The set of pre-TSQ will be a function of the initial set of parameters, and of TSO. Perhaps it is poss. to design a set of primitives \rightarrow PPM is adequate for almost all prob solns. (it was started w. a suitable "TSQ").

One Q is: to what extent is CFG form used? We do use $If \alpha then \beta else \gamma$ or just $If \alpha then \beta$. The boundaries of α & β have to be some way indicated. We have a function of α & β : This is easy to do in LISP, (I guess) — $If \alpha$ & β were given, then If in Basic is a ways is for sequence control. If α "IF" = false, we go on to next step. In Lisp I guess If is a ways function — a switch tests 2 fences.

How would it be done on a stack machine? Maybe put α & β on stack (2 slots), then do If : We can end up w. either β or γ on TOS.

It may be best by giving β a large no. of probs to solve first to whose solns have some common features, just ~~that~~ PPM would have to be able to use these "common features": On the other hand, it may be necessary to see a detail how 4TM is supposed to solve various problems in the TSO.

When we find a soln to a problem, we should always continue to look for higher soln (if poss). (Since w such in cc order, it is will have to be a soln of higher cc as well.) To reason why β is good, is that it will give good "concepts" for subsequent problems. Also simply having several solns to a problem is usually an economical way to get good / new cones & good wts for cones.

So: In drawing up the detail. $F \leftarrow ZT, 11$, random models: Cost of many can be a contributor too. cc. To extent to which disc 5MB/yr can be used is unclear. See 101.09 for discussion, Numbers. 40 MB/sec write, 20 MB/sec read. $25ns + 17ns = 42ns \rightarrow \approx 25MB/sec$ for read/write

19:14 ~~Ready~~
20:13 "

4TM

o: spc
4.0

It would seem that ~~the~~ Random Search would not be v.g. in this respect, because a change in concentration of pc of part of a system and even would seem to not produce an immediate effect in Random Search.

If Search (of perhaps T <= T Search) would have to be more scholarly Candidates: we would probably want to have a special set as soon as it is by G and is found. This might be possible in Random Search as well!

There is a particular variety of GA that does local search of this kind - Has a particular Name - (published in Italy?)

It would seem that normal search would be ^{not much} better than Random Search in this respect. Since an A in pc (of a particular region of candidate space) merely \uparrow the time share of that region - but this does not normally produce a rush to investigate that area.

In T <= T Search, if we try to do it in a pc order, we could get this "localized investigation of ~~the~~ by G pts".

On second thought, \uparrow in G doesn't really \uparrow pc of a region - But it ~~is~~ is able to change order of trials!

So: Present Prob: In view of 104.37, we'd like to be able to get a quick feedback from successful ~~the~~ Cands. (of local search).

Would it be a good idea to write up a more detailed descn of each of ~~the~~ ^{Search} systems. (T <= T, //, Random).

I would also like to refer to $\times 2 \times D$ speed of // a random over T <= T, if possible

At first glance, it would seem that the normal methods of doing // a random Search would be correct. But if a particular cand led by G is we wanted to give more wt. to it, we could do this by (slightly) increasing the pc's of adjacent contexts (via Row 6 structure). - So maybe we'd double or quadruple the pc of a cand ~~is~~ \downarrow by 3 or 4, the cc of solns (if any) near to that cand. After, this is a "business as usual" type of soln, we'd want a more "local search" type of reaction. In Phase 2, I expected to use long domy Search

~~the~~ ^{Big} ~~the~~ - It may be that this "immediate RR" cannot be used so effectively in Phase 1. But we have to get to Phase 2 to do it properly.

To put ^{more} emphasis on ~~the~~ to "neighborhoods" of by G cands, is the kind of "elitism" - Its values is uncertain.

One way to get this "local concentration" done by G cands: If a by G cand has occurred, in Search, the "prob" threshold of the system must be at a certain level; ^{just} ~~near~~ ~~to~~ ~~the~~ ~~level~~ ~~of~~ ~~the~~ ~~by~~ ~~G~~ ~~cand~~. If we \uparrow the pc's of nearby cands, ^{as well as out of it by G cand,} we should ~~the~~ ~~cc~~ ~~on~~ ~~mean~~ ~~to~~ ~~bring~~ ~~them~~ ~~up~~ ~~to~~ ~~the~~ ~~new~~ ~~level~~ - which will probably cause some of them to finish & be evaluated. ~~the~~ ~~cc~~ ~~on~~ ~~mean~~ ~~to~~ ~~bring~~ ~~them~~ ~~up~~ ~~to~~ ~~the~~ ~~new~~ ~~level~~ - which will probably cause some of them to finish & be evaluated. If these cands have by G, this will further \uparrow the pc's of nearby cands, etc. So we could \rightarrow (06.00

: A possl. serious expense in Random Search is ^(Nos) generation of random ^(Trials). In most cases we don't need v.p. random nos., & so cheaper sources could be considered. Another possy is H.W. generation of random nos. At present, it isn't clear to me & (1) How many bits of randomness I need for each token choice, (2) what kinds of non-randomness are tolerable. — It may be that an exhaustive (non-random) seq. would be o.k. // See how much time it takes for PB35 a/o PBCC to generate random Nos. (subtract time to do administration; by timing a "null" loop.)

T. big Q's Time to do random search v. time to do || deterministic search

Also, just how much time is needed for pc. calcn. in $T \leftarrow 2T$ search?

My impression is that usually for every random method, there is an non-random method that's about as good or better! — but not really easy to find!

Q.K. How to get pc's exactly; maybe not so expensively! At each point in "L" we should.

Cumulative number case counts of each of ~~tokens~~ ^{TOKENS} (since + binary of L). ~~After~~ After we have

No. of bits is a dem. of L. \approx constant kernel.

been running a long time, $\frac{d}{dx}$ of those case counts will be about constant, & will be our

"kernel" — So we average out a "kernel", which we update every once in a while.

Still, it's not clear how we get token pc's out of this!

T. density distribution of id - all $\approx \left(\frac{d}{dx}\right) \Rightarrow$ over TOKENS ; it is quite different

t. "kernel"

"L" is motionology "L", say $P_L(x)$ is t. density at pt. L of ~~token~~ ^{token} τ .

$K(x)$ is the kernel at L (it is \approx indep. of L). The pc of symbol τ is, $P_{L\tau}$ or:

$$P_{L\tau} = \int_{x=-\infty}^{x=L+\infty} P_L(x) K(x) dx$$

say $N_L(\tau)$ is t. total no. of times τ has been ~~produced~~ ^{correct} following symbol, upto point L.

$$P_L(\tau) = \frac{d}{dL} N_L(\tau)$$

T. quotes are because it's not a derivative but a "difference",

$$\approx N_L(\tau) - N_{L-1}(\tau) \text{ is exact form, } \dots$$

13:12

So, we do know P_L function \vec{P}_L (t. vector w. token components) empirically as t. identity of t. tokens for each L value. The convolve of t. kernel $K(x)$ & t. density \vec{P}_L can be done rapidly by fast Fourier xfm. It takes time $\propto N \ln N$.

couple of "Notes": (1) T. convolution is essentially a "smearing process" & doesn't have to be done very accurately (for conf. generation, at least)

(2) It totally is ~~totally~~ a convolution because of kernel with N values \approx is $\propto L$.

Re (2) I ~~was~~ don't know the L distribution so that t. kernel was constant. My have to modify amplitudes w/ diff. dilatant of "L" direction, \rightarrow I think this will work. The dilated points need not be put at integral L values

Unforty, this makes t. FFT ~~imposs.~~ (perhaps), I could move from to \approx closest lattice pt. \leftarrow yes,

— the kernel doesn't have to be so exact.

37

SN It would seem that immediate Feedback sub-~~st~~ during L search (would could) be very imp! — i.e. one finds an apparently ~~way~~ by of direction & pursues it — T. alternative is to have to wait until t. next $T \leftarrow 2T$ round (or equiv. delay).

If this F.B. usually is by CC, then perhaps only data relevant "update" when ~~unusually~~ by G-card has been found. \rightarrow 105.00

4th

SM: On Action Algorithm Problem: -11 -29 AAP

0 : ON CROSS VALIDATION: This is one ^(concept) ^(bottle neck) of f. Univ. D.F.: It is not clear how to divide corpus to employ it. - i.e. (understand) ^{modeling} ^{test set}: One way: Use all possibl. divisions of N into k+l: We will get 2 main error ~~maximize~~ ~~over~~ test set for each k, l pair (assuming "Best fit" for each ~~set~~ set of models for each k, l). So we have curving this ~~maximize~~ ~~over~~ ~~test~~ ~~set~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~ ~~is~~ ~~not~~ ~~the~~ ~~same~~ ~~value~~ ~~for~~ ~~all~~ ~~models~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~. So we have curving this ~~maximize~~ ~~over~~ ~~test~~ ~~set~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~ ~~is~~ ~~not~~ ~~the~~ ~~same~~ ~~value~~ ~~for~~ ~~all~~ ~~models~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~. So we have curving this ~~maximize~~ ~~over~~ ~~test~~ ~~set~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~ ~~is~~ ~~not~~ ~~the~~ ~~same~~ ~~value~~ ~~for~~ ~~all~~ ~~models~~ ~~for~~ ~~each~~ ~~k~~ ~~l~~.

In many cases, it is not so easy to divide up data this way, But in some, if data is unstructured, One can pick k ^{of} modeling set from N in $\frac{NS}{(N-k)!k!} = \frac{N!}{k!l!}$ ways. Total no. of ways = $\sum_{k+l=N} \frac{N!}{k!l!} = 2^N$. This "size" may be so large that one way wants to use only part of it - using "Random Sampling". \rightarrow 13

I vaguely remember doing something like this, but ending up in something equivalent to using all of data for modeling test set ||

11 12 EN On t. A. A. (Action Algorithm) Problem: Perhaps Goodway! Try to get Pd of t. A. A. over table. This should be directly summarizable to UDF methods. \rightarrow 29

13 07 If Random Sampling is done, one might concentrate it on small l and l=1,2,3 etc. This may be a common way of using Cross Validation. (e.g. l=1).

14 10021 More on RANDOM Search: At first keep track of which cards have been "hit" (by t. Macro) better. But don't get stuck in local minima. If Ray ever hit a perm within a certain time, ~~then~~ we start storing such states. This means we don't consider cards w. PC's below a certain level. This level will depend on how much RAM we have (assuming we don't use Disk Storage). Actually, Rays may be quite expensive - we may be quite happy using only 10K most probable cards (Roberz sounds like a Bad kind of SUMAC!).

16 As for keeping track of cards: we already have cards in Lex order (?). So we can have a register that tells where the params of the SM are.

28 Quick jobs: Root of cards parameter of cards could be rather cheap: but computation in deterministic Search is slow. 34

29 12 SM: Say one has a strategy w. certain param(s), giving yield γ_0 for t. corpus.

In param space, take random pts about $\bar{\theta}_0$ & see how narrow to $\bar{\theta}_0$ regions that gives $\geq \gamma_0$ yields.

Also, given yield γ_0 : How large is region in all of the space that has yield $\geq \gamma_0$? What is say it is $V(\gamma_0)$. From what is $\frac{dV}{d\gamma_0}$? - t. density around γ_0 ?

34 29 One way Mt corp would work: We keep pers "L" in Lex order: we start w. initial symbol Λ - we go to put Λ in L and we get a distribution of possibl. following tokens, in t. "L" direction (i.e. "directions" Lexical direction). If t. corpus α of size N, then we pick a random symbol to follow, using a lookup kernel of width α N; in t. "L" direction. For initial symbol Λ , we can only look in t. "positive" L direction.

As we generate more symbols of f. card, we can look in both directions more. Having chosen one token, say X, after Λ we look up ΛX in t. file L & get (probabilisticly) next symbol, etc. \rightarrow 10300 initials specifically 103, 25

HDD Speed of Read/Write

4PM

For 8000 ~ 30MB/sec read or write for 7200rpm IBM 70G.

chris user ~ 25 MB/sec. Use of caches for read & for write is crucial.

Saw some other w... ~ 40 MB/sec. This is at least 16 M by ~~other~~ chunks.

Wedge Corral w R200JB dH 68 MB/sec w. 128b chunks
51 " " " " 84b "

So assume 60 MB/sec. ~~ns/byte~~ This is 85 clocks = 5 GHz CPU probably. it will take > 85 clocks to do a function?

Feb 2004 Qn another work page 40 16 MB/sec write but up to 60 MB/sec read, 70 70 MB/sec read. $(\frac{1}{40} + \frac{1}{70})^{-1} \rightarrow 25 \text{ MHz (read/write)}$
So, for Qn, can't disc read bytes to CPU factor Read/Write Process Read?

If, say 4 times as fast, we can do swap & compute sequentially

Otherwise, we've either by doing it in 11: To some extent, Disc is slow, Disc has a big Cache. We have to do a read & write to swap.

How, if the instructions are all forward, the time for a read will be of time of instr.

In the more direct the performance TM, we will be able to build more complex storage systems

Another possibility is to have several HDD's in parallel feeding the CPU: Or pick off data from several read heads in 11 - (More than is currently done: Chris says 16 or 32 heads in 11).

How: I will probably start w. TEST, then look into 11 or Random Look if it looks like the factor of 2D could be achieved or even a factor of 2 misses cc!

We write arrange so that ~~the~~ time spent on each read would depend on rate of write with swapping occurred.

If ~~the~~ read/write cache memory is properly used, the disc should be containing data either reading or writing at all times. (Almost all - since read rate is write rate.)

Just how the HDD system manages heads read rate is write rate (is unclear) - i.e. it would seem that both would have to be synchronized with bytes going by on the disc; at a constant rate!

unless every other bit is written, or every 4th bit is written. or every 16th bit is read.

The No of heads on a HDD can vary (a lot); but normally only, say 8 (or 16?) or (32?) might be used at one time. Perhaps 64 could be done w. a 64 bit CPU?

Empirically, several recent of many HDDs have 16 heads: 0 Idarones 10CS (15 for comm!) 5, 6 etc.) How 16 is 4 times I've seen. To 5th Converse: I don't know how many heads they have.