

S F T M

387-388

01

Some ideas:

1) A β netm (Boolean Net TM) is ~~some~~ sort of Boolean Net w., at any particular time, a particular $f, 0$ relationship. This relationship can be modified ~~to~~ to be various "near by" relationships by suitable ^{total} mutations.

To operate it: we try various mutations ~~then~~ Test each according to some criterion. If it is "good enuf" we return to mutation - if not, we try a new mutation from f . old pt. as center. If "good enuf" (or even just better than before) we do new mutation from this new pt. as "center".

2) Q! If we have a ~~net~~ net that is "Universal" in some sense, is it always poss. to endy place it at some pt. \Rightarrow its "closeness function" is any particular function that one wants? (The usual ones for induction are those involving some mixture of cost & profit).

20

3) T. system of 1) isn't really so good. We would perhaps like to try a mutation ^{then} obtain its Govc. ~~It~~ We do this for many mutations. We remember t . Mutation, Govc Pairs, hvr . T. set of best mutations are then tested further, using a Govc of greater cost, a t . Best ~~mutation~~ is then retained & used as a base for new mutations.

30

\rightarrow I think a ~~net~~ system that could do 3) was the one with a very large, memory of relatively slow access time, used by many n units in 11 , each having its own \approx reasonably sized cache memory. [I wrote much about such a device in t . first $\frac{1}{2}$ of 1973] It can do Back tracking if needed.

So 2) is a Very Imp't Q.

\rightarrow

In t . system of 30, since each n unit is a sort of unc., we should be able to implement any "closeness" funct. we like. \therefore Perhaps just use profit, & ~~we~~ make cost part of t . Govc! - Tho it would be nice if we could reject by cost trials before even trying them!

N 15.73 SFTM: .01:387.40: In this 387.30 system, we can save all of t. cache unit memys in "archival" (unchangeable) memy, so we can track arbitrarily. Also, we may want to go back & actually have "offspring" machines reconstructed, having & starts at certain "Branch pts." of a voln. of t. system.

.05 4) Note, hvr, that 387.30 is not at all SFTM. SFTM was a very large β net w. very fast transistors in it as components. It was very fast, asynchronous, perhaps partly analog, is very inexpensive

5) Humans may implement 387.20 by using & abovm short term & long term memys

6) T. System of 387.30 may be adequate for most any TM, that I will initially consider. After I get one running, I can try to get it. micro etc into SFTM form (like .05).

7) ~~More~~ ^{almost purely} A/sequencial machine usually gets more computing done per cost. Probably a sequencial with unit (1 bit at a time) is best — particularly inasmuch as one need only use as many bits accuracy as is needed. Another big waste is "look ahead".

Planning

944-851

73TM

UD 2073

Planning
Discn w. G. Sussman

See also his "Outgrowths of some ideas"

Given a problem: To ^{Plan} express the problem in

a form \rightarrow A "solution plan" is suggested.

A "solution plan" is a very large set of poss. solns. It takes the form of a rather imprecisely defined soln.

Solving the problem then consists of narrowing down the poss. solns.

One common method is to take a poss. soln. that conforms in a central way to the plan. Then, after trying the soln., it is clear that certain sections of the soln. do not conform to the plan. Also the way they don't conform suggest modifications to make them conform.

In the large ~~plan~~, dozen, the plan was divided into sections (sort of "subgoals"), so that one could try to get each section of the ^{proposed} plan/soln. to conform to the corresp. section of the plan.

If certain subgoals appear, after a while, to be impractical, or if, all subgoals are achieved & the soln. is still no, we go back & try to devise a new plan.

I think the big ~~discn.~~ discn. is too much tied up to the idea of subgoals. The idea of a plan can be more genl. than that, & I should write a more genl. discn.

01 D2473 : On Planning in general: A Plan by definition,
02 encompasses to a (usually) set of poss. solns. to a problem.

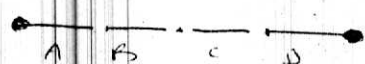
T. Language of Plans is quite imp. - ^{at least} 2 aspects of imp.

(1) Easy to reduce plan from problem. This is based upon good ways for plan to prob. Also ~~is~~ relevant experiences in the "area" of the problem. So - from a desc. of the problem in problem lang. \rightarrow a desc. of one or more plans in plan lang.

(2) Easy to solve the prob. once ~~you~~ one has a plan in the plan lang. This could be because the plan lang. facilitates natural ways to break prob into parts (i.e. sub-goals is one way to look at this -

e.g. problem \rightarrow plan \rightarrow plan \rightarrow sub goals; then each sub-goal \rightarrow plan \rightarrow subgoals or goal \rightarrow soln. or goal \rightarrow plan \rightarrow soln. (Some goals or plans are very close to solns. for a properly experienced TM). partial net


A prob can be broken into subgoals in a non structure (I don't know if this is most goals) as in SP.

A subclass of this is ~~more~~ ^{very} for it to be broken into a set of problems,  that may

or may not be sequentially constrained (e.g. the exact nature of B can't be determined until A is solved)

T. defn. of "plan" 01-02 seems fairly general.

The idea of ~~its~~ facilitating the breakdown of the problem into a task net (i.e. SP) ~~is~~ seems rather non general, but ~~is~~ ^{has} perhaps received all or almost all the attention of AI workers (perhaps following GPS) up to now.

A "non-wooly task net" concept of working via a plan:  working of the problem consists of narrowing down the specificity of the plan, which

≡ initially a large set of sol. sets. ~~Minimization~~

In task net approach, the plan is a sort of "factorization" of the set of solns. One solves each factor by narrowing it down. More generally, factorization

How to construct plans

Task

2 degree

to cause all plans

Solns

May be CAM paper

of this form is not possible, i.e. one must somehow narrow down the soln. set "as a whole" (or perhaps not as a whole, but not in factored form).

This ^{subjectly} narrowing of plan in not factored form is often done on a specific (factors) of a larger plan. It can be done if

the goal is recognized as being of a particular type & being amenable to certain (non-factoring) approaches.

SN

On SP: I suspect that normally one does not solve large task nets using probly curves for each ^{known} sub task in the net.

Instead, there is a certain set of net forms that are readily solvable (perhaps by table, analog or MC means) in the form of any actual SP net is "reduced" to one of these forms as an approx. Some methods of "reduction", are to express a set

of all or sequential tasks as a single task, or any single input, output object ~~sub~~ sub net as a single task with suitable probly curve.

Viewed in this way, if a large task net is decomposable into series, parallel & subnets, one might be able to solve it, by obtaining a ~~prob~~ curve (probability of soln, failure v.s. time) for each of several set of tasks - like an Ohm's law soln. of a Resistance network.

Furthermore, if one could do the analog of ~~Δ to Y~~ Δ to Y & Y to Δ x fms, one could perhaps solve all nets!

D2473: Planning

3.18

To obtain a more general approach to planning in solv. finding: Consider a general idea of 2.37-3.18.

Look at various planning methods. Also various soln.

methods. See how they fit into the general data of a plan of 2.37-3.18. I think they ought to all fit.

In particular look at ⁽¹⁾ Michael Miller, Galanter's Pribram - "Plans in the structure of behavior" for ideas

(2) Polya for prob. solving methods to fit into the general schema.

Remember N is S method of Getting Plans: Simplify the problem \rightarrow in some way; solve the resultant problem. Use the soln. as a plan to solve the original problem.

The ~~of~~ planning lang. will be closely linked to the methods one uses to solve problems.

Another closely related idea - to which prob. solving & planning methods are ~~are~~ related - is the form in which info is perceived & stored. Minsky's frames idea is relevant here - for perception, storage, memory, ~~first~~ low order abss.

In a good planning language, it should be possible to work out various further constraints on the soln. - in some v.g. plan langs, it might be possible to work out the complete soln.!

One very general kind of plan is a very non-detailed desc. of a soln. - e.g. "to go from home to airport, I will walk along streets A, B & C". Another kind of plan would give a set of goals - with each

D2573 Planning.

01: 450.40 goal specified either exactly, or with much inexactness, while staying within t. plan lay, one might sharpen up t. goals (Make them more specific), or decide out. order in which they are to be done. One might find that ~~one~~ goal α must be achieved ~~and~~ before one can specify t. next goal w. any exactness.

This characterization of Planning seems v.g. What I can do, is to look at various workings of problems (w. protocols, but perhaps just introspective), & try to characterize them in a common way - of seeing just how t. set of possl. solns., as gen. by t. initial plan, is reduced in each (successful) step of t. soln. of t. problem.

Using ^{is using} Subgoals, 'working backwards', use of diagrams to suggest & discourage subgoals, Thms. (An analog of diagrams (in Geom) is making up lots of randomish examples for logic ~~prob~~ proof problem), can all be discussed in terms of narrowing down t. "Plan" to a more specific (^{narrow}) Set of ^{possl.} solns.

Certain "plans" can be complete - in t. sense of one being able to prove that ~~whether~~ all possl. solns. must be included in it. One then might be able to narrow down t. plan & still have t. resultant plan be "complete".

H's (Genl)

379-382

History of my TM work.

N 0
HIS; N873
Gen

01: 378.40 ^{story} → Now it seems to me, that many times in the past, I was at the pt. attraction of 378.06 - .90. Just why didn't I continue? In fact in ~~1962~~ 1962 "Tug says for Mark Ind" was pretty much the forg. idea. I wasn't as confident ^m understanding of CMT - also I may not have had PMTM - but why didn't I continue? I did take off on ~~SP~~ Stochastic Part, also some shortcut ideas on superfast TM ^{is} more recently (this year) on seq. property Machines. But still, what other side tracks were pursued? }

For an analysis of ~ why I didn't continue on ZTB 143 (Tug seqs etc). See 70TM (K6) 220.01 - 221.03 For more analysis of what I spent my time on 1965 - 1972 See MI folder.

In 1964, PMTM and ^{inductive} GIM were the main track. I had devised some Tug seqs. at that time - ^{overall scheme} I may have devised some earlier as well - perhaps that's what "Dart TM" sequence was.

After browsing a bit thru the work ~~of~~ of a 1965 (before five): Let us compare ~ situations:

- A) At time of ZTB 143 (Early 1962)
 - B) Just before five. (late 1965)
 - C) Now. (late 1973)
-) 3 1/2 yrs
) 8 yrs!

- (1962) A) I had idea of $\approx TM_1 \approx TM_2$ (\approx BSTM = JS)
- 2) Also that Gray search. was very close heuristically to BW search
 - 3) Ideas on a good bug seq. of Arith, alg, etc.
 - 4) Basic ideas of CMI.

- (1965) B) In addition to A) I also had
- 1) A more exact idea on how to xfm ^{but} soln. of BW search into a genl. soln. for Gray search.
 - 2) Some genl. ideas on how to do searches (Stochastic Part \approx SP) ~~is~~ ^{is} Also how to xfm most, if not all AI probs into the form of a Gray or BW search.
 - 4)

01:379,40

1973 C) All of B) available plus via Willis.

- 1) Assurance that CMT is rigorously correct & that my method of doing Grammatical induction w. CFG's is correct. (This is probly impt. for ruff extrapol. of hours).
- 2) Work on ^{Grinstead} ~~ilbourn~~ ^{G-S} method & other ^{summer} shortcuts to TM, involving sequential property (This was done after Willis' paper).
 + SQP m's, ~~of~~ ^{of} this sort, will not work, hvr. ~~is~~ I know just why they wouldn't.
- 3) Understanding of SVT soln. to induction problem - also estimates of ^{expected} error. — The little rigorous work, is no idea of speed of convergence.
- 4) Some ideas on IPC: Digital logic, Memory hierarchy Utilzn: IPC of Human Brain. These are not immediately of value, but will be probly impt. if a super TM is to be constructed that will take optimum advantage of hardware available.

$$\frac{((1+x/3)^{x/2} + 1) \times x}{13}$$

Actually, from Mid 1965 (~ \$122) untill ~~beginning~~ mid 1969 (~ 4 yrs) I did relatively little new work on TM.
 Mid 1965 to f. end of 1965 was a slowing down.
 Jan 1969 to Mid 1969 was a re beginning.

Actually Oct 13, 68, I started reading W's first report.] serious work - Willis
 On Jan 4, 69 I had perhaps finished reading much of it.
 On July 11, 69, hvr. I had only done 30 pp on TM paper.
 What happened from Jan 4 to July 11? [Cornell talk Feb 69, constructed Cooper Aug back Mar 69, Ap 25 - Jun 11

Gend.

440-442

D1973 Previous Gen Ref: 407.40 (Also R 439.25-40)
Gen! On t. SQPM idea!

R 439.40: See Blous 374.10 - 40 for most recent summary of objections
129.02 for previous " " "

TM440

Anyway, when looked at in t. forgo. way, in terms of these "Batches", this SQPM method is pretty much like just mutating old perms to find good trial new ones.

As such, it lends itself well to ~~the~~ parallel processing, i. y. discuss. of ~~the~~ procc. efficy, use of heav. memys, IPC, etc.

There appears to be a big diffrnce, hwr, betw. t. present concept, i. t. old SQPMs: In t. present thing, t. machine looks at t. Macrosymbol

(= batch of corpus), before trying to devise codes for it. In t. old SQPM's t. trial code was devised ^{for t. corpus} before t. Machn could look at t. corpus, i. then t. ext code was rejected if t. output didn't conform to t. corpus.

Hwr, in t. "new" Machn, it seems that there is something that has t. old sq. propy. We have a seq. of corpus coding Methods (= sort of perms). These methods are obscurely mutated, re combined, etc, in a "randomish" way to form ne t. next trial coding method.

So our "object to be extrapolated" is this batch of perms (= coding methods) i. parts of perms. As t. corpus grows, this object also grows. Also t. wts. assigned to various parts of it changes as it i. y. corpus grows.

Re! Back tracking: As t. set of perms i. parts of perms, i. continuation rules i. their wts (= sort of naming perm) changes, TM, tries to re code t. past corpus in terms of these t. modified bunch of predn. objects. If a significantly new code is found, t. wts. i. objects in t. "Prediction body" are reforming of t. Predn. Body suitably modified. This can take a lot of cost in a machine, i. it certainly can do t. same in a human!

My mind isn't clear on just what t. pt. of this "New View" is: Perhaps that new trial perms (corrupt to new trial codes) are found by mutations etc. of old codes i. sub codes — but that this is done in a sequential manner — so that we have a one pt.

441.91

pl: 440.90 a term of t. "prdu. body," that works up to that pt. We normally modify this prdu. body (mutation, recombinn. of parts) by adding symbols to t. "prdu. body" desc. I'm not sure if we devise new definitions in this normal "sequencial coding". We may, hvr, do Backtracking — which is to apply newly devised successful (or almost successful) abs to t. old corpus to recode part of it.

One of t. problems in any ^{random} mutation scheme w/o ~~recombinn.~~ recombinn. scheme, is to keep it sequencial, so one doesn't have to ~~try~~ try recoding t. entire old corpus, whenever a new mutation ^{is tried}.

Perhaps it is poss. to show that ^{usual} backtracking, results in a universal / ^{independent} device — i.e. one that can devr any CPM.

Note that t. "prdu. body" plus some error rate, constitutes a complete desc. of t. corpus, so, one should be able to apply ideas of CMT to it wrt t. problem of whether or not it's a univ. ind. device.

Well, my impression is that in order to be univ., t. system ~~will~~ need scan over only a finite (the very large) set of mutations & recombinations — but somehow, as time went on, t. CB would have to A. This corresponds to trying "only" (a finite num) 2^L poss. codes for t. corpus, but t. C.B. to ∞ .

From a practical standpt., hvr t. no. of tryable codes is $\ll 2^L$, so one will have to ~~try~~ try $\ll 2^L$ & also let CB be $\ll \infty$. 2^L is sort of like ∞ from a practical pt. of view — so one will take some path in CB, code space! i.e. CB will be some funct. of t. no. of codes tried — also a funct. of t. part of t. code being tried.

For most, perhaps ^{almost} ~~all~~ all, funcs we mite consider, CB is ~~finite~~ of approximately known size — definitely bounded.

Some actually usable functions (like t. inverses of some easily defined funts), have short descs but take an enormous time to compute.

[e.g. $F(x) =$ the middle n bits of x^2 — where x has n bits. F^{-1} is lots of trouble to compute. — also I don't know if it's single valued — (the one could specify t. smallest value. — also we don't know if it exists for all values of its argt.!)]

01.441.40: ~~TM442~~ Hvr., it would be well to look at one's mutation, recomb. ~~TM~~, backtrack scheme, to estimate to what extent it's "univl." Perhaps a system might start out not univl., but after it had a certain eng. seq., it would have certain new mutations & obs. \rightarrow so that then it would be "univl."

80 \rightarrow Actually, I think that in the Genl. approach I expect to restart soon:

[ZTB 141, (376.18 - 378.40 Genl), 387-383 SFTM (Also see (391.28 BM)]

Universality is Not a problem: I will have this ^{Eng.} seq. of problems: If TM is unable (or even unlikely) to ~~backtrack~~ solve any problem, I will know it, & add something or modify the system so it will be able to do this. Since universality is not a difficult thing to get, it is likely that I'll have it ~~early~~ early in the Eng. seq.

As for this business of Backtracking: It would seem to be desirable for optimum prodn. w. a gn. SSZ (or corpus). Hvr., strictly speaking, it is not ^{absolutely} ~~necessary~~ ^{for universality} ~~providing~~ the TM remains univl. I.e. at any stage of TM's "tracking the corpus" - say the rule of the underlying Fdss is suddenly completely changed.

TM should be able to derive this rule, no matter what this new rule is. Hvr., the CB may be very large if the new rule is very distant from the old one - or more specifically, the new rule is very distant from the most likely "near pairs that would be tried by the ~~TM~~ ^{it prediction} body" that was attained in following the old Eng. seq.

I think the idea is to keep the "predubody" good in the sense of being likely to track near by (to the previous pair tracked) pairs, but still not so narrowly a priori - so that more distant pairs have reasonable chance w.o. excessive CB.

PLAN : T.S. 248-406

VHM =
Van Heerden Method 248
= Shortest rod.
From itself

On the Q. of what to work on re. SUMM:

1) Should I spend more time on the SVM ~~re~~ SVM itself
i.e. the 2 post. forms of SVM (max profit of cylinder code for PEM. v.s. Max Σ profit of PEM)?

2) Working on Max ~~Method~~ v.s. Direct v.s. Standard methods.

Whether or not SUMM is "true", 2 is of much value:

- a) In ironing out just how to apply SVM - how to evaluate profits of various types of PEM (The much of this, I have done before - I should try to find it a index it)
- b) The basic idea of 2 is to see if it is a counterexample of SUMM - i.e. that the direct method may give a better score than the Max Method.

There is some Q about the relevance of my comparison of the Max & direct methods: Should the comparison be for large S ($\approx SS \geq$) or for $S \approx 2N$? - or should the N value in the Max method be automatically chosen for each prodn. so as to Max the post score? - so that in the Max method, even N need not be part of the decision of the PEM.

On second thought, however, I know that in a certain specifiable sense, Max method ~~is~~ will always give better results than the direct method - so under all conditions (any values of N & S) Max ought to give a better score.

So perhaps consider $S = \frac{N^2}{2}$. I think in cases where S is just a bit larger than N , Max method may win, because direct has those extra N params to specify.

01: 248.40: It may be that I can get some good comparisons
 betw. Direct, standard & Max methods w.o. a lot of
work. One of the main problems is punctuation - t.
 goal. problem of coding steps of nos., etc.

10 Re: T. most immediate goal of the "Review Paper" on
 the necessity & sufficiency of CMI: In order to get others
 to work on this, I'll have to do more than just present
 these results. \ddagger One thing I can do, is to show how
 CMI can be used to solve various problems: e.g.
 Grammatical induction, - how to use it for Abduction (search of
 new sci hypotheses to try), show how it unifies a/o
 quantifies a lot of other work on Pattern discovery, induction,
 etc.

Sep 18, 73: T. work of 148.01-40, 250.01-99 has been
 more or less sloppily completed. No proofs, but strong conjectures.
 300.01 - 302.40 is a review of the results. 250.00 remains to be done.
 I will now take a bit of time off to Reference & paper on Grammatical
 Induction.

Tag. seqs.

There are some recent (Busby) refs. Also, I would do well to try to make an index of older references. There are various direct references to tag. seqs. in the set of Indexes in MF₂₀₀. Also, see notes on my indexing of the Burnt papers after t. '65 fire.

on Tag. Seqs.

Re: The ~~Tag~~ Tag Seq. of 2143 - i.e. arithmetic learning.

For the "with learning" part, the problem was ~~learning~~ learning x fn from Algebraic to Polish notation. ~~It~~ As we usually think of it, t. Alg notation involve partial ~~ordering~~ ordering of rules wrt precedence.

If the symbols used in Alg in Polish notation, are mainly t. same, e.g. $+$ \leftrightarrow $+$; \times \leftrightarrow \times , etc., t. problem is much simpler than if other notations are used.

I.e. say t. input is $4+3=7$; i. we want t. machine to be able to write a pgm. \rightarrow if its input is $4+3=$, t. output will be 7. The permissible

~~symbols~~ symbols in t. pgm. are 0, 1, 2, ... all integers, +, -, x, /, (,), =. | +, -, x, / are Polish operators so $+(3, 4)$ means "add 3 to 4". So we want " $3+4=$ " to be translated as " $+(3, 4)$ ".

~~We~~ we are given the 2,0. pair $3+4=, 7$

One way to solve this general type of problem, is to "work backwards". I.e. we have many examples, = occurs in all of them is not relevant to this set. If t. examples are

all of t. form $A(\pm)B=C$, then we factor them into $\alpha = \{A_i + B_i = C_i\}$; $\beta = \{D_j = E_j = F_j\}$

We then try to solve each of these 2 sets individually, by trying to combine A_i & B_i into machine instructions, to get C_i . A natural thing to try is

$\begin{pmatrix} + \\ - \\ \times \\ \div \end{pmatrix} \begin{pmatrix} A, B \\ B, A \end{pmatrix}$ - only 8 trials needed to get ~~them~~ both right answers. - An average of 4 trials

Tug says

01:380.40:

If it is suspected that the "+" in x input is the "+" in the machine (arg. might be similar, we try

+ (A, B) then + (B, A)
 Same order Reverse order.

So we get to write answers written off

Other impl. ideas: (1) the idea of "an object" i.e. if (3+4) occurs in an expression, we can rewrite the expression as "7" i.e. the expression will be equiv.

(2) Idea that expressions w. fewer symbols are easier to evaluate. — So whenever we can do (1), this reduces the no. of symbols in an expression is usually desirable.

Idea (1) & (2) should make the idea of "precedence"

in alg expressions, come easier.

So, in this Tug sep, TM, system, I have a sequence of "hours" along with a seq. of "solutions". Hours are for TM₂, "solutions" are for TM₁. At any particular pt. in time, a "solu" is an (I, 0) operator.

So here we have a very nice example of a TM₁, TM₂ combination, but TM₁ cannot = TM₂. I am TM₃, supplying hours for TM₂.

Eventually, TM₁ can begin getting into probabilistic problems, problems of improving programs, problems of devising hours. So we can then try TM₁ = TM₂.

Perhaps at a not very advanced pt. (say when TM₁ is able to extrapolate using CF Grammars) we can present TM₁ w. TM₂'s entire seq. of hours, & ask TM₁ to suggest good trial hours from this seq.

Here, for serious hour guessing, TM₁ should be at the pt. where it should be able to "reason out" why certain hours are good, & construct new trial hours based on principles

01: 250.40: See 341.28(BM) for Genl. plan: A More detailed descr. of what t. plan of work on TM prepar: Gen 376.18 - 378.40 is fairly Good. Other discussions of General Modes of Constructing TM's should be locked into → See INEX 281.28 RT for some refs. Also stuff on "Tug seqs." is good. Gen ~~375.01~~ 375.01 - 376.17 is also good.

As for t. Mechanization: see SFTM 387.01 - 388.40. This is pretty much t. old (Early 1973) idea of having a large, slowish common mem, used by many arith units in ll, that have each their own cache memy & cannot write into t. common memy.

This latter assumes that t. main problem is more or less search w/o any interaction betw. trials. To t. extent that this is true, .08 - .10 is a v. fine thod. If, out. of how hard, trials are very much determined by t. results of previous trials & /o compns., then it is not clear that that system is at all good.

On t. other hand, for essentially sequential problems, serial hardware that is very fast, is much better than lots of slow parallel hardware — So modern computers could be far ahead of t. hum. brain in this respect!

Marvin suggests that in most present AI sys, very few trials are used (say < 10), & it may well be that I work is essentially sequential.

Anyway If one has a lot of TM-type problems to do, at t. same time, it is probably poss. to multiplex various sub-operators ~~≡~~ betw. t. ~~various~~ various problems — so All in all, it is t. "IPC" needed by each problem, that is important.

My present impression is that I should ~~sp~~ first spend a lot of time on t. Tug seqs. T. seq. & t. soln. should be written out in English in as much detail as poss..

Note that it may be necy. to introduce "side" tug seqs. so that certain concepts, heurs, ect., are available to t. rite times,

I now, presumably, have available to me, more tools for doing tug seqs. than I had in 1962. Hrr, it would be well for me to have some of t. problems in solving this tug seq.

in Mind as I review t. 1962-1973 yrs. — i.e. I will find suggestions for solng.

→ Perhaps Parallel w. t. direct work on T.M., write "Study Paper" or "Position Paper", etc. for CIAP. Consider Meeting, Informal "colledge", visiting people, ~~etc~~ writing lots of letters, etc.

Reviewing this old material is very stimulating. Some recent ideas on what the "Main line" should be!

378.40 gan
378.40 gan

Agree 376.18 - 378.40 - but some new ideas!
Working on the problem, it is not necessary to start at the beginning of a task.

Instead, work on probs. that are interesting, but I have ideas about. Later I can try to fit the probs I've worked on into a top. seq. & make top. seq. "beats" be to. team.

Some problems: (1) My narrowest view of "interlocking" is that it is properly presented it will take few trials. (2) That the machine should use concepts that a human could use. - e.g. the idea of a "quantity" in which traveling? t. soln. of equations.

(3) In solving equations or systems of eqns. There are many ways to solve eqns. We may want TM to be able to do general of them. T. decay of each in values several imp. concepts that would make that decay likely.

e.g. One way of solving systems of eqns involves successive approx. for a computer, this can be very fast, because computation of pts. is very fast. Some ideas needed (a) T. idea of "elasticity" in numbers, so TM can know that an approx. soln. is almost adequate. (b) T. idea of continuity - that small changes in params will produce small changes in the system.

(c) Some topological ideas - that it is ~~spurious~~ changes in params will produce small changes in the system. Usually one or many paths of values of ξ that pass thro all pts. from a to b.

Using ideas (a), (b), (c) it is possible to devise successive approx. methods, by trying to solve within a manifold a manifold of x. manifold with idea of Taylor expansion, w. ~~of x. manifold~~ that order costs probably & one into be able to dev. ~~of x. manifold~~ Newton's method for more clever methods

for economy on cost & machine should be as sequences as possible. e.g. Sequential heuristic means that only accuracy as we needed, need be processed. Also "Look Ahead" in units is wasteful of cost.

SN

01: 401.40:

Also, it might be possl. to use work of other A.I. workers.

Any work on learning ppgms might be used - e.g. Sussman's & Goldstein's ppgms.

Or, perhaps more commonly - Given S's or G's ppgm as a working device that can solve certain kinds of problems - what is an adequate set of abstractions that ~~is~~ a machine would have to have, & what kind of type seq. would it have to have to ^{to perform} learn those particular tasks in those ways?

To reason I might want to choose hard ppgm. tasks & t. ppgms that do them as "things to learn" is that these kinds of tasks & t. way they are solved, were meant to be patterned after t. way a person might do them or try to do them.

It would seem like a good idea to get some experience w. simple problems (like with learning), before getting involved w. t. more complex kinds. Hurv, in t. case of with learning, I'm not sure I can define t. problem in a way that will make its soln. contain abs. ~~is~~ that will be useful in more diff't probs. T. problem of a child learning with seems very complex.

Hurv, I could define t. ^{with learning} problem as I have in TS 380-391. This might be enuf to get me to more diff't probs. - if not, then add to t. ~~is~~ the seq. or change to a diff't initial problem.

Anyway, if I do have input in alg. notation, it would be neccy for t. machn to learn precedence rules; a/o that certain things ^{always or almost always} can be treated as objects - 7. meaning of parens, etc.

After t. machine learns alg. notation, try teaching it polish - both forward & backward, then perhaps a part of APL, then

Some ordinary machine ppgms. (one that expresses t. same thing)

One approach would be to start w. Sussman's & Goldstein's thesis ideas. ~~is~~ Another would be Marvin's Frame Theory, which is meant to be a model of child learning up to a certain pt. of sophistication. In t. case of S & G's work, I'd have to outline t. main hour pts., (406.01

N. 28-73 PLAN T.S.

01: 402.40 then try to devise the seqs & leading from ~~to~~ a simpler set of hours to those pts.

In the case of Marv's work, perhaps the outline of the entire the seq. is more or less implied — along w. the seq. of hours used. — but I haven't looked at the "Frames" ~~the~~ report yet.

One idea from "Frames": That in choosing a problem field, we must be careful to try to ~~us~~ make it a manageable small universe, — so that we can get essentially all necessary hours from that small universe.

Anyway, in choosing problems, at first try to stay away from

① Human Logic & Grammar ② Vision: Because probably these have very many special-purpose tricks that don't generalize well to other fields. Also, one should be able to find very seriously difficult creative discovery (say in math ~~the~~ & science) in which these 2 things are not impt. (altho vision obs. may be impt. in physics, say).

P. W.

Casey Wtng. Problem 384-399

14,73

(PW) =

Pawn Wtng Problem: see 397.01 for REV (JPM) This is essentially a continuation of the old SVH - since it got mainly into this problem.

Previous wts 350, 349, 345, 329, 325, 27, 324
317.10 318-320 on CPM's

The idea is: Say I have 3 pems, $P_{1,2,3}$. I have their costs or k cost or profits of deriv. I also have t. profits of t. corpus w/ each. Using t. full CMI then I compute t. best (linear or other wts to use for these 3 pems)

SN: Even tho there is overlap of component Derivs, t. all pems method may be correct. It may have to be renormalized, hvr. T. normz factor may $\rightarrow 0$ if so, will have to be inserted w. care. The idea is that "t. all pems method" is probably "unbiased". It has extra codestory corpus, while other unbiased "all pems methods" use less than all possl. code - i.e. they leave out many, but in an unbiased way. I think t. All pems method is both unbiased & has a normz factor of 1.

Essl. Criteria for a good combination rule:

- 1) t. resultant pem has max expected value of future profit of t. corpus (if this is zero for all future combinations considered, take some sort of Geom mean profit/symbol)
- 2) ~~best~~ - smallest expected mean sq. error in future or smallest total sq. error in future (if that) is zero for many possl. combins)

Consider t. next 100 symbols of t. corpus. There are 2^{100} possl. seqs.

Via CME, t. i th one has probty q_i ($\sum q_i = 1/2^{100}$).

each of t. 3 pems has profit $r_{j,i}$ ($j=1,2,3$) for this corpus.

So of t. 3 corpi, t. one w. best expected future corpus profit is that for which

$$\sum_i r_{j,i} q_i \text{ is max.}$$

If we give P_j wts a_j then what we want to select

$$a_j \text{'s} \Rightarrow \sum_{i,j} a_j r_{j,i} q_i \text{ is max, subject to t. constraint that } \sum_j a_j = 1$$

01:384.40:

Consider 384.37. Let's expand q_i in the form
 given by some version of the "all pairs method":

Let $R_{l,i}$ be the i -th part of the corpus with the l -th pair.

" α_l " " " wt. of the l -th pair in this particular "all pairs method".

So $q_i = \sum_l \alpha_l R_{l,i}$

384.37 gives ~~$\sum_l \alpha_l R_{l,i}$~~ ~~$R_{l,i}$~~

$\sum_{l,i} \alpha_l R_{l,i} r_{j,i}$

384.34 gives

$\sum_{l,j,i} \alpha_l R_{l,i} z_j r_{j,i}$

If, in .13, show that ~~the 3 pairs considered are included~~

~~those considered in the "all pairs method",~~

we had not 3, but M pairs (where $l = 1/M$ as well as $j = 1/M$)

o those M pairs were the same as those considered in the "all pairs method" — then I guess $\alpha_l = \alpha_l$ for all l —

but can I prove this? we then have .13 \rightarrow

$\sum_{l,j} \alpha_l \left(\sum_i R_{l,i} R_{j,i} \right) z_j$
 symmetric

the symm matrix has positive real eivalues — so rotate the system
 until the matrix is diagonal. we get $\vec{z} = \vec{S} \cdot \vec{z}$

The solu. for \vec{z} is normalized $\vec{z} \cdot \vec{S}$, not \vec{z} , unless \vec{S}
 has all unity eivalues — which maybe it does!

However, even if constrained .20 is not true, so we still have only
 3/pairs to work with, $\vec{z} = (\text{norm}) \vec{z} \cdot \vec{S}$ is the solution.

False! it would be true if we had the constraint $|\vec{z}|^2 = 1$.

Let's go back to 384.37: say

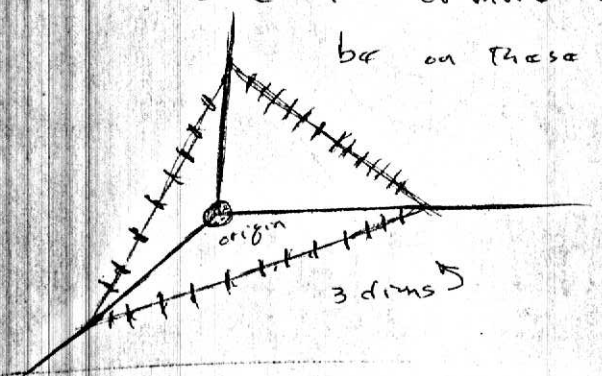
$\sum_i r_{j,i} q_i = b_j$

01:38540 : Now, in general, to solve $\vec{a} \cdot \vec{b} = \max$, w. t. constraint $\vec{a} \cdot \vec{I} = 1$ is a ~~problem~~ with the further condition $a_i \geq 0$ for all i .

In 2 or 3 dimensions, a solution is always $a_j = 1$ if $b_j = \max$ of all b_i ; otherwise $a_j = 0$.

05
06
Another solu.: If n of the b_j 's are all equally max, then the corresp. a_j 's can be any set of value \Rightarrow Their sum is 1.

In ≥ 3 dims, $\vec{a} \cdot \vec{I} = 1$ means \vec{a} has its nose in a certain plane. On the other hand, t. locus of all vectors $\vec{a} \Rightarrow \vec{a} \cdot \vec{b} = k$ is a plane \perp \vec{b} . If we consider all planes \perp \vec{b} , they will all intersect t. plane \perp \vec{I} in st. lines. In general, the plane corresp. to a plane of max k , will pass thru one or more corners, so t. intersection must be on these ~~lines~~ lines.



This result can be obtained directly ~~from sense~~ in n dims: Consider $\vec{a} \Rightarrow \vec{a} \cdot \vec{b}$. We know t. b_i 's. Clearly, if b_2 is t. largest b_i , then $a_2 = 1$ for $i=1$; 0 otherwise. If several b_i have max value, a linear comb. is ok. (1.05-06)

This is an interesting result! It says that one can always choose just one pen, ~~ie~~ that is at least as good as any combination of pens.

Note that this does ~~not~~ ^{almost} contradict t. results I got w. linear regressn., because t. max method was almost one of the $\&$ linear pens I considered. It didn't have a gaussn. error distribn., but it would if there was either a lot of data or a large no. of coils (I don't remember which).

Well, perhaps Max is "Best" if averaged over the future rather than t. past corpus!

In 384.37 note that r_{ji} and q_i both include both the known corpus & t. unknown future 100 symbols.

Wall, not so surprising - that one poem should be "Best"; But not that this is "Best" wrt t. "correctly known" ~~future~~ probly distribn. of t. future Corpus.

→ Suppose that all one knows is \approx krusts $\&$ of pems $\&$ prosts of a sub corpus wrt. these pems.

One could consider ~~t.~~ t. average (over all future corpi of length 100) of all pems w. krust brst of subcorpus $\approx \alpha_1$.

Compare it w. similar average of α_2 $\&$ α_3 . Perhaps It would seem that probly t. one w. min α_i would be best, but t. cross correlns. betw. γ . 3 pems are not "consulted".

Perhaps a better way to look at it: we regard t. 3 pems as each 3 difent estimators of γ . future. $\&$ If they were "indip" - it might be easy to estimate their relative wts. for a better estimate than any of them. Hvr, if they are not "indip" γ . prob. is much harder.

Perhaps we can discover γ . amt. of lack of indep. by Sampling probly values from γ . 3 in γ . past corpus.

Say use random samples or use entire past corpus.

O.k.! Say each poem P_i ~~is~~ is a sequence of values ($i=1/3$)

q_j^i ~~(j . symbol of P_i)~~ $j = \gamma$. order of symbol in t. known corpus. ($j=1/50$) say

Say q_j^i is t. true probly. (which we don't know).

So we have t. 3 vectors \vec{q}^i . We might be able to estimate their relative wts. for optimum estimate of γ . unknown probly of γ .

1st symbol ~~is~~ by knowing only γ . cross correlns. betw. γ . 3 vectors, \vec{q}^i .

A counter example: Say P_1 gave $q_j^1 = 1$ for all j ,
 P_2 " $q_j^2 = .5 \dots$
 P_3 " $q_j^3 = .2 \dots$

It would seem (unless P_1 has by krust), that P_1 is v.g. $\&$ that P_3 is very poor. - I don't see how "cross correln/ studies" would give P_1 any more wt. than, say, P_3 .

One approach: Assume that t. corpus is a deformation seq. - so probys are all 0 or 1. In this way t. (error)² for each component of \vec{q}^i will be $(1 - q_j^i)^2$.

01:389.40: In this way, we can give each pair wt \propto its (variance)⁻¹.

It troubles w. this!

1) We haven't taken cost of pairs into account. We may be able to do this by computing expected future variance of each pair.

2) Cross correlations betw. t. pairs were not considered - but this can be done.

Getting e. into coils is (as I'm modeling it now) ~~is~~ identical to t. linear regressn problem. T. linear regn. soln. is ~~not~~ exactly correct for a n-dim. Gaussn. distribn, i/suitable a priori of coeffs. (uniform) - we get t. Max method soln.

The correlation matrix of interest is t. 3x3 matrix

$$M_{km} = \sum_{j=1}^3 q_j^k q_j^m$$

The constant vector is $k_k = \sum_{j=1}^3 q_j^k \cdot 1 = \sum_j q_j^k$

t. "y" value is always q_j^k

T. wts. are then $\vec{X} = M^{-1} \cdot \vec{k}$

(SN) A "criticism" : we could just as well have used $q_j^i - 1$ instead of q_j^i ($i=0|3$), but we'd have n eqs. - but

then, $q_j^0 - 1 = 0$ always so $\vec{k} = 0$ & we get nothing!

Whoops! But then we would have to allow a constant to be one of t. pairs - so it may be ok. (this is equiv. to a constant term in regressn)

While 08-18 takes cross correln. into account, cost of pairs

is not considered. Perhaps we should do this by modifying t.

\vec{k} coeffs. In my proof that Willis' information/error convergence

implies ~~it~~ \leq sq. error convergence, I got a relationship

of proportionality betw. $f \cdot z$. I can use SVH ideas

to estimate ~~variance~~ \downarrow of mean (-ln probab) error. - how can

I relate this to t. of pairs ~~or q_j^k~~ $1 - q_j^i$ error?

In general, we may always want to introduce a constant term

- but if we use q_j^i as t. variables, t. constant term will correlate

w. $q_j^0 = 100\%$! thus destroying t. whole thing!

(SN) Another criticism of linear regressn. is that t. sum of t. coils.

~~$\sum_{i=1}^n x_i = \sum_{i=1}^n x_i$~~


(or $\sum_{i=1}^n x_i$) $\cdot 18$ for def of \vec{x} may not ≥ 1

N-16-73 PW

390.90

Another Guess at a Soln:

Just do $390.08 - .18$, ~~is~~ obtain

\vec{x} . Then mult f . ~~weights~~ coeffs of \vec{x} by t . proper costs of ~~the~~ den of corresp pairs! 

In the case of linear regressn., we chose t . wts. so that y . resultant pair has least ms error in t . known corpus. A corresp. \rightarrow

SN $\left\{ \begin{array}{l} T. \text{ coeffs may occasionally be } < 0 - \text{ i give probly values that are } \\ < 0 \end{array} \right.$

\rightarrow method would ~~not~~ choose coeffs so that t . prob of t . known corpus wrt t . known Σ corpus, would be max. - Hvr. again, we did not consider cost of den. - so a pair that got $q_i = 1$ for all j ~~that~~ would get $w = 1$, even tho it ~~cost~~ was completely a.h. (i.e. large t . cost). $\rightarrow 392.20$

In the linear regression approach to PW: what I would like

is that in ~~the~~ 390.18 t . coeffs of \vec{x} is ~~the~~ M are obtained w. t . infinite future corpus. The ~~is~~ ad-hockness due to non-zero den costs perhaps makes t . coeffs of M as derived from t . known ^{past} corpus, about rite - i.e. no reason for systematic error. But t . coeffs of \vec{x} are always excessively large due to ad-hockness.

No!!
 \rightarrow see p. 38

One way to estimate how ^{much} \vec{x} components should \downarrow !

Say S is t . length of t . known corpus, i ~~is~~ AWARD

~~is~~ $-\log_2$ prob of pair den. Then

$$\rightarrow X_i \rightarrow X_i \times 2^{-\frac{A_i}{S}}$$

$\rightarrow 2^{-\frac{A_i}{S}}$ is about how much each of P_i 's probly estimate is "hy" due to ad-hockness.

\rightarrow for all S symbols of t . corpus

$$P_i \rightarrow \left(2^{-\frac{A_i}{S}} \right)^S = 2^{-A_i} \text{ - or a}$$

cost of A_i

T . ~~components~~ of M can be similarly corrected, but $M_{2m} \rightarrow M_{2m} 2^{-\frac{A_i - A_m}{S}}$

01: 391.40: $z^{-\frac{A_i}{S}} \approx 1 - \frac{A_i}{S} \ln z$ if $\frac{A_i}{S} \ll 1$,

similar ~~approx~~ approxns for $z^{-\frac{A_i}{S}} - \frac{A_i}{S} \approx 1 - \frac{A_i + A_m}{S} \ln z$

There may be a simple approx. effect on $\frac{1}{M} \sum X$ that is gn. by these corrections.

06 \rightarrow Another impt idea: In my recent work on SVH, I got several ~~corr~~ different corrections to obtain expected future profit/symbol. These ~~are~~ more exact corrections should be used instead of $z^{-\frac{A_i}{S}}$. The difference is, I believe, a sizeable factor, under most cond. \rightarrow see 399.18 for some Q's

09 Also note: T. cross correlns in \vec{M} 's components can be obtained by sampling - i.e. it is not necy to do all of t. cross products.

20: 391.15: It is poss. to write $\prod_{j=1}^{50} (z_i q_j^i) = \text{profit}$

then try to choose t. z_i 's \Rightarrow this prod is max.

\rightarrow one can "correct" t. q_j^i 's using t. method set 391.30-40 (or more exactly 392.06-10).

There might be a good way to do t. expressn. of 20 - but at first glance it looks like a real diff. thing. If we only have 3 or 4 pairs, then we could do a lot of computation w. different values of \vec{z} - a H.C. problem. If t. hill is sufficiently smooth a non-bumpy a soln. may not be hard to get - one can use various tricks to speed up t. work.

32 \rightarrow One rapid way to evaluate t. product of 20: add $\ln q_j^i$, but choose i 's randomly w. probly z_i . This method can perhaps find a ruff max, \Rightarrow then more exact evalns can be made in t. "good" region - wrong! (see 393.23)

Another way: MC when always choose j at random ($j=1/50$) - uniform distribn. If one has tables of $\ln q_j^i$, one can really do this very fast. (fit pt. add may take too much \rightarrow time, so use fixed pt values of $\ln q_j^i$. So we just \rightarrow 393.01

.01:392.40: Keep choosing these random values and adding them on, using an exponential update (x window) & slowly change the $\frac{1}{2}$ values.

→ There ~~is~~ is a basic Q. of whether the method of 392.32-33 actually gets the correct expected values of \sum products.

→ While the forgo method does give a lot of noise, remember that even if one does it "exactly", there is still a lot of noise due to finite SSZ - so one can perhaps choose the "x window" τ value ~~is~~ is the scanning parameters so one gets a maximum expected rate of hill climbing per cost of computation

Note: To do the forgo analysis, one need only have a set of sample pts for P_1, P_2 & P_3 - they need not, however, have the same j values for the same i values. Hvr, the mode of selection of samples must not have a systematic bias.

Actually, the forgo ~~MC~~ method can be applied to any type of Goro - not only max profit, but best ms error or any other

→ Well: the method of 392.32-33 is wrong!

We want $E(\ln(a_1 q_j^1 + a_2 q_j^2 + a_3 q_j^3))$ - that's what we would get
 US $E(a_1 \ln q_j^1 + a_2 \ln q_j^2 + a_3 \ln q_j^3)$.

To get .24; $\ln(x) = (x-1) - \frac{1}{2}(x-1)^2 + \frac{1}{3}(x-1)^3 \dots$

.24 $\approx E(\sum_i a_i q_i^2 - 1 - \frac{1}{2}(\sum_i a_i q_i^2 - 1)^2)$ ~~is wrong~~

The three partial derivatives wrt a_i 's give 3 linear eqns in 3unks.

The coeffs ~~are~~ are of the a_i 's are probly $E_i(q_i^2 q_j^m)$, or some similar thing - woops!

in .24, $\frac{\partial}{\partial a_1} = E(\frac{q_j^1}{a_1 q_j^1}) = 0$

Hvr, since the q_j^i & a_i are always ≥ 0 , I don't see how this is possl.

As a result, it must be that we are at an edge - i.e.

either a_1, a_2 or a_3 must = 0 or 1. Perhaps we can use this

opt. to show that all but 1 of the a_i 's must be 0!

Woops! I forgot the constraint $\sum a_i = 1$

→ 394.25
394.01

N.17.73 PW

01; 393.40: The expansion of 393.24-.32 can be just carried to 2nd degree - in which case t. equs are solvable, if not too many pams are involved. If we do it to > 2nd degree, t. equs can be solved by h.c. methods. We need t. high order cross correlus of t. various pams (e.g. $E(q_j^i (q_j^i)^t)$ etc. These can, as before, be obtained by sampling.

10 Since we also solve 3/equs in 3 unks in t. "linear regressn" method of 391.01 ff, it may well turn out (that t. soln. of it as well as t. quadratic case of 393.24-.32 are identical! This would be easy to check.

18 The corrections of t. forgg. soln for excess of pams ($\equiv A$) > 0, can be done ~~before~~ to make corrections to t. correlus & cross correlus. as in 391.30-.38.

Alternatively, it may be possl. to insert these corrections after one has t. solns. for $\vec{A} = 0$. In such a case, it will be a lot clearer as to just how t. \vec{A}_j 's effect t. results. Hvr, eq. 26, if true, suggests that this ^{latter} may not be very easy.

25 : 393.40! ~~perhaps~~ t. 399.33 equs and up with $E_j \left(\frac{q_j^i}{q_j} \right) = \lambda \cdot \text{const for all } i$

26 \rightarrow $\frac{\partial}{\partial \lambda} \left(\ln q_j + \lambda (2\lambda - 1) \right) = 0$
 $\frac{\partial}{\partial \lambda} \left(\sum_{i=1}^3 2\lambda q_j^i \right) = 0$
 $q_j \equiv \sum_{i=1}^3 2\lambda q_j^i \rightarrow 396.25$
YES!

\rightarrow In t. SUH discn. of A_{ij} (t. excess of deriv of P_{ij}), we get an upper bnd on an estimate of t. cost/symbol of t. future corpus. We can, then, expect a pair of 1-sided error in such estimates. An impt Q. is - what is t. effect on such discs as 391.30-.38 is in general, on t. optimality of t. recent results using corrections like 391.30-.38? An analysis like 18-24 mite make this a lot clearer! - but eq. 26 (if true), suggs that this mite not be possl.

01:394.40: Well, O.K. say we do have means for combining Pems linearly in
 any other way, so that a Gorc ($M.S. \text{ error}$ cost) is extremized.
 If we do take the cost of Pems into account via $\approx 391.30 - 33$
 of what significance is this? — Clearly we will end up w. $\approx 2n$ ~~approx~~ approx
 $\text{cost/symbol} \leq$ any of those expected from any of the component Pems.
 Will this "combination Pem" really on the average, be any better than
 the components?

Remember that since, if we have n Pems combined, we have
 just adding in an ad-hockness assoc w. that n! If the combin
 rule doesn't gain us more than we lose via that n — then we will
 do better by either using the "best" Pem or \approx a combination
 of $< n$ Pems.

20 I think the exact value of the ad-hockness added by this
 combining of n Pems can be obtained via the method that I
 used in analysing linear regression (q.v.). To do this, it would
 be necessary to know the $\frac{\partial^2}{\partial a_i^2}$ (cost Gorc) — or correspondingly for
 any other Gorc. (Again, note that w. the Max Method, I didn't
 use this simple Gorc, yet the calculus was not difficult.)

In general, this $\frac{\partial^2}{\partial a_i^2}$ can be obtained by obtaining $\approx 2n+1$
 empirical pts. — n being the no. of dims. of \vec{a}

Anyway, it is clear that while I have described but
 2 combination methods for Pems (actually, 1 method (linear) w.
 2 different Gorcs), there are an ∞ of such combin. methods,
 the Max expected cost being the "ultimate" Gorc. Which are "best" or even
 "reasonable," I don't know.

An important thing in the Pem combination is getting
 the "A" value of the combination rule. Perhaps we could
 devise a combination rule like the Max Method, so that A value would
 be zero. Note that the "A" value is not "zero because
 it is optimum" — but zero because the rule is fixed before 396.01

01: 396.40 Review of 384.01 - 396.40

T. Basic idea of P-weighting is: Given a bunch of Pems, how should we best combine them into a new, pem, at least as good as any of the components?

One method is to use linear wtng of t . n component pem probys. The wts, a_i are chosen $\Rightarrow \sum_{i=1}^n a_i = 1$ i. t . probst of t . known corpus is max.

If q_j^i is t . probty of t . j^{th} corpus symbol wrt t . i^{th} pem, we chose $\vec{a} \Rightarrow \prod_{j=1}^S (a_i q_j^i) = \max$. [S is t . corpus length]

Because the q_j^i 's have been obtained w. ad hocness ($A_i \neq 0$), the q_j^i values must be reduced (391.30-38). The ~~max~~ correction $q_j^{i'} = q_j^i \cdot 2^{-\frac{A_i}{S}}$ is not unreasonable, and these $q_j^{i'}$ values, rather than q_j^i should be used in .10. Better correction factors under various counts (392.06) are discussed in t . recent work on SVM. T . difference in correction factors can be appreciable.

The $n-1$ degrees of freedom of \vec{a} give an additional sh-ness to t . result (395.20) i. q_i 's can be computed in a manner similar to that which I used in my ^{most} recent analysis of linear regressn, i. t . Max method

Some Mathematical techniques of approxn, involving Taylor expansion of $\ln(1-x)$ are giv. on 393.24 - 394.26. Also, because approxns are used, some statistical sampling can be used to obtain various ~~est~~ correlts & cross correlts. (292.09).

An alternative ~~criteria~~ Gove using MS error, is discussed in 389.38 ~~to~~ ^{upto} This max gives t . some results as t . more exact Max probst Gove if only quadratic cross-correlts are used in t . approxn (394.10-17)

A Goul. read discussn of t . Signific. of Pemo Results (396.02-.40): That t . log. combn. methods will often be important i. that these methods of dealing w. sh-ness may be good approxns. Mvr. other methods of combining pems should be considered, 396.10 is an empt type. watching humans is another.

A more rigorous analysis of the $q \rightarrow q \times 2^{-\frac{A_i}{S}}$ correction is needed.

396.39 spec
397.40

Note a serious ditty in t. maxen of $\prod q_j$ NO!

If one of the pairs gives $q_j^i = 0$ for an event that occurred, that pair will get zero w. - even tho its other q values are v.g. (even if they are all 1!!) This arg. is also true to some extent if q_j^i is very small for some i, j

This may be some sort of inherent weakness of use of $q_j^i \approx \sum x_i q_j^i$ (i.e. linear comb. of pairs).

- Or could it be a fault of t. Gove?

On second plot col ff is all wrong. If $q_j^i = 0$ this does not cause any trouble w. t. i's pairs.

On t. other hand, if $q_j^i = 0$ for some i, j , then t. post of y. corpus for that pair is, indeed 0, so we would get very small expected future post/symbol for that pair. Hvr, using t. corrections 391.20-396 (like $q_{22}^i \rightarrow q_{22}^i \times 2^{-1/5}$), it would be OK.

We might run into trouble if we tried to use t. more sophisticated corrections of 392.06

25: 396.90 2 solns: $\vec{z} = 0$ and $\vec{z} = \vec{M}^{-1} \vec{B}$ NO!

26: This is wrong! We want to Max $\vec{z} = \vec{M} \vec{a} + \vec{B} \cdot \vec{z}$ wrt \vec{a} .

We do end up w. something like $\vec{a} = \vec{M}^{-1} \vec{B}$ or $\vec{a} = \frac{1}{2} \vec{M}^{-1} \vec{B}$ so its OK!!

747-449

SP

Stochastic Part .

D2473

SP Planning ← see separate copy

make copy copy

3

01

≡ initially a large set of poss. solns. ~~XXXXXXXXXX~~

In task net approach, task plan is a sort of "factorization" of the set of solns. One solves each factor by narrowing it down. More generally, factorization

of this form is not poss., i.e. one must somehow narrow down the soln. set "as a whole" (or perhaps not as a whole, but not in factored form).

This working out of plan in not factored form is often done on a specific (subgoals) of a larger plan. It can be done if

1. goal is recognized as being of a particular type
2. being amenable to certain (non-factoring) approaches.

How construct plans
→ task net
→ delegate
to search
all poss. solns?
May 66
CAM
paper

18

SN

On SP: I suspect that normally one does not solve large task nets using probly curves for each ^{known} sub task in task net.

Instead, there is a certain set of net forms that are ^{readily} solvable (perhaps by table, analog or MC means) in task form of any actual SP net is "reduced" to one of these forms as an approx. Some methods of "reduction", see to express a set of 11 or sequential tasks as a single task, or any single input, output object ~~as a~~ sub net as a single task with suitable probly curve.

Viewed in this way, if a large task net is decomposable into series, parallel & subnets, one might be able to solve it, by obtaining a ~~probly~~ G curve (i.e. probly of soln, failure v.s. time) for each 11 or serial set of tasks - like an Ohm's law soln. of a Resistance network.

Furthermore, if one could do a. analog of ~~probly~~ Δ to Y & Y to Δ x fms, one could perhaps solve all nets!

6.01

D 2473 SP

It would seem to me that the $\square \rightarrow \downarrow$ is

$\downarrow \rightarrow \downarrow$ x fms are always possl. - Tho the resultant

G curves are not or directly into some functional forms as those of the component tasks,

In particular these x fms can be very worth while if there is a central or sibling dependance on the G functions of the subnets - so one would try to align these subnets to get a net w.o. such dependancies! If the dependancies are outside the set, then one is in trouble!

The v. sav-par. idea still needs some thought,

consider the $\gamma \rightarrow \Delta$, $\Delta \rightarrow \gamma$ Q. Just what would it mean if it were true? - say $\gamma \rightarrow \Delta$:

Say we had the subnet: $A \begin{matrix} \gamma \\ \beta \\ \delta \end{matrix} C \rightarrow \begin{matrix} \alpha \\ \beta \\ \delta \end{matrix}$

This would mean that dw unit of work on A, say, is equiv to $f_\alpha dw$ on α , $f_\beta dw$ on β & $f_\delta dw$ on δ .

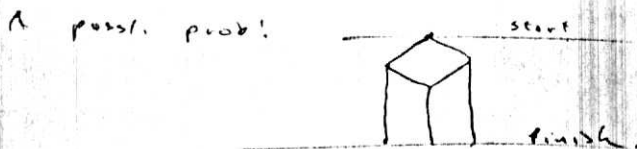
Any motion up or down on the G curves of α, β, δ is mapped into corresponding motion on the G curves of A, B, C & vice versa.

But: while doing A & C is equiv to doing β
 " " B & C " " " "

I don't know what α is equiv. to - if anything!

Perhaps the xfm is $A \begin{matrix} \gamma \\ \beta \\ \delta \end{matrix} C \rightarrow \beta \vee \delta$.

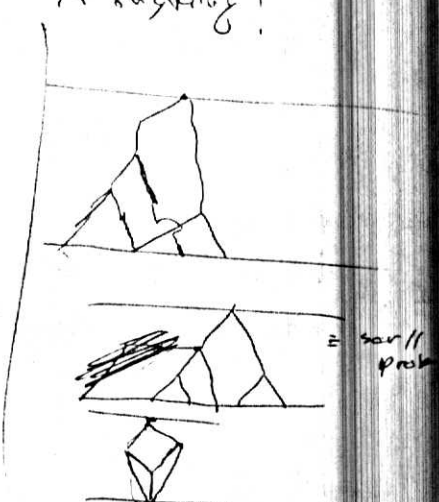
Also $B \begin{matrix} A \\ \gamma \\ C \end{matrix} \rightarrow \beta \vee \delta$

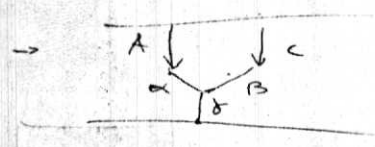
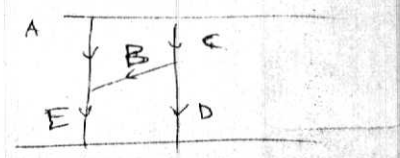


or



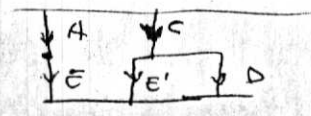
Note diodes, $E \rightarrow B \rightarrow E$ is a soln., but not $A \rightarrow B \rightarrow D$.



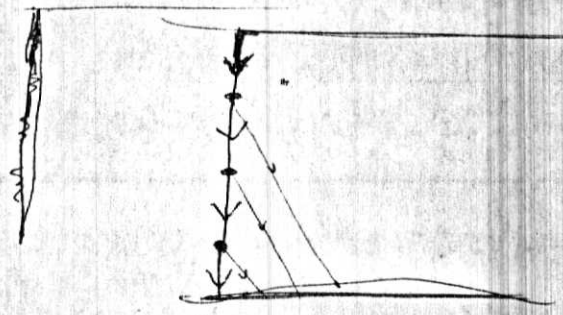


soln. of **E** \equiv soln. of α & β
 soln of **D** \equiv soln. of α & β
 " " **B** \equiv partial soln of α & β .

This is the same as
 with E & E' being 100% correlated
 (i.e. identical).



Can perhaps: A single task w. a gn. $G_{1,2}(t)$ curve
 (\equiv probly of ¹ success, ² failure, after effort t) can be approx'd
 by a sequential set of tasks w. \textcircled{a} uniform completion times,
 \textcircled{b} probly of success or failure || not so clear!



VHM (old)

345 - 17

(22)

(24)

24.75
VHM = SVH (old)

01! 294.40! Re: T. "Thru" of 299.03: If we can just show that

t. error in using Z -krest of a corpus, is bounded with some known factors of t. true strip of that corpus, we can use Z-krest as an approx. to strip - (say for a Perm or an operator or grammar) - \hat{z} with this bounded error in krest, we can be assured that for z -log enuf extra (over complete CMI) szz , we will get O.K. convergence.

This would legitimize the use of krest for a strip of "patterns" \hat{z} gives us idea as to ~~an~~ upper bnd. on error (non-optimality).

D2173: Desere Willis: He said he had read something like krest is within ± 1 of krest, \hat{z} that he would try to find t. proof \hat{z} send it to me. He gave me t. idea that this relation was true of a certain particular Machine. It certainly is true of ~~the~~ t. machine generated in his thm. 12 (Perm \rightarrow FOR \rightarrow Thru).

Re: Any way, At. 299 Thru: If we make t. thrm say t. no. of 0's = 1000 x t. no. of 1's, then our units are 1000 0's \hat{z} 1, 1 - which makes t. grain size quite large. Another way to do it, hrr - ~~mark~~ t. no. of 0's = $(n \times \text{no. of } 1's) + k$, with $-500 \leq k \leq +500$.

Here, we end up with a set of codes that seems to be "unbiased" toward any particular P (set of \hat{z} perm(S)), yet it is ^{very probly} a univ. code.

\rightarrow [Hrr, \hat{z} codes w. \approx 1000 times as many 0's as 1's are ~~not~~ good particularly good for something! - Just what is this something - just what perms?]

\rightarrow Anyway, from (W Thru 12), Say we have a univ. M_u \hat{z} a Perm P_0 . M_u can decb. P_0 in k bits, info. sense that $P_0 \rightarrow R_0$ (R_0 is on FOR) \hat{z} $R(\leq k) M_u$.

Consider t. set of all CPM's P_i . Say we want to know best cost of a gen. corpus / wrt a gen. univ. M_u . We know krest of C wrt M_u . What is \log_2 best of C wrt M_u ? We can use M_u to decb. $\sum R_i$'s. ($P_i \in P(R_i)$) (The R_i 's have t. cond. that their sum krest $\leq \log_2$ best + 1)

01:12.40: Perhaps it would be best to try to prove the opposite!

Some possible Thms:

03 1) $\forall k > 0, \exists$ a u.m.c. $M_0 \Rightarrow$ for all $A^{(Q)}$ (or for at least 1 $A^{(Q)}$)

$$-\log(p^{(M_0)}(A^{(Q)})) < (k_{cost} \text{ of } A^{(Q)} \text{ wrt } M_0) - k.$$

i.e. $\frac{-\log p_{cost}}{B_{cost}}$ & k_{cost} differ by an arbl. log amount.

06 2) for every u.m.c., \exists a corpus $A^{(Q)} \Rightarrow -\log p_{cost}$ & k_{cost} differ by an arbl. log amount.

Say 1) is true i.e. \exists one u.m.c. w. this property, then 2)

is implied. Yes!

Because: Say 2) is false, so \exists one u.m.c. \Rightarrow $k_{cost} + \log p_{cost}$ is bdd for all $A^{(Q)}$

Def $B_{cost} \equiv -\log p_{cost}$: There is no possy. now of confusing B_{cost} w. "the short length of the shortest code" — since this is defined to be k_{cost} . I think I'll transfer let $B_{cost} \equiv b_{cost}$.

So if $k_{cost} - b_{cost}$ is bdd for all $A^{(Q)}$ for a gn. u.m.c., it is bdd for all u.m.c.s. Because — say its bdd for M_1 :

so: $k_{cost} - b_{cost} < k_1$. Consider a u.m.c. M_2 . k_{cost} of x thru. from M_1 to M_2 is α_{12} ; x thru. y other way costs α_{21} .

If k_{cost_1} is the k_{cost} wrt. M_1 of a gn. $A^{(Q)}$, then its k_{cost_2} (wrt M_2) must be between $k_{cost_1} + \alpha_{12}$ & $k_{cost_1} - \alpha_{21}$.

Similar bdds for $b_{cost_{1,2}}$ make it clear that

$$k_{cost_2} - b_{cost_2} < k_1 + \alpha_{12} + \alpha_{21} \text{ — which is a bdd.}$$

OK, so just try to prove 1) (03), since it implies 2) (06).

Also the falsity of 1) \Rightarrow the falsity of 2)

So: try to find a u.m.c. for which either $B_{cost} - k_{cost}$ is bdd or " " on Bdd.

See R 38.18 for "proof" of

bddness of $k_{cost} - b_{cost}$ for UMSM.

01; 345.90; The idea of ^{dyframaly} biasing the probs of 0 and 1 so that 1. "mass
 02 "group size" is very small; Say we have 1. no of 0's =

03 $(k \times \text{no. of 1's}) \pm \frac{1}{2}k$ $(-\frac{1}{2}k \leq k \leq \frac{1}{2}k)$. So we use these codes,
 04 \rightarrow select the code of least cost — i.e. the shortest code.

A method close to this: Use 1. "Model of 1. World" model of CMI

Here we send a biased stack, sep. into a unac — w. probs $\epsilon, 1-\epsilon$; ϵ very small ϵ .
 The output ~~is~~ still has the same distribn. as w. probs $= \frac{1}{2}, \frac{1}{2}$.

This \rightarrow biasing gives, for normal CMI, a very sharp wt. DM —

particularly for large l , of ~~the~~ ~~code~~ codes where

no. of 0's = $(\text{no. of 1's} \times \frac{1}{2}) \pm \frac{1}{2} \times \sqrt{\text{no. of 0's}}$. If we ~~are~~ selected

out \pm rather best code from this \rightarrow sum we would get 0.02, etc.

(S) V H

349-45

(73)

(74)

Stooges VHM, (1986)

SVH

321.12) In t. ~~conv~~ down. string of "A" bits, clearly some bits will have more effect on prodns. than others. E.G. 1s. 1s in the coifs in linear regressn — so one would want to give less penalty to t. ϵ on the account of base bits.

Hvr., Normally, one will use only n bits in t. coifs, so that t. total no. of bits in t. down. (coifs bits + ϵ^2 bits) is minimal. If one does not minimize here, then A will, indeed, be ϵ seriously large.

"General SVH Thm:

If $A = 0$, then t. "expected value" of t. next symbol will be $\frac{F(S)}{S}$. This exp. Val. is wrt. all possl. coifs w. a pri assigned to t. coifs via CMI.

Another (perhaps) possl. interpretation of "expd Val" would be CMI

wrt various pairs w. similarly assigned a pri

$A=0$ simply means that one know, a pri, nothing about t. corpus that was coded. This ought to be easy to show,

or modify t. thm so that it is true.

For $A \neq 0$ I will perhaps have better insight after doing 012

SN

This SVH Thm is sort of nice (if true in any way). "Completes" t. CMI Theory in a neat way. Via CMI,

one can, by 1 CB, approach "true" probs as a limit.

Hvr., ~~one never knows how close one~~ while one can make ^{arbitrarily} limits of error at any pts. These limits do not $\rightarrow 0$.

Using SVH Thm, hvr., a CMI w. 1 CB, we have an upper bound

$A=0$ is so we can estimate, not for error in proby,

but ϵ , t. expected value of next from t. next symbol.

Actually, ϵ ^{in some circumstances} can be of more value than just knowing ϵ an upper bound on t. error in proby

01: 349.40 : $(349.27)^{space} \rightarrow$ Note 7. form of SVH given on 321.01-12

It may be much more or less diff. to prove it in this form!

D373 : first note that SVH may be restricted to Pairs form

which a mean Entropy exists.

Second : Can I find mean ~~value~~ Expected value of ~~prob~~ prob of next symbol } or of best next symbol

of ~~prob~~ prob of next symbol } for Pairs : (b) any other pairs.

(c) Pair that has a best value of $\{ \text{corpus best} + \text{Pair prob} \}$ or some other Gorc.

That have a mean Entropy :

Of all Pairs, the one w. the largest σ for its $\Sigma(\text{best})$, is Pair seq.

$p = q = \frac{1}{2}$.

~~Dep. on~~ $\ln(1-p) = \ln \frac{1}{2}$

Its it a max or min!

The distribution of prob is $\frac{1}{\sqrt{2\pi}} \frac{N!}{(N-n)! n!} -$ which is a gaussian distrib. of width $\propto \sqrt{N}$.

All seqs are equally likely. (\approx prob of $\frac{1}{2}$ seq).

How, if $p \neq \frac{1}{2}$, then I get $p^n (1-p)^{N-n} \frac{N!}{(N-n)! n!}$

$\left(\frac{p}{1-p}\right)^n = e^{[n \ln p - n \ln(1-p)]}$ — which moves the center of the gaussian,

but leaves its width invariant! i.e. $e^{-n^2} \cdot e^{kn} = e^{-n^2 - kn}$

will have a width that is indep of k . (Seems that I should have

known this anyway — its the same as the result on 'errors in estimation of mean of a gaussian variable.)

Now — This seems like a fine result — but is it relevant?

Well! just what do I want? : say I have Pairs : Its most

likely seq. of length N , has a prob of α . How many of

the 2^N seqs of length N have a prob betw α & $k\alpha$?

($k = e^t$, say)

Perhaps a no. (indep of N) will have a prob betw.

$\ln \alpha$ & $\ln \alpha + \epsilon \sqrt{N}$.

Wrong! If $p \neq .9$ (probly to 0), all 0's is the most

likely seq. I most likely no. of 0's is .9N.

01: 350.40

Heavens! What looks like a fantastically silly mistake!

On 279.40, I got the idea of essentially $S = \frac{F(S) + A}{S}$

S = expected future n cost/corpus symbol

S = No. of symbols in corpus plus far

A = n cost of pen.

$F(S)$ = n cost of corpus wrt pen for $1/S$ symbols of it.

Then on 293.07 I got $S = \frac{F(S)}{S-A}$

T. Augt. at 292.25 - 293.07 seems quite clear - here, n r. example given on 292.25, $\frac{F(104)}{104}$ is "correct"

$\begin{pmatrix} 100 \\ 101 \\ 102 \end{pmatrix}$

since $F(104)$

is ~~not known~~ probably not known,

$F(103)$ is known

is $\frac{F(103)}{103-3}$ is about right, but a bit by.

$$\frac{F(103)}{103} \times \frac{103}{103-3}$$

Perhaps ~~on~~ 292.25 should ~~have~~ be modified. Saying that

"~~say for pen~~ is: next 3 symbols are 000, is not a v.g. way"

SN

Note: MM is not a sequential coding method

to code them. If we are getting k corpus symbols/bit of code, at rate p (i.e. $F'(S) = k$ at that pt. for a no. load code) NLC

then we want to multiply $F'_0(S)$ by $\frac{S}{S - \frac{3}{k}}$

i.e. $F'_0(S) = F'(S) \cdot \frac{S}{S - \frac{3}{k}} = F'(S) \cdot \frac{3}{S - \frac{3}{k}}$

3 bits of code would give us $\frac{3}{F'(S)}$ symbols of corpus for corpus length, S .

∴ factor $\frac{S}{S - \frac{3}{k}} = \frac{kS}{kS - 3} \approx \frac{kS + 3}{kS} \approx \frac{F(S) + 3}{F(S)}$

if

so $\frac{F(S)}{S} \rightarrow \frac{F(S)}{S} \cdot \frac{F(S) + 3}{F(S)} = \frac{F(S) + 3}{S}$ i.e. f.

results of 279.40

While 120 ff is not all so rigorous, still, it suggests that 293.07 is too lenient in penalizing A , i.e. that 279.40 might be

just rite!

i.e.

$$\frac{F(S) + A}{S}$$

over perhaps

$$F'(S) \cdot \frac{F(S) + A}{F(S)}$$

01: 408.40:

$$\frac{F(S)+A}{S}$$

seems o.k. I'd have to go into t. argt.

more for t. $F'(S)$ correction.

$$\frac{F(S)+A}{F(S)}$$

does seem

like a reasonable distribution factor for $F'(A)$.

The arguments on when $F'(A)$ can be used, & t. various

to elements on t. old $\frac{F}{S-A}$ formula (which is wrong), are

probably all applicable to ~~the~~ formula.

15 3

(2674): Since W's Method used as a limit) is SVH method, both

obtain limits for probly, for fixed Q , perhaps one can compare them easily, for fixed L .

W's method for fixed L :

2 corps: $A^{(0)} \rightarrow 0$; $A^{(2)} \rightarrow 1$

Actually, $A^{(0)} \rightarrow 0$; $A^{(2)} \rightarrow 1$ are also poss.

consider all codes for t.

got ratio of \leq costs, using a C.B., T.

Let $C \rightarrow \infty$.

But forgot about varying C.B. for awhile.

In SVH, consider all ~~seqs.~~ ^{nonprinting} seqs. of length L or less, ~~(A, A, A, A)~~.

Perhaps re-examination of obs prob values may show inclusion theories!

Actually, I'm not sure that every non-printing seq. defines a pair (\in CPM).

22

(S1U) why confine t. ~~to~~ defn. to n.p. seqs! What we want to do

is get t. machine into a certain state. If it had to do printing to get there - O.K. - we can simply ~~neglect~~ neglect t. printing, & concentrate on $L, 0$ that occurs after we get t. t. desired state.

T. input seq. that got it into that state is a defn. of that

24

State & i. of t. CPM that that state realizes

25

(That old stuff about pairs & i pairs may be irrelevant - but I should look it over again) - Anyway - say any ^{non printing} n.p. seq. defines a pair.

Def

Then t. ~~WMM~~ (Kolmogorov version) \leftarrow SVH method says, "pick t.

pair defn string, S , $\Rightarrow (N_S + \text{cost of t. corpus wrt } S)$ is min."

Avr, I think that in this case S is always t. null string! ^{C.B. t. best pair} that defined by UMC!

To see this - put it in t. form of pcost rather than bcost.

so we want $S \Rightarrow \left[2^{N_S} \times (\text{cost of corpus wrt } S) \right] \alpha$ is max.

Now α is simply t. \leq cost of all codes for t. corpus that start w. S !

So clearly, since the null defn pair ~~is~~ (i.e. that defined by t. UMC)

has some other codes in addition to these, it will have a greater α .

If we take the (Witt's version of SVH) \leftarrow WSVH, we chose t.

01:409.40

pen \rightarrow

$$\left[\sum (\text{prob of last pen}) \times (\text{prob of corpus wrt that pen}) \right] \rightarrow = \alpha_w \text{ is Min}$$

$$\approx \sum \text{prob of all strings that can get it into that state or an}$$

"equivalent" state that gives identical probs to corpus

How, again α_w is the \sum prob of a large bunch of codes for (. corpus, and again, the null string (\rightarrow ~~unc~~) has more codes than those

\therefore greater α_w .

Yet the SVH method does seem to have some meaning & validity!

Perhaps not! — T. main meaning in SVH being the use of

$$\frac{-\log_2 \alpha_w}{N_{\text{corpus}} (\approx \text{length of corpus})}$$

is an upper bound estimate for ξ

10 t. expected mean burst/symbol of the future corpus.

i.e. it is True but trivially identical to Willis' method.

\rightarrow So 409.25 - 45.10 looks like a demolition of the old SVH method! There are still, perhaps, some Q's about!

1) whether Just how a pen should be defined — whether by the longest poss. n.p. string, or by any n.p. string, or whether primitive strings can also be legal (see 409.22 - .24 for a bit of this latter pt.).

20

on C 88.30: For large corpus of RW data, it

is likely that ξ will not exist for the corpus as a whole —

i.e. different part will have different ξ values. Just how

one should deal with this is unclear. How, it seems

clear that one can't simply use $\frac{F(S)}{S}$ [in this RW case,

A 20 since data out of data $\approx (\approx S)$ is very large &

we can choose our unc "A priori".], $F(S)$ may be more reasonable — but it is not clearly defined.

Punct, Coding

99-100

~

1974

1. punctuation, coding
This is
: / Pratty much in spirit of previous writing on Ruz stuff.

Some new ideas, hvr.

TM 99
which is an improvement on Huffman.

- 1) Use the method of Willis' Runs to code a) Bern seq. b)

Bern seq. w. changing probs as a funct of previous corpus (Markov chain)
w's method mite be Huffman's, for. ~~that~~ Most likely code - but w. "block length"

- c) Bern seq. w. varying alphabet at each symbol pt. (Markov chain)
- d) Coding w. defas. e) PSC coding.

In general, I think R. methods used in I & II were wrong for these problems. ~~but~~ I.e. T. method for coding the Bern seq was wrong at the beginning - but the ideas are correct and to be useful in further work.

2) Suppose we are coding a time series of nos. by linear regression.

T. code consists of t. following seq. of nos.

- a) The first N values of t. seq.
- b) " N coeffs
- c) The σ
- d) The seq. of corrections involving σ as scale factor

(This latter is not exactly right but is O.K. ~~for~~ for illustration)

Initially, N is not known exactly, but a probly distribn. is known.

Similarly for t. no. of signif. ~~bits~~ bits in t. values of t. seq., t. coeffs is σ .

If any of these are partly computable

& say, probly, from previous data in t. code, then this

info should be used to narrow t. probly distn. & ultimately

reduce coding cost.

T. code itself; start with a ternary alphabet: 0, 1, Δ

Δ is punctuation. Its signif. will vary at various pts. in t. coding.

T. relative probs of 0, 1, Δ will vary, often at each symbol pt.

At first, we must write out t. first value in t. time series,

T. first bit must be 1, so we don't write it. We start w. t.

2nd bit. This can be 0 or 1, or Δ can be no 2nd bit -

in which case Δ is an "end of no." symbol. The

relative probs of 0, 1, Δ will depend on how many bits we expect

to be in t. no.

01:59:40: Anyway, say we write 01101Δ (meaning 101101 is 1. first value)

[There is some Q about using floating pt., in which case, each no. would be expressed as a pair of integers.]

If there is only 1 coil, Δ is 1. next symbol. If not, then 5 binary bits follow — it is imposs. for Δ to come before 1. final binary bit. Anyway, so we continue coding, using Δ & $\Delta\Delta$, etc, to designate various punct. types. Hvr. T. next thing is that 1. probly distribution over $0, 1, \Delta$ should use all info available — both a priori (external to 1. seq) & all previous info in 1. seq.

Instead of using $\Delta, \Delta\Delta, \Delta\Delta\Delta$ as difent. punct. types, we could use $\Delta_1, \Delta_2, \Delta_3$; but 1. end result in code cost would be (i. same).

In the case of 1. coding of 1. errors (wrt Δ), this can be done directly: knowing Δ makes it poss. to compute 1. probly of each poss. error that is written out.

~~There will be problems involving coding costs.~~

T. above deriv. of coding, gives us costs in a very direct way for a "known" coding method. If new defns are devised, hvr., it's not so clear as to what 1. coding costs are.

E.g. say in 1. linear regress. situation, we want to define a new variable \equiv 1. product of 1. ~~and~~ z & s (y). We can write a punct. symbol, indicating 1. pair of integers (z & s), tell which legs are to be mult. to get 1. new variable. Just what is 1. probly of this punct. symbol, hvr.? T. problys of z & s , resp. can be $\frac{1}{N}$ each, (if there are N coils thus far), or it can be other things if there is any previous info on this.

Hvr., 1. large coding method (i.e. 1. one for linear regress.) is not universal, so I would, certainly, need a larger vocabulary of instructions (or a different set of instructions).