Nips

50

∴ How Rev 2 differs from Rev 1: (to be printed overleaf. Also in letter to Correspondents)

A table of contents ~~noted~~ has been added, ~~~~ & a description of each section.

The updating technique has been modified to include information on trials that have failed.

~~~~ An advanced ~~spiff~~ search scheme is described that appears to be ~~much~~ significantly in important systems

better than Lsearch.

Some short explanations have been expanded and various ~~typos~~ small errors have

been corrected.

---

10 { Another Nice Expo for Rev 2:    Before we explain the new update/serch schemes!

Describe it in general terms. From previous experience with various PST's on various

problems, we ~~decide~~ ~~~~ ~~~~ find that Fr is most promising PST for the present problem,   ~~~~ Gₙ
Fr to Pro                            ≅ short time.

In some books 13.16

We apply ~~Prot~~ ~~Fr PST~~ to Pro ~~~~ present problem for ~~a certain of time~~, then

~~then~~ in view of its degree of success, or lack of success, we revise our estimate at to

~~~~ which is the best PST to ~~~~ ~~~~ Gₙ.  ~~~~ we then ~~~~ apply these
                                            can or may not

Revised PST (which may be the same as Fr ) to ~~~~ Gₙ  ~~for a certain of time~~ ≅ short time
                                            was solve the problem or our time runs out

This process is repeated until ~~we/exc time runs out~~.
                    first,
          But given a corpus of historical date on the results of various PST's working

20      on various problems, how can we ~~estimate~~ find the PST that is most promising

    ~~for the~~ present problem?

---

13.16 → " Consider the set of quadruples · · ·

    At the end of section 3 (on 0₂ probs)

27 {    It should be noted that ~~the~~ time limited optimization
                    only
   problems · are / well defined when the G(·) function ~~~~
                                                        13
29                                                                              ( This alternative between ~~use~~ ~~of~~
0  "linear"                                                                     application of a PST and
          It will be noted that if x maximizes G(x), then x will also maximize   choice
                                                                               revision of our ~~estimate~~ of
   H(G(x)), if H is a monotonically increasing function. In general, however,   the best PST, continues until
                                                                               we solve the problem or
                                                                               run out of time.

    $\int_0^\infty G\, p_0(G)\, dG = z_0$  ;    $\int_0^\infty H(G)\, p_0(G)\, dH$

35   If    $\int_0^\infty G\, p_0(G)\, dG > \int_0^\infty G\, p_1(G)\, dG$        $\int_0^\infty G\, (p_0(G) - p_1(G))\, dG > 0$

    Then   $\int_0^\infty H(G)\, p_0(G)\, dH > \int_0^\infty H(G)\, p_1(G)\, dG$   need not be true    $dG = dH \cdot \frac{dG}{dH}$

                                                                                                    $dH = dG \cdot \frac{dH}{dG}$
          $\int_0^\infty H(G)\, (p_0(G) - p_1(G))\, dH > 0$   need not be true

          $\int_0^\infty H(G)\, (p_0(G) - p_1(G))\, \frac{dH}{dG}\, dG = \int_0^\infty H(G) \frac{dH}{dG} (p_0(G) - p_1(G))\, dG$

140

N.ps

Say we verbally modify $G$ by a monotonic ↑ funct: $\alpha(G) \cdot \boxed{} \ G^\gamma$

$G$ is from 0 to $\infty$.



viz, hy $\gamma$ :

Drop this for a while.

Going back to OZ probs ($\approx$ 289.00 → 06) & we are 289.22 → 40

Assuming Linearized $G$; we pick $F_\varrho$ w. h that the best expected $G$! ≈ 89.30.

This is a style pick — we then just run that $F_\varrho$ for t. specified time. We will not necessarily get a v.g. "$G$" value, but on "average" for these OZ problems, we will get the t. expected $G$'s".

Say we have a OZ problems we really need to optimize at n t. given time. We might divide up t. time into 10, say; work on t. most promising PST for $\frac{t}{10}$; then do a re-optimization of 288.35 then work on PST that looks most promising, etc. If we spend ½ time on ↗ & ½ time on optzn pms.

Instead of $\frac{t}{10}$; we can timeshare betw. 288.35 and various $F_\varrho$'s.

So, (except for the Linearization postulates) t. foregg. would perhaps solve OZ probs. In both Env & OZ, hwr, this "advanced update" will not work unless TM is able to get good model(s) of $O^2$ in 289.35 & in t. corresp. $\bigcirc^2$ for ENV.

So do explain that ① Normally & "Trainer will decide when to switch to WON. & ② that best way to decide is to do/problems via WON & viz Lsrch & Ssearch

WON seems to be doing better.

A hy/ly educated TM will usually be able to chose a very appropriate PST for most OZ problems, & so it wouldn't much be jumping from one PST to another. It may be wise to spend less time updating in such situations.

This, in general, updating is "Self-improvement" in t. larger sense, & so it would be worth while to spend 5% of available time on it.

▶ In doing Inv probs each cand could take a fairly small amt. of time ·· (Not necly!) — We are still trying PST's, & a PST has to look at t. problem & try to generate a soln···· this can take a lot of time! — In fact, a PST could be a search (or Lsrch) & takes much times as we'll give it!

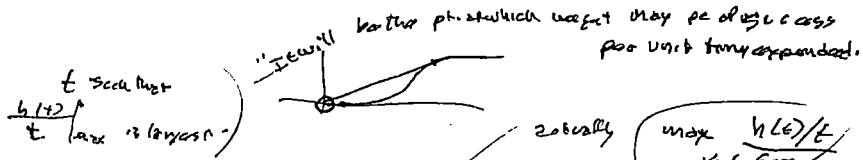.00: 288.40 : Problem $G_j, t_j$ will ~~...~~ Find an x with $G_j(x) = G^{j,l}$ within t ment

~~...~~ ⬜ That $F_l$ will obtain an x / m time $t_j$ such that $G_j(x) = G^{j,l}$

~~Declaring~~ ~~...~~ $O^2(G_l^{SM} | \tilde{G}_j, t_j, F_l)$ } Then $j, l$ product is overall ~~...~~ past

$h_l^{2'}(\tilde{G}_j) = $ optimization problem solutions obtained using the set of $[F_l]$

φG     We want a PST, # $F_l$ such that the expected value of $G$ is as large as possible.

"t=will" be the point which we get may be of success per unit time expended.

t scatter
h(t) / t large is longest·)

actually ( max $h(t)/t$ vs t· Gore )

⬜ 1) Write about "when to switch from L such to W/O IV". : Mostly involves ⬜ } one could use both & switch when one was ahead.

being able to find a good $O^2$ functions

2) perhaps mention that max "$h(t)/t$" is / Gore, not $\partial/\partial_l$ (from GHT1).

3) Is t. working model I'm using actually optimum or near optimum? It seems very "Greedy" i.e. ⓧ work on "Best $h'(t)/t$" to peak; Then recalculate Gores, loop to ⓧ

I mean Best $h'(t)/t$ is best initial (Greedy) ~~strategy~~ strategy.

4) It would be nice to have pictures of $K'$ cards
   a) Ask Grace  b) Maybe pretty candid!

5) possibly include Concept Net !

6) Perhaps Make big Zip file for allowing Download & blog papers ( so people could
20  Download them in one Full swoop!

22  Actually, probly density will be a funct of $F_l$, $\tilde{G}_j(\cdot)$ ⬜ $t_j$ and $G$.

G |
   |                   a plot such as this for each $F_l$.
   |_____         Say we use time varying Gore.
       → t
                  For each T value we have a probly profile of $G$ :  ↑ p(G)  ↑ Pc
                                                                        ⌢
                  To select t. "best" of those.           → G

Probly density of G
~~...~~ known at time T.

30 — If G has been "linearized"   is  $\int_0^\infty G \, Pc(G) \, dG$ + rate Gore ?

"Linearization" means $Pc' = Po$ of $G$ ; is equiv to $\frac{Po}{2}$ of $2G$.

i.e. $Pc \cdot G$ is what we are interested in. So t. "$\Sigma$" $\alpha$ "$\int$" $Pc \cdot G \, dG$ is
what we want to make.

If $G$ is not linearized, t. optim problem is "not well-defined". I really need more
into better & compare a good soln. In t. absence of knowledge of Linearization,
I can act as if G was linearized — which is not good, but may be t.
best I can do.   I would recommend, getting Linearizn. info.

Perhaps Title of section: ~~Probabilist~~ Conditional Probability Distributions.

(285.22 spec.)

0 (282.00~23): Start out by defining d. funcs, s-functs:

There is only one way to express D-function. ~~Section~~

(Well no:     say $x = f(t)$, $y = f(t)$ implies $y(x)$ but can be many valued ∴ not a function.

For $-functs there are ~~may~~ several ways to express them!

1)

_____●⊤

[SN] [AZ] can assign a pc to any funct of a vector   [k vars, k scalars]
$(y = F(x))$   AZ assigns a pc to them.   This is a ~~xxxxx~~ pd that is fixed by
choice of primitives & formalizing of pc assignment.   I want to be able to try various possl. pd's on
all $y = F(x)$'s.   By adding just 1 more input, (essentially "R") we get a (perhaps)
universal s-funct.   Consider a scalar funct of a scalar.  $y = F_i(x)$.   we are in AZ,
[fixed]
a/y & on all such $F_i$.   Here, in $y = F_i(x, R)$ t. function depends on R and
t. pc of $y$ depends on t. (AZ assigned) pc of $F_i'$ mult by $z^{-|R|}$.
In AZ, the legal R (indeed, any input) must be a prefix set, so TM can tell where
that input ends.  A stop symbol is a common way to end — but ~~xxx~~ does that really
give t. pd that I want?

_____

[SN] Looking at §2.1 "~~xxxxx problems~~" (in EDSIA report)
                                    improved updating and search techniques

It starts out looking for a L-search soln. — Also @ section §3 does same.

Try this:  On p 14, after eq (7)   as $\prod(h) \prod(1-h)$

"When we find a writeup ......... at time t."

26   This O' becomes part of the updated GCPD (General Conditional
       probability distribution) and can be used to guide L-search.
insert at (14.20)       ~~xxxxx~~  While it is, indeed, possible to run an L-search in this way,
       we will describe a search technique that ~~xxxxx~~ seems to be much faster ~~xxx~~
       ~~xxxxx~~   Than L-search.
                               / p 14
               Continuation:  Given a problem, $(\vec{G}_n, \xi)$ ...

32   ~~insert-text~~   Change to: ~~needs~~ a  §3.1 (oz probs)       [ change Name at §2.1 ]
       We want to find $O^i$'s such that

35             $as^i_o \prod\limits_{j, \ell} O^i(G^{j,\ell} | \tilde{G}_j, t_j, F_\ell)$ ●

       is as large as possible:

       $(\tilde{G}_j, \xi_j)$ describes the $j^{th}$ optimization problem: to find, in time $t_j$,
       an X such that $\tilde{G}_j(X)$ is as large as possible.       (in view of $O^i$)
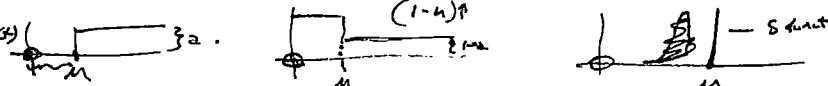       $O^i(G^{j,\ell} | \tilde{G}_j, t_j, F_\ell)$ is the probability density ~~xxxx~~ ~~PhPET Ed ofxxxxxx~~

:  I want to discuss how  s-funcs are to be realized in QATM!

Have a **new** sect'n on "Stochastic Functions".

Do give example of a 3 param ⌒ h'(t) curve as a way to get a s-funct from a 3 output d-funct.  See section A ~~very good~~ much improved update techniq.
Refere to this section for examples.

$\boxed{SN}$:  Since I really don't use α in h'(t) of t update function ( §2.1 of input)

Why not use α = 0  h(t) ⎡___⎤ ≥ a.   ⎡___⎤ ≥ m  (1-h)^t   ⎰⎱ — s-funct
                                    M                M

would I get φ pc for some races?

I would!  For  h'  I would almost **always** get ~~φ~~ p=φ except for when $t^{ind} = t_d$ at
which point consider.

So it looks like I will ~~have to~~ model α, even tho I don't care about its value — it does
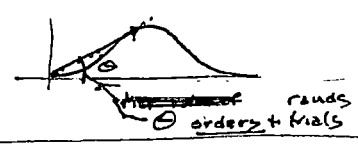contribute much to "Goodness of fit".

**So** while Eq (1) p.5 is correct, ~~that~~ Appendix A tells how to get ~~types~~
of d-funcs, not to structs of eq(1).  W

We can start our TSA w. d-funcs ( MT & TSA), but eventually, we will want
to work on s-funcs  At one time I thot ~~that it~~ was necessary to get far to do
prod'n of s-funcs. — ~~░░░░░░░~~ so it could do updating — But recently,
its become clear ~~that~~/updating can be done w. **pure** ~~____~~ d-funcs.  I still don't have a clear
idea how to construct vector outputs of functions, using economical codes, so all vector component functions
share com cost.                                                                    → 288.00

$\boxed{N.B}$ $\boxed{GHT-I}$ ~~assumes~~ (that cc ordering is optimum) assumes t. pc's are **uncorrelated** —
i.e. that knowing one trial faild, ~~t. other~~ (odds of other pc's invariant.  The ~~update~~
~~░░░░░~~ re-optimization of the $\sum_0^z \pi \bar{0}'(QA)$ takes advantage of correlns between t. cands.

In WON/GHEI,  t. ordering of t. h's is by  max values of
$t_{or}(\theta) = \frac{h'(t)}{4}$ ·· Each h' has a t value so this is **Max**.
t. largest θ for a cand, ░ its ordering index
                          ⌐becomes

If we then declare t trial to ~~does~~ extend out to that max θ pc's and stop, then
then $\frac{pc}{cc}$ obtaind is $\boxed{exactly\ right}$ for t. @ $\boxed{GHTI}$!  The only part that causes it
to fail, is the correlns of t. pc's of t. cands.

( we say "suppose" because the Gambling house Theorem assumes the ~~t.~~ success probabilities
of the trials are **uncorrelated** — i.e. when one trial fails, ~~t.~~ the ~~░░░░~~ success probabilities
of other ~~░░░~~ trials do not change.  In our update scheme, we will take advantage of these
correlations ~~░░░░░░░░░░~~ to speed up our search.

Rev2add2.tex ← Title of file.

<u>Transcript of 284½.06 — 285.1</u>

Section 3.6.2?

oa ...ill How to chose among them? There are two desirab in a h(t) function

First consider $a = \int_0^\infty h'(t)dt$ This is the probability that the associated (PST, $F_\ell$), will <u>cover</u> (will) solve the problem of interest. We want $a$ to be as large as possible.

Next consider $\mu = \int_0^\infty t \, h'(t)dt \; / a$ ⓞ   "divided by"   for those $F_\ell$ that

eventually <u>do</u> solve the problem, this is the expected time to get that solution. We want $\mu$ to be as small as possible.

In the present case, the first Gambling House Theorem

{ \fn " At a certain gambling house there is ..... Pr (b$_k$ gives least expected time to win " ⓞ_{N_0} (Sol 86 section 3.2) suggests that we will

minimize expected total solution time if we schedule our $F_\ell$ trials so that their associated $h'(t)$'s are in $a/\mu$ order: largest values first. — We say "suggests"

because $a/\mu$ is not <u>exactly</u> the same as probability of success divided by time to complete a trial.   Replace by .37 — .40

Obtaining a set of PST's of high $a/\mu$ value is a time limited optimization problem that is solvable by the techniques of section 3.

To solve our original inversion problem, we first try the PST of largest $a/\mu$.  If it has not been solved by time $\mu$, we re optimize

equation (282.21) using the additional information that the   Continuing to work on

present PST has failed up to time $\mu$.   If, after this re optimization, this

PST still gives the best $a/\mu$, we continue working on it

or in latest $\mu$      — otherwise we switch to a more

promising PST.

It will be noted that the foregoing technique is not at all, "Lsearch"

In fact, it seems to overcome a serious deficiency of Lsearch. If there are

many trials that are identical (or nearly identical, but which have different codes of the same length, Lsearch will test <u>all</u> of them — which is quite

wasteful.  The technique just described will usually test only one of

them   when one ² candidate is abandoned because it no longer

looks promising, candidates that are identical or very similar to it, will usually

also be abandoned.

We say "suggests", because the Gambling House Theorem assumes that the success probabilities of the trials are uncorrelated — that when one trial fails, the probabilities of success of all other trial do not change.   In the update schema we will describe how to take advantage of existing correlations to speed up our search   } insert at .13½-.15

# NIPS.

00
.01

because it no longer looks promising, candidates that are identical to or very similar to it,
will also usually be abandoned.

02

$\boxed{3N}$ To get $h$, $h'$, one way is to get the 3 vector $z, \mu, \varepsilon$ as a function
of $\varepsilon$, $G_n$, $S_n$, $F_L$. (presumably, we know how to do "d-fords")

For computational purposes, $h$ & $h'$ need not match exactly — just so they
match $z, \mu, \varepsilon$. We would like $h(0)=0$, and $h(\varepsilon)$ should be a ↑ funct of $\varepsilon$.
We write use some functional form all the time — say $h \sim x^n e^{-x}$,
and use TLU to get values for both $h'$ & $h$.

0

For $h$, $\sim$ tanh and $\frac{1}{1+e^t}$ or other squashing functions from ANN can be used.
Because of this form of $h$, would it be good to use ANN for models? —
Perhaps use fixed forward or more accurate (exp. features) methods.
Since the models will not be / need not be very accurate, we can afford very
approx. models — but we have to be careful that the errors of the models
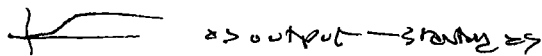do not screw up the results! — So we have to find just what approxable
models are critical, & which are not.

Perhaps mention using $\varepsilon$. 2 param $h'(\varepsilon) = \varepsilon \delta(\varepsilon-\mu)$.
10
We may try fitting these 2 params as a separate problem, then try to find
$\alpha$ as a function $F_L$, $G_n$, $S_n$.

A FET ~~transistor~~ could give [curve] as output — starting as
.22
square law, and ultimately saturating.

2.84½.34 — .35 is not exactly correct: Try: If continuing work on the present PST still gives
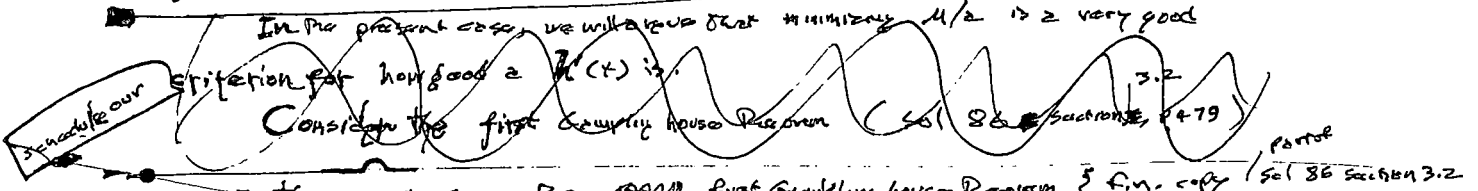best μ/d, we continue working on it — otherwise switch to a more promising PST.

30

Nks

○ ∴ Expo following lines of 13.28-.29

"When we find a suitable #0², then for ~~any~~ now problem $\bar{G}_n$, $S_n$ and arbitrary $F_\ell$,

we can find the probability $h_{\ell,n}^{2'}(t)$ that $F_\ell$ will solve that problem ~~time~~ at time t.

Then it'll $F_\ell$ : Given a problem $(\bar{G}_n)$ and a good $0^i$ function, there are an infinite number of comments stuff or not. This is already in paper!

of ~~interest~~ functions ~~for~~ for which ~~a~~ $0^i$ can obtain associated $h^i(t)$ functions.

────────────────────────────

○○ ~~There~~ How to chose a many them? There are two deserata in a $h^i(t)$ function:

first consider ~~that~~ $\pi = \int_0^\infty h^i(t)dt_i$, the probability that the associated ~~problem will ever~~

PST, $F_\ell$ will ever solve the problem of interest. We want this to be as large as possible

Next consider ~~that is~~ $\mu = \int_0^\infty t\, h^i(t)dt$ ~~~~~~~~~~~~

~~~~~~~ when $F_\ell$ does ~~a problem~~ a solution, it is the expected amount of time to

~~get that~~ solution. We want $\mu$ to be as small as possible.

In the present case, we will argue that minimizing $\mu/\pi$ is a very good

criterion for how good a $h^i(t)$ is. 3.2

schedule our Consider the first Gambling house Problem ( sol 86 sections, 1-79)

In the present case, the ~~small~~ first Gambling house Problem & f.n. copy / sol 86 section 3.2

then give references to sol 86 section 3.2 & ~~~~~~ suggests ~~~~~~ that we will

minimize our expected total solution time if we (schedule) $F_\ell$ tools so that their

associated $h^i(t)$'s are in $\pi/\mu$ order : largest value first. We say "suggests"

because $\pi/\mu$ is not ~~exactly~~ can the same as the ~~~~~~ probability of success divided

complete by the time to ~~make~~ ~~~~~~ a trial.

~~~~~~~~~~~~~~~~~~

~~~ Obtaining a set of ~~~ PST's of high $\pi/\mu$ values is a problem that can

be solved by the time limited optimization technique of ~~~ section 3.

original inversion

~~This set of ~~ PST's ~~~~ To solve our ~~original~~ problem, we

first ~~~~~~~~~~~~~~~~ try the PST of maximum $\pi/\mu$. ~~~~~~~~~~

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

□ PST /to look timing If the problem has not ~~~ been solved by time, $\mu$,

we optimize using ~~~~~~ the additional information that can

~~~~~~~~~~~~ equation (2.52.21) ~~~~~~~~~~~~~~~~

~~~~~~~~~~~~ ~~failure of~~ our present PST ~~~~~~~~ failed at time $\mu$.

continuity to work on gives

34 If ~~~~ our present PST ~~still~~ has the best $\pi/\mu$ , we continue

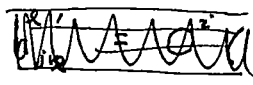35 working on it — otherwise we switch to a more promising PST.

It will be noted that the foregoing technique is not at all Lsearch.

which is quite wasteful.

It ~~seems~~ to ~~be~~ able to overcome a serious deficiency of Lsearch.

If there are many tools ~~that give~~ identical or almost ~~identical~~, but which have different codes, can

Lsearch will ~~test~~ all of them. The method just described will usually test

every one of them. This is because when/~~any~~ worked a candidate is abandoned

Nics

∞

: A rewriting of 282 07 - 29 in more readable form (eqn)

Note that

~~though~~ (13.B) uses only information from **successful** ~~problem~~ attempts

to work problems. This expression ~~can~~ *must* be modified to include information from failures as well.

Before explaining how to do this — some simplification of notation:

Let $\quad h_{j,\ell}^{i\,\prime} = O^i(t^{j\ell} | (\tilde{G}_j, s_j, F_\ell))$

This is the probability density (according to $O^i$) that PST, $F_\ell$ will solve the problem $(\tilde{G}_j, s_j)$ at time $t^{j\ell}$.

Let $\quad h_{j,\ell}^{i}(t) = \int_0^t h_{j\ell}^{i}(t^{j\ell})\, dt^{j\ell} \quad\longrightarrow\quad$ (according to $O^i$)

This is the probability that by time $t$, $F_\ell$ will have solved the problem $(\tilde{G}_j, s_j)$.

$1 - h_{m,k}^{i}(t)$ is then the probability that $F_m$ has failed to solve $(\tilde{G}_k, s_k)$ by time $t$. (according to $O^i$)

Consider the expression! (put space between $i$ and $/$ — so $/$ is above $i$)

○

$$a_o^{i} \prod_{j,\ell} h_{j\ell}^{i\,\prime}(t^{j\ell}) \prod_{m,k} \; (1 - h_{m,k}^{i}(t^{m,k})) \qquad (eq\ 282.21)$$

The first product is over $j, \ell$ pairs in which various $F_\ell$ s have been

successful *(eqn)* at times $t^{j\ell}$.

The second product is over $m, k$ pairs in which $F_k$ has _failed to solve_ *(eqn)*

problem $m$, by time $t^{m\,k}$

We want to find a $O^i$ such that (282.21) is as large as possible —

the $O^i$ that makes most likely, the observed successes and failures.

20

30

\indelible${G}$

10: 282.29! A possi approach! To discuss option of $\frac{n}{2}$: Ran do won; Ran
go back do Lerch as Alternative ney. Discuss General Trouble w. Lerch due to "many
~ conds" problem.

(N.B.) on pp 13 & 14 find out how to put space ~~between~~ between $i'$ and $/$ in $h^{i'}_{m,e}(t^{n,n})$

Having obtained a suitable $O^2$, ~~Ran~~ at least we will ~~even~~ discuss
two ways ~~w.wa~~ next the updating can proceed.

[SN] If ~~ib~~ is best to get $F_e$'s in $\frac{n}{2}$ ~~~~ $\frac{n}{2}$ order first, Ran
Brouck into 2 methods.

Another possy, is to (derb Won only:) Do analysis purely in terms of
~~~~ $u$ & $S_n$ & $S_t$ & $S^{2n}$.

Thou, give $b = a x^n e^{-bx}$ as an ant & give its $z$, $u$, & $z$.
So~~ Discuss eq. 13.1 (say) w. $O^2$ giving $a, u, b$ as functions
of $F_A$, & $G_n$, $S_n$.

Perhaps appendix w. outline of proof of "GHT #1."

[SN] ⊗ In discussn of Grammar v.s. 3 1 U models of unvarse 1 S. Functs!
T. "Grammar" unabal users functions that have (occasionally), limited Domain! i.e. there are certain
inputs for which they "have no opinion" (these may be Partial recursive ~~doesn~~ opinions —
So one often doesn't know, for sure, then whether t. machine has an opinion or not) Anyway:
Having no opinion is not at all t. same as having pc=0 for all outputs w. that input. — The
Exact interpretation is to be obtained by looking at t. proof that t. grammar & 3 1 U models are equiv.

In ~~generally~~ ~~wa want the~~ functions that have zero output for certain inputs can be "outlawed" — so
if a funct does have zero output, we regard it as having "no opinion" about that input.
( This, of course, is subject to change of opinion" as $\leftarrow$ B3 ↑ )

NIPS

2, 9, j :   i, m, n.
8, i, j  :

>o : 279.40  :  P13 of rayout:  for "eq. 13.1"   ($\approx$ equ 5.5)  write. :  $\overline{\phantom{w}}$ derive this

$$z_0^{i'} \prod_{j,\ell} O^{i'}(t^{j\ell} | (\tilde{G}_j, s_j, F_\ell(\cdot))) = z_0^j \prod_{j,\ell} h_{g,j}^{i'} \qquad eq (13.1)$$

_After 2 lines following  eq (13.1) . . . . .  $(\tilde{G}_j, s_j, F_\ell), A_j$ Pars. $\leftarrow$ 2 line 13.27_

07   Though (13.1) uses only information from _successful_ problem solutions, we can
discuss in modify to include information from unsuccessful _failures_
as well.  Before explaining how this is done we will simplify the notation a bit.

Let      $h_{j,\ell}^{i'} = O^{i'}(t^{j\ell} | (\tilde{G}_j, s_j, f_\ell))$

This is the probability density that PST, $F_\ell$ will solve the problem $(\tilde{G}, s_n)$
at time $t^{j\ell}$.

Let      $h_{j,\ell}^{i'}(t) = \int_0^t h_{j,\ell}^{i'}(t^{j\ell}) \, dt^{j\ell}$

This is the probability that by time t, $F_\ell$ will have solved the problem $\tilde{G}_n, s_n$.

$1 - h_{m,n}^{i}(t)$  will then be the probability that $F_m$ has _failed_ to solve $\tilde{G}_n, s_n$ by time t.

Consider the product expression

$$z_0^i \prod_{j,\ell} h_{j\ell}^{i'}(t^{j\ell}) \prod_{m,k} (1 - h_{m,n}^{i}(t^{mk})) \qquad eq(282.21) \quad italic$$

The first product is over $j, \ell$ pairs for which $F_\ell$ has been _successful_ at times, $t^{j\ell}$
The second product is over m, k pairs for which $F_k$ has _failed_ to solve
problem m by time $t^{mk}$.

We want to find a $O^i$ such that 282.21 is as large as possible $\leftarrow$ Grace has.
The $O^i$ that makes most likely the observed successes and failures. $\rightarrow$ 284.00 spec.
$\rightarrow$ 285.00

Next correct to details of 13.28 - 14.23.  on 2, m : Also give expression for $\sigma^2$.

$M_0$     $\int_0^\infty \alpha x^n e^{-ax} \, dy = \alpha \, n! \, a^{-n-1}$

$M_1$     $\int_0^\infty \alpha x^{n+1} e^{-ax} = \alpha(n+1)! \, a^{-n-2}$

$M_2$     $\int_0^\infty \alpha x^{n+2} e^{-ax} = \alpha(n+2)! \, a^{-n-3}$

$M_0$ :   $\int_0^\infty \alpha (x)^n e^{-ax} \, dax = \alpha n!$

Normed $\Gamma$ distrib: $\boxed{\dfrac{1}{a^{n+1} \cdot n!} \, x^n e^{-ax}}$

$\alpha x^n e^{-ax}$   slope = $\alpha(n x^{n-1} e^{-ax} - x^n a e^{-ax}) = \alpha \dfrac{x^{n+1}}{x} e^{-ax}$

$\dfrac{M_1}{M_0} = \boxed{\dfrac{n+1}{a} = \mu}$

$\dfrac{M_2}{M_0} = \dfrac{(n+2)(n+1)}{a^2}$

want $\mu = \dfrac{n}{M_0}$

$\sigma^2 = \dfrac{M_2}{M_0} - \left(\dfrac{M_1}{M_0}\right)^2$

$= \dfrac{(n+2)(n+1)}{a \cdot a} - \dfrac{(n+1)(n+1)}{a \cdot a}$

$= \dfrac{1}{a^2} \cdot (n+1) = \dfrac{(n+1)}{a^2}$

$\boxed{\sigma^2 = \dfrac{n+1}{a^2}} = \dfrac{\mu}{a}$

$\sigma = \dfrac{\sqrt{n+1}}{a}$

$(n-1) - xa = 0$

$x = \dfrac{n-1}{a}$   This is slope=0

BIG error! 30.00

$\mu + \sigma = \dfrac{n+1+\sqrt{n+1}}{a}$

MX : Mechano XLta.

NIPS

○! : Using ~~typing~~ Grammar Models that are "Non parametric" — or just Grammar Models
that someone (supra SSE.

STATISTICAL Machine XLTN

: On Mach xltn; statistical. ~~[strike]~~ Franz Joseph Och: As he is doing it, he puts in a bunch of sentence pairs from $L_1$ & $L_2$. ($Q_i, A_i$) a Prom. is expected to take a new $Q_j$. — So, it is a QA problem. — But as it is, it probably doesn't do incremental (very. —

~~Also, probably i. ...... noted~~ But there may be some "incremental (mg." by successive improvement of the models.

Anyway, this looks like it mite be very useful for getting ideas for the arcus General QA problem. (Not nearly Linguistic; — Solas. tending to be broad statistical distributions, "Glamorous" quite different from Rosen in NL ) Also, I'd want y. system to learn a second lang more rapidly, after it has long to KHt a first pair. (Also, perhaps, t. two parts could share a lang in common".)

I have a bunch of Och's papers on D:\PS\ M Translation—statistics). I didn't copy all of this, hrr. / I downloaded 1 pm! he has 2 more that I didn't download.

Q: Can I use his pap (w. innov modifin) to learn Algebra i go to ~~t~~ Phrase? >

2 papers on Alignment: I started w. shorter paper (both 2002): They use HMM (hidden Markov Models) — like in Voice Recogn. — ~~Phrase~~ (I downloaded a paper on re-latin of MT to V recog.). [I had it printed up: Put into Grammar files]

➙ An 'Impt. Characteristic of MT is that words in $Q$ correspond to words in $A$: Certainly this is not true in many other kinds of QA problems.

— What about power of mapping from ( word seq ) to ( word seq )?

It may be that t. alignment problem(s) of ..23 are critical to understanding what they are doing.

➙ Another impt. difference betw. my QA Tm & his MT / long pm! That I expect good results from small ssz. — Much smaller than this pm. uses.

Other Point: word to word mapping i his use of excessive ssz, what are t. xfnns that are learned? Could they be generaliz so I could use them in much more general problems!?

➙ Is his enormous ssz necy? Th. enormous ssz may be due to using

Q: 275.17 : **OOPS** — critique on.

Re reading 4. first few PP: I think he's over hyping how good it is in / optimung / modifying its

PC's: In "Optimized such over Such methods"

In fact he only has 5 pc modifn. mst3, & it is by no means clear that th system using

them (or any other similar in my [scan? method etc.]) is at all "optimum" in any way.

of these 5 mst3, only ~~the~~ one, boostg, has been found useful.

→ That update routine of ~~th~~ time limits (Time-dependent Rins bar). $P_2$ is actually **not too** . ~~///~~ It multiplies ≥ cc by between 2 & e each room". — Tell J. Ruiz !

(Actually, my original/my analysis was wrong.)   Say $\tilde{T_2}$ is time spent by # on cond. in $p < p_2$

so   $T_2 = \boxed{///}(1+\tilde{T_i})P_i$   This is conditional $\bar{T}$ on number of trial.

$\tilde{T_i} = P_2 + P_i \tilde{T}$ ;  $(1-P_2)\tilde{T_i} = P_i$ : $\tilde{T_2} = \dfrac{P_i}{1-P_i}$ ;  $1+\tilde{T_i} = \dfrac{1-P_i + P_i}{1-P_i} = \dfrac{1}{1-P_i}$

$\displaystyle\prod_{i=1}^{n}(1+\tilde{T_i}) = \prod_{i=1}^{n}\dfrac{1}{1-P_i}$   · if n is large & all $P_i = \frac{1}{n}$ then $\left(\dfrac{1}{1-\frac{1}{n}}\right)^n = \left(\dfrac{n}{n-1}\right)^n = \left(\dfrac{n-1}{n}\right)^{-n} = \left(1-\frac{1}{n}\right)^{-n} \approx e$

Hm if $P_i = \frac{1}{2}$, $P_2 = \frac{1}{2}$   [scribble]   $\prod = \left(\dfrac{1}{1-\frac{1}{2}}\right)^2 = 4$   So it looks like factor is e minimum,

& can be >> e   say $P_1 = 1-\epsilon$, $P_2 = \epsilon$   $\dfrac{1}{\epsilon}\cdot\dfrac{1}{1-\epsilon} \approx \dfrac{1}{\epsilon}$ ∴ arbly large.

So, if one of the $P_i$ is very close to 1, factor can be **very large**.

So, ~~if~~ | factor of / 3 is optimum, a factor >> 3 is bad.

My impression is that usually all of the $P_i$ will be << 1 so th factor will be $\pm e$ .

---

~~2 papers on "Bellman ??" : I've/ started / reading th shorter one.~~

~~They use HMM (Hidden Markov Models) like in Voice Recogn (Bra??'s paper I copyd on relation of M.T. to Voc. Recogn.)~~

I have found NIPS 2§8.06 → .28 : (A dct. on **How to use factors**). Integrating this into
§2.1 ~~dua~~ (Improved. Updating Technique) may do difficult. I **could** just add it as a footnote
or appendix.   Best rewrite this part of th section !

As for the + **won** method; I could just add it as an extra section, saying that it
will probably work, but that I am not sure it will work for all problems, ~~so~~ I have not
yet discarded the **Lsrch method** .   Is there any point in this report where I discuss G-HT #1 .
($\frac{P_i}{cc}$ is optimum ordering) Perhaps not !   So I will probly have to **introduce** GHT1.

In discussion of || updates, note that the d.f.'s of other PST's will not change much at all
$T$ gets > μ — since no soln for $T < μ$, is not at all surprising ,   $\displaystyle\prod_{k=1}(1 - a_k \int_0^T x^n e^{-bx}dx)$
   ↳ a new & of current Pst distribn.   can be too consuming to compute !

If I use $x^n e^{-bx}$   th products for th successful cases will be simpler,
but th product for th **failures** will be quite complex !   what is $\int_0^x x^n e^{-x}$  See Gr R92 : 2.32 ff
   $= x^n e^{-x} \cdot \sum_{k=0}^{n \to \infty \text{ also works}} \dfrac{-x^{-k} n!}{(n-k)!}$ ) $= \dfrac{x^n e^{-x}}{n!} \sum_{k=0}^{n}\left(-\dfrac{1}{x^k}\cdot\dfrac{1}{(n-k)!}\right)$

Add extra section on √ **won** method !  Give here Args about GHT#1, etc.
discuss "|| updating", but mention that I have included both methods because I am not
yet sure that WON will work ! Args for its effectiveness are good, but ~~I can~~ I don't/yet
certain that it will work

7.27.03
NIRS      RP (report)

0:     Refs to ~~revisions~~: 263.00 - 40 ; 265.00 - 40

Comments by J. about:

1.10.03 ① on the "sections". ② He _says_ OOPS will improve its "trace" specificity SS of SRHS — but it is unclear (in present OOPS model) how Refs would occur! — Re: 33 + entire
Q of how OOPS is to do McQ long.
He mentions "OOPS RL" (Reanf. long) — tries to Max fut. reward.

Comments by J. on ~~step~~ Second version of report!

10

1) In 11.27.02 Letter to J: # removes 3 & 4 area of import. — (J. didn't remember
my remark #3, so I guess I didn't put it into/report ~~there~~ + tier transp. of
So if my remarks # 3 & 4 can be inserted in to report f perhaps
2) subscripts) But see that they have not all ready been included — since
those remarks were _before_ first revision.

A Good Approach! 1) List sections, etc, that I have to write & write them.
2) Make tentative table of contents. 2) Make list of corrections, Additions I want,
— Write them & insert them. || 263.00 # more or less does this.

Sections 1) Table of contents w/ summaries
20     2)

3G

0 :277.00 I want to do 2 things in next month: ① Review status of TM, record decisions + goals.

??!! (See Sol 89 on Plug) ② Finish up — or make next revision of Report.  → ≡ RP

The following text will (alternately) discuss work on these 2 problems.

‡GTM is pretty to interest grad. students

03 RV: See §6 of report: pp 18, 19, 20 :

Also see Sol 89 for "present status, what needs to be done, etc."

Some myth kinds of TSQs that mite be tried.

1) Elementary Algebra

2) Symbolic Integration. Try to learn techniques of present symb. int. 'expert systems".

3) Start w. a well developed expert system (in Literature)

Try to devise TSQ to learn to do what ES does

10 4) E learns English about Algebra — Then gradually move to other subjects that TM knows about: a/o have to learn via data in "English" (or formalized subset of English).

NOTE 1. writing TSQ's has 2 goals: ① Teach TM ② Teach TRAINER how to write TSQ's. How to teach TM. Maybe discuss "Hints", Skinnerian TSQ's, "Ideal CJS at each aim.

List of TM.

( July 20, stuff.

16 See 253.0p.—90 for discussn of Sources of TMs. to be grads at MIT, say.

18 5) Try to start how Then → List more problem in MLrys are shd how this Approach addresses them. ( I think ∃ previous note entries) see 253. 27 — 40

20 Main probs: ① Great variety of probs solvable.  (IV) Computability subjectivite  This Declines to ML problem  Quite diffrent! N.B.
② ③ ② Adequate induction Algms.
④ ③ Adequate srch algms.
⑨  ④ The "educate from baby to college: Serve pably, lidk Exp starrovers approach is bad.

34 Solved by ②Inc. (ray) ② Great variety of probs solvd ③ Good transfer lrng. algms.

:00!  :  So we have to set of functions of 275:34 ~ .40 : Since this set of functs
.02  &. orbs to compress, we want to minimize its cost, so we devise a/enumer, that
assigns pc's to the (functions) we use.

T. system would work just as well if we had S-functs in line (.02), — & in some
cases, it might be easier to work probs in that form.

.05 ——— Which brings us back to y-system of 271.00 ~ .20, first we wanted to show was Universal. (fine)

⊙→ T. Force s-ff gives good understanding of Section 1 of the report where I used
Recognition functions (R) & stochastic operators (R(?)) — Recently I'd been
thinking of it as very AH wonderful, but it's not at all AH : it's very General) — Also its
a way to go from a formalism for d-functs to one for S-functs.

.09
0  :267.40  — On ANL again. I just wanted to get some CJS's for learning up to general Polish notation
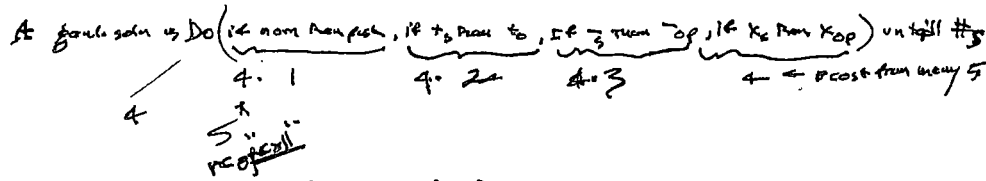2.68:27  evaln.  & After that, get CJS for k rather than just 3 binary operators.
272:16  Then have to learn 1 & 3 and arb. operators.
An alternative TS would learn 3,7,+; then 3,7+8+, then 3,7+8+5+ oct.
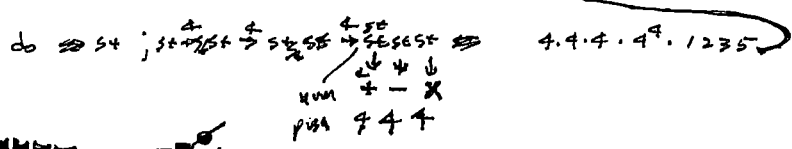Then maybe 3,7 —, then 3,7 —5+ oct.
The solns to those problems could be unique & easy for one to discover if S uses only d-induction.
If I allow S induction, there are more poss TSQ's & more poss. solns.

A formal soln is Do(if num then push, if +₀ then t₀, if ÷₃ then ÷op, if k₆ then k_op) un till #₅
       4·1      4·2      4·3      4 ← pcost from many 5
4    S⁴
     pc of call

So rc ≈  4⁵·1·2·3·5  = 4⁴·5! = 256·120 = 27.120k

probably more) ←————————   ·····→ So 27,120×16 = | 433.920 k.
                                              = for final soln for arby expression
do ⟹ 5+ ; 5+5+5+ ⟹ 5+5+5+ ⟹ 5+5+5+5+  ⟹   4·4·4·4⁴·1235 ⟹    only 3 binary funcg hvr.
                  num + — X
                  push 4 4 4

Do I have all details correct?        If num then push :
Problem is "6" soln is "push"      4  5   4   rc = 80
push 6 on stack;   put soln in store as a poss call.

Actually, t. soln is just "push".            for t. problem  6,8, t. soln is "push, push"
for problem 6,8+s, soln is push, push top.   No! The programmer could just do unconditional
                                              "push" so is a poss. soln. "push" is
At first content of store will be "St"s only  not.
if solns. to probs.

Bo rather even "total funct in just + & —", then +, —, X , then + — X ÷

NIPS

: I'm not so sure 274.3# -.40 is of any value! ☺ . → Woops! It looks O.k. now → (.34)

Perhaps wait until I have found s-funcs that are useful & I want to convert them from

Gramm to 3IU or from 3IU to Gramm. — e.g. 272.09 - .16

_____

.07 → [SN] ·OOPS uses incremental lrng only ▬ via boostg : It does not use definitions

shared between problems. If true, it's a very serious criticism of OOPS!

Hvr. OOPS does have Memory of previous solns. to problems & I think J.

expects to be ▬ "editing" ▬ them to use in trial solns to new problems. —

But this is not saying much! T. only thing he's done so far is boostg plus a sneaky

insb. to store & use "patterns" of Token type. So he has not really solved t. problem

of how to do "Transfer lrng" ▬ except, perhaps, between problems in a same "set"

___ like for t. grammar problems, he did try to reuse many of to old successful PGMs.

● → But did he retain ▬ tokens derived from previous problem solns? (even in a same

▲ set of probs "(e.g. all in t. "same type" grammar problem)

I have to reread t. OOPS paper again. Many things are unclear as to how it

works! — But in general, his facilities for ① Recalculating PC's ② use of previous →279.00

.17 solved problems; is Minimal.

.18 Review: I want to do 2 more things now. ① finish up ' report. ② Defn of

"state of THL" project: what has been done, what needs to be done, what can be

.20 worked on now ▬ as a "well defined sub-projects" for Phd work say → →277.00

.21 for David Lindsay, ▬ & MIT prospects ( & long term v.s. short term goals

.22: 274.22 [SN] N.B In order to allow Universality, the individual d-funcs ▬ probably most be able

to say "no output" for certain inputs. I'm not sure of this, but ti proof of 274.12-.26

demands have t. funcs have that capability. If all d-funcs have legal outputs for all inputs

I don't know if Universality is possl. or impossl. In my earlier analysis of

this problem, I may have assumed output for every input & perhaps proved that

universality was impossl. ▬ (not realizing that it was t. crippling assumption).

.29 My treatment of "Recognition functions" in my present (Report ( section 1, pp 6, 7 (...maybe 8)

.30 Does find a way to create ▬ funcs that have output for only certain inputs, but t.

.30 technique would work just as well if PC's were d-funcs rather than s-funcs.

.32: 272.16 In fact, t. Analysis of "R funcs" in (29 ff) could be a useful way to think

& bcet. 272.09 (.22 ~3)

.34: (OO) ▬ Well, this makes 274.3# -.40 more ▬ ! The issue is that one has 2 d-funcs — each one has

.2 been ▬ designed to work only on certain "kinds" of problems. (T. funcs themselves tell us which kinds

of problems). Each d-func has a wt. T. resultant (potentially universal) s-func is a

probbistcly wtd sum of all those d-funcs. For a given Q, it has to be normalized

since t. sum of all PC's need not be 1 (it usually ≤ 1) — not nearly ≤1 — it can be 1 usually; because of (.39 -.40)2 certainly not.

.39 Q: do t. sum of all wts of all t. functions have to converge? In 274.19 → .22 they don't — may

.40 are 2⁻ⁱ for i=1,2,3...∞ — they diverge rapidly ☺.

Nips  (MATH)

Grammar ⟷ 3·Input Univ  equiv.
Prot §§ ·12, 26  | Also see 271.2c~·29

00 :  [SN]  If, say ordinary arith. were inconsis: this need have no serious effect on R.W.
T. inconsistency would be obscure & would be very rarely relevant. We would still
use Mathm t. usual ways, because we know Prot t. inconsistency is of "depth 1000", say.
But in general, Math is just a formalism, that happens to correlate well w. events
in R.w.  If Math were inconsis: — we would still know how to use it to successfully
__product events in R.w.__

·08 : 273.40 |   Move details on t. "Many $\phi_j$ model"   Each $\phi_j$ is a Baysian hypoth.
If / produce __no__ output for certain Q's.  The predictions for a given Q are renormed to
 ~each way
t. take into account t. $\phi$'s that are "silent" for that Q.

0      Maybe different·××··08; Baysian hypoths have each a pri, __and__ a p.c. of data
if that hypoth is true.  T. $\phi$'s on t. other hand, give p.c.=1 to each of its outputs.

                                                       ┌ of 3IU
·12  O.k.  271.21~·24 is exactly correct :  ═══ From input #1 & each R input,
defines a d. function t. Q input. For a given R, there may be no output for
certain Q's — for R.more, since t. funcs are partial rec. ═══ we can be
uncertain as to whether there is a legal output or not.
        For each Q there can be a different penalty set for its R's that give output.
·17·18  If, for a given R, t. machine asks for another "R bit, then that R "has no output",
                                          prefix
·19         In linem  ·12~·19 , for each #1 input to 3IU, we get that set of d. functions
·to      on Q :  One function for __every__ value of R ---- T. set of R values for which
         t. funcs are defined do __not__ constitute a prefix set. This set of R's is __all__ ═══
22       __finite binary strings__  Various subsets of course, do constitute "prefix sets". For a given Q, t. set of R's that give output,
                                    ══ form a prefix set. → 275.22 →
         ┌ from ·12~·22  It seems clear that any a 3IU ▨▨▨ s. function can
         | map to a __set__ of d-funcs — (that have outputs for some Q's but often
         | not for others)·]
26       → To do the other way:  It would seem ══ obviously poss|. — ══
         The model (of ·08~·10) with each / hypoth only having prob'ty of 1 or $\phi$ for each poss|. output.
                        "fashion" ← Bayesian
         Prob by 1 as a d. function & having certain output, or occasionally __no__ output.
         Since this model does define a s-func & 3IU can deliver __any__ s-funcs exactly,
0        this model is simulateble by 3IU.
21            But here is another method of simulation that's of interest, because its more "__direct__".
         Say we have a $\phi_j$ see w. p.c.'s $a_j$ resp. ($a_j$ are not usually / power of 2  — so we can limit to Bay sum)
                                                        int.
         by expressing a $a_j$ as t. sum of int. powers of 2, we break up t. associated $\phi_j$ into a larger set of
         identical $\phi_j$ w. different (power of 2) wts.  So for each Q we have a set of p.c.'s that
         are powers of 2; for all different outputs for that Q  $\sum 2^{-(s_i)} \le 0$. so we can (by inverse of
         Kraft inequality) assign t. l. to strings so they constitute a prefix set.
         or  $\sum n_i 2^{-i} \le 1$   $n_i$ = no of p.c.'s w. wt. $2^{-i}$ .
         This enables us to assign R values to each $Q_i$, $A_j$ pair: → 275.00, but move closer, 275.34 and 275.22~3)

Frank Tipler : Physics of
Immortality

N1P3

20 (spec 272.40©) : Actually, I'm not so certain about t. objections 272.36 - .40 : it may not exist at all,
01  ∴ if it does exist, it means only, that # 310 → Gramm isn't exact inverse (usually) ⊆ of Gramm #310

Apart from t. doubts of 272.36 - 271.01 ; It looks like AZ ¿ OOPS are both
(potentially) universal S. functs.

To use Prams as S functs may not be so feasible for t.  QA problem :

.05  We have to get a single d funct that gets all of t. A's "correct" .... which is
→ a BUG, since t. same Q can have several diferent A's.
      This seems to invalidate t. equiv. of "Gramm" to ZIU. Maybe related
10  to t. "diddy" of 272.36 - 273.01         — in Grammers

11      .06 is wrong. For each  A_i , we  sometimes are' lucky & usually
( but not always ) have t. some  $\phi_j$ .  In general t. Best (cheapest pc)  $\phi_j$  will
give a hy A_i for most Q's ; but for some Q's it will give y. "wrong"
answer or simply have no output at all ( i.e. stop before output or loop w.o. output)
for each Q , we want at least one of t. present t. array of  $\phi_j$ 's to have t.
"rite" answer.                                                           [UB: 272.09] - .16
      So we want a set of  $\phi_j$ 's such that .... (See 271.00 - .15 for how it's done)
 T. cheapest pc  $\phi_j$  gets most of t. A's rite, T. next nyoest pc  $\phi_j$  get most
20  of t. rest, rite, etc., etc.
21      T. for $\phi_j$ 's Mindful of | 272.09 - .16 |
                                       d. functs
.22  : 272.16 → In 272.09 - .16 I have ll nodes of t. corpus — ( Q_i → A_i functs) ;
Initially, these  $\phi_j$  are informationally indip. — Ray each codes sequence, but
do not share cores.   We then decide to code this set of  $\phi_j$ 's as an s-grammar,
to minz. t. entire reast of them, so this set of  $\phi_j$ 's can be used
for predn., as 271.06 - .20 .  T. present idea falls once just what I want to
optimize — just what I want short codes for — because I have a  top Goal
to t. entire procedure .... how to get max pc of t. ~ QA corpus = $\sum_{q_0}^{\bar{0}} \prod_{j=1}^{n} O^j (A_i | Q_i)$

30      It looks like grammar : 310 are interconvertable, but I am still not sure!
        ZIU
Another diddy in ZIU, if an output has a pc of >.5 ¿ close to 1, say, t only way it can
do this, is t. have many codes. (To get as close as 1 - ($\frac{1}{2}$)^k , one needs at least k codes.)
It could have many pc's close to 1 — say essentially d-prodn.
→ I had better write up t. proofs before I forget them! — Also tell how to deal w. various
apparent Anomalies.
        First note : In "proof" of 271.21 - .24 : For each possl. value of k! certain
Q_i's will have output & others will have no output. Each Q can have a different pla Rix set of k's
(This has to be true! Say t. rite A for Q_i as pc...

3 IU.29   dofn ≡ ∃ input UMC : (S-function)

30: ~~37~~1.40  :  Re: ANL : what I've been trying to implement : A method of ANL in which
I make trials. @ ~~start~~ some trials will be of ~~some~~ value in coding.
.03   I want to be able to use Pieces ~~that are~~ " somewhat useful " objects to
obtain new trials of hyper expected yields.  (.09) seeming to be of value here : it is
a more formalized version of .00ff
e.g. " Push, push Xop " is of _some value_, so I should keep it in mem;
just _how_ to use it to help create new cards, is unclear ( ≡ t. problem
of 5.56 @ ).

09  (.03)         I _may_ want to do several ll codings of t. corpus : Each one  uses  ~~uses~~ different trials —  → See
                                                                                      269.23 - .40
has difrnt Successes, has difrnt ~~□~~ Defmitions. The ~~■~~ defns have pc's
that are peculiar ~~■■~~ to each ll code.  If we find a conc. that seems v.g in
~~one~~ ll codes we will usually try it in ~~■■■~~ The other ll codes as well.
Basicly, what I'm aiming for in .09 ff ~~■■■~~ is to find difrnt ways to derive difrnt
                                                                       [ Also note 275.22
kinds of Regularities : ( t. treatment of ll ~~of~~ codes at 272.24 - 273.21 ) ( 273.11 - .21 in particular )
                         └ is very relevant here

'17         HA! I had forgotten how's got into the problem of 270.08 ( t. universality of ~~■~~ parallel
S. functs ) .... It was from 269.23 : I had an S-grammar giving P.d. on strings that were
'19   d. functions ( Q→A ).  So _this_ trivial answered t. Q of how I got t. PC's of this sort of
0     ── d functs
                    This leads to t.  (impt.       Q of ~~is~~ t. (AZ) systems universal? ( say they have Turing —
                                   interesting)              OOPS
- complete set of uses
         Also ( t. unclear ) Q of " universal in a useful way ".

24    [     Anyway, t. " BIG Breakthru " is 271.21 - .24 : that any ~~■■■■~~ that gives > 0 pc
      [                                                          P.D.
      [ to every poss'. d.funct ( Q→A ) must be universal / in a sense that its output ⌐ pc is
      no_more than                                            S funct
      ~~only~~ a constant less than t. pc of any finitely derivable S funct.

         → It would be good if I had _easy_ ways to _switch_ betw. these 2 ways to represent
      S-functs          S-lang            d-
(Dot).29  @→  Say I have an S-lang for funcs : say it's GFG.  To get it into 3 input UMC  ( 3 IU ) form
0                    ~~XXXXXXXXXXXXXXXXX~~  Since I have a P.d. on t set of objects, it is poss'. to
derive a set of R string. assignments to implement that P.D. ( w. Jis did it for a finite set of d-funcs.
I think  Cover & Lvang (1978)  did it for a countable no. of d-functs.
   So, t. Grammar bern plus a " free app'n " to xpm pc's into psm. lengths, gives t.
derived  3 IU.  So t. 2 methods of deriving t. final S-funct, are exactly
equivalent in p. cost of derivn.  # This is a very _formal_ way of doing it — may not
35     be _practically_ achievable, hvr.
.76          To ~~go~~ from 3 IU to S-lang over d.funcs : ( 271.21 - .24 ) does it, but _not exactly_ :
      For each R value, we get a d funct over all Q. imp'ts. But incremental, R₁ and R₂ may give same
.39   outputs for Q₁ , but usually they will not give same outputs for some other Q₂ ...
      This is different from ~~t~~ Grammar → 3IU of .30 - .35          → 273.00 spec.

## Nips

so: 270.40 :  Say I had a large/finite set of k — cardinality

d-foncts, à = {Q, A} corpus.

Say, for each oft. Q, A, at least one function would do it. — So of [φj] are
d. set of functions, we can represent t. set {Q, A} j=1|n , by a set of
n φj funots that do t. correct mapping. In some cases, > one φj will work so
we have pairs or triplets (or whatever) for several {Q, A}

.08  If we only have one φj for each Q, A pair, we have a Ramsey of n symbols
& we can assign pc's to t. φj by Lap's rule (or based Laps' rule).

O →  If we have >1 φj for some QA's, we can make t pc of t.
relevant QA's by sum of t. pc's of t. φj's that work. We can then assign
pc's to t. φj's so that pc of. corpus is max. (I dont know if this problem

k5 —  is hard or easy, ) — That this is correct is evidenced by t. proof of (.21 – .24)
        If k > n, then it would seem that we are not too my much induction.

If k << n then we might be doing good induction. We can compute entropy
of corpus. It is minimal when we have most wts on only a few

ω  Qj's.



21  A possl. way t. system under work: Consider t. 3 input ump model of a  {3 IU}
untvl. pb.  We can make a set of φj's from it!
For each value of R, we have a unique φj. The wt is 2^{-|R|}.   for all {Q, A} taken together, t. same R value.

24 → That's it! T. 2 systems are identical! GREAT! It simply proves
that t. 2 systems are Equiv. → see 272.24 for 2 useful interpretation.
Hvr, for t. 3input Ume model, we are a single φj that we desrb corresponding to
R = Λ (nul). Each value of R enables t. desrn of a new φj, and (R | it it's bcost, over
most of R=Λ". ( This bcost is t. first input to t. 3 input Ume.)
total
30  So t. set of φj's that simulate t. 3 input ume have relatively small total costs.

Easy to get from 3 input ume to set of [φj]'s: But can we go the other way?

Say we have a set of 3 φj's  push, push ((+)(−)(×)/op) . The bcost of coding all 3 is
not so big!  r push, push is common, Then t. 3 distnt ops. This total bcost is needed
for t. corpus coder of .00 – .20. [t each of t. Oj's is improved — say by
"If ts then push push xop". — This has a hy pcof working (≡ Getting rite Answer)

A big problem is usually: How to get by pc's of t. set of S-functrs'd?
This is easy to d in t. 3 input ume.  One way it is done in t. 11 stands is
by 272.17: Use an S grammar to generate t. strings that represent t. [φj].

NIB

This instr is quite difrnt. in spirit from what I remember about Soart ANL!

It may be that we wouldn't get as rapid ↑ of recog of solns (scaling) as we did w.

t. old soorb. soln ≠ 2/0 that "context" will be broken in a more natural way,
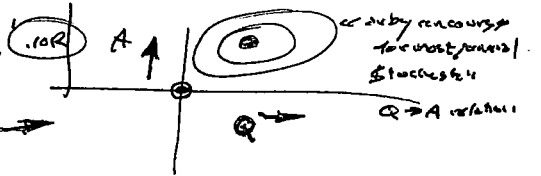
"Modifying t. Grammar" can allow anything (universal Grammars).

Essentially, this (Growing S. Grammar) represents t. wtd. set, [Oʲ] → stoch of operators → see 272.09-16

They T. can observe A's of t. grammar are functions that map Q to A.

→ So we have a P.D. on functs that map Q to A. Can this be universal? — i.e.

Can all S functs be represented this way?

Goal. S funct: for each Q we can have an arb. dist'n on all A's. .10R



Look at Q & A as x & y (2 dims) ——→

Now consider all possl. d functs relating Q to A.

Consider each contour of constant pc on .10R to be a functional relation

betw. Q & A. So ▦ 10R would then express a pd on that set of functions (= contours).

There is a diff. here, hvr: In the Graph of .10R, T. contours make A a 2-valued funct of Q — i.e. not really a "function".

More general, for a given Q₀ there may be >>1 A that have same pc.



T. foregoing diagram is about a single contour —— i.e.

I think many/different functs because Q & A can have t. same probty, so both f₁ & f₂ could have same pc.

So Maybe all stoch functs can be represented this way:

But (are universal) is it a good way? (Actually it is unclear that t. older way " represented by .10R, is particularly good (

BIG

Another freedom I haven't used in t. P.d. of d. functs betw Q & A, is addition of pcs by being on >1 d function! — (tho this may not be so useful from pt. of view of (stoch) !

[ Perhaps .26 is t. main Q: Are either methods Good for t. kinds of probs I'll want to be solving?

I have worked on (.09)(.08)(t. Universality Q) some time ago. (at ASTIA perhaps) — even both.)

One problem is: if I did have a d.f. over all d functs, what would it cost (recog)?

20 ∴ Note "t. Stack" as part of "context": It summarizes input parts of recent history

~~SAT I very prob wrote this Esp'l... recently... Inf QATM~~

03 : 268.32 : Unclear as to how serious this is! While "update" is ~~not~~ exactly a QA-induced problem, it can use cncs usd by QA.

35 ....┌─ **Formally** any sequencial predn. problem can be regarded as a single QA, 
.06 │ but that t. A's possl. need not be finite: ie. t. machine does n't have to stop (ever)
.07 └─ after printing t. part of A to be extrapolated.

It would **seem** that .05 **does** tie sequencial predn into QATM in t. useful way. **It** tells how t. cncs for QATM are linked to Rare oft. sequencial predn. problem. In t. more General Version of QATM, .06–.07 **is** t. way we evaluate t. individual $O^J(A_i | Q_i)$'s

So saying "Updating is just another QA problem" is a bit ~~deceptive~~, but only because ~~it needs~~ t. statement needs expansion to be properly understood. Anyway— for **me** its O.K., à I now understand some things that were prob'ly confusing me in t. ~~#~~ past.

Ⓞne of t. things is re: Jess' remark about updating being a regular problem of QATM! That this recursive idea **could** ever be meaningless nett. BUT, it **is** O.K. if Rare are other Q "A's in t. corpus — t. more there is, t. less chance for "error" due to "self definition" (or "self optimization") . — In general, t. result of this self optimization will depend much on t. A priori info .... but **much less** so, if there are a lot of other probs in t. corpus.

23 ~~Or~~┌─ It **is** easy to get a stochastic **QA** from ▬▬▬ Paul's Grammar of 6.7.90 Paul 14.03 ]!
    └─ We put in a Q : eg. ~~....~~ 3.7, +s ──── of 90 P14.03 –.12

The "machine" is t. current s. "**CFG**" that generates poems. So it gives a s. lang. as output, for any input, like. [ Actually, this is a peculiar "s-funct" we have an s grammar not gives p'd on strings (≡ opns) →

After we have gotten several poems ~~......~~ in t. tray, that solve problems (ie. Got hy pc's for "correct" A), we use those examples to modify t. grammar so it gives even more pc to t. corpus. We can do this by making new (tokens) defns) à modifying continuous params of both old à new ~~def~~ tokens'

→ At first those "successful examples" are given npm names, ~~........~~ à assoc. pc's. — ~~#~~ say $\alpha_1, \alpha_2, \alpha_3$ are 3 of those I solns, then ~~......~~

$\alpha_1$            $\alpha_2$
push push top, push push –op , ect.

We **augment** Paul's Grammar by $st \genfrac{}{}{0pt}{}{\to}{\to} \begin{matrix}\alpha_1\\\alpha_2\\pc3\end{matrix}$ , since this gives

each example in corpus a pc of ⅓, (which is >> t. amt of info in "3" + "7" (2 à 7 could 16 bit randomness (or 8 bit — for speed)).

~~Since~~ As t. corpus f (and/or our CB f) was continued to modify t. grammar to get hyer pc for t. corpus.

(margin right) of 90 P14.03 –.12 · Pd on output strings set poems we have t. Question of 270.08 –.09 ·An artifact understanding of this is 274.12 –.22 , 275.22 , –276.05

N:φ>

∞  : I've been running thru this a bit too rapidly, missing many points.

1) Num is not a terminal      , so is    cand → Num, +ₛ, -ₛ, ×ₛ, # goal.

   presumably one ~~se~~ could have operators that recognizes that some thing is a number.

2) we can have different ~~parts~~ of theory for strings, cands, ops, etc. : this would be reasonable.

3) §G.04 discusses "the final ANL" form

i've 67.25 is a different formulation. (73.20

4)  4.36 of ~~~ ANL  (of 8/84)

_____

[SN] Actually putting every ~~test~~ problem soln found in ~~~ "memory" (≡ Srtus) does not seem
like a good idea! Perhaps (at first) just solve a lot of problems, w.o. using "calls".
__Then__ look for regys in solns.  T. regys will (at first) be t. solns to t.3 ~~~
a'g. operators, w.o. "conditions".  (like ~~~ "push, push top")

So we solve all ~~those~~ problems, & in doing so, we "have all those little "pgm fragments" that solved t.
problems.  At first, we don't know which ~~~ fragment to use in a problem, but we __do__
have a proby distribu. over fragments.  (say 4 fragments).
T. "4 os grammar" has o= push, push is one word ~~tree~~ & ⟨ +op / -op / ×op ⟩ are t. other 3 "words".  "fragments"

Using those 4 words (& ~~~ & pf. on them) plus t. other possl. words, TM ~~~ discovers
p. "longer fragments" that have hgher ps.  like:

       & push, push if +ₛ then top else ⟨ -op / ×op ⟩ )    ← 2 pgms. But would ≥ of time

So _____↑ gets declined.
                   gets declined.

            push, push if +ₛ then top else if # -ₛ then -op else ×op . ← always works

_____

28  ↔ (NB)  In induction, t. functional forms can be (& usually are) order indif.
       So funct. is of form Π Oˢ(Aᵢ | Qᵢ).

∞       But approxns to Oʲ can's are dependent on order in which data is processed.

    ⟨⟨ ⟨So t. update problem is then sig. crit... diffrnt from t. "first order" ("exterval") ⟩ ⟩⟩
32       problems that TM works on!      → 269.00   (so not such a big deal!) ⟩ ⟩ !

       [SN] An alternate Approach to t. problem (in Lsrch) of many good cands being
       about same pc (of being __best__") is being hgly correlated.
           If we use for our "guiding" P.D., "t. pc that t. G of a cand will be > Go", this will
       be o.k. — but then t. pc's are not normzd in a useful way (i.e. Σ pcᵢ need not
       be bounded), so we don't get a useful CJS (which is t. ~~~ Goal (ck Lsrch)

: On looking thru SAARB notes for work on ANL: details of implementation
ā CJS's of solns.

1) SOP 96.30   Disc. w. Peter Bergman.
   " " 98.10   "    "    "
   (8/5/90) → SOP 96.25   Disc. w. w.  : ? ... up ANL w. CJS's, ~~?~~ list of defns,
   expected directions.   So maybe ANL of int'l was ~~?~~ !

2) OSCA notes.   SOP 63.33   & ibid: 103.19ff, 111.21ff

3) AH!        SOP 14.03 : ~~Fu~~ W's Grammar
              (Also ibid 44.25 (?)

2 or Mantissas of "roots"   * ibid. 11.25     8° : 80     15.30 (1280) | 16.09 (320) 17.26 (320)
                                      11.20   14.20, .25

21.29    50.625k   ( or only those 47.25k ? )   T. Prog. is a 17 symbol Grammar. ( ibid 14.12)
  .25    78.608k                                 The item "Call", its not clear how · subtasks are numbered.
  .05     3.136k                                 Has its storage successful pgms ā then accesses
                                                 them later.

Th. Grammar ibid $\overset{(SOP)}{\#}$ 14.03 –.12 seems o.k. ( but no ÷ — I don't see why : simply because
÷ φ gives trouble, is not enuf!   It could just cause "end of trial".

Anyway, T. first problem in t.seq is "∉ & no. "17" say :
t. soln. is " ∫ ð enum then push "   3 symbols.   re vals = 4,5, & :  rc = 80
                  1    2     3
This become   s·m #1 s·m

To do addition   first   st → st; st.           first st → call 1   }  4 × 4 × 80 = 1280 = rc
st →   [ x +₅ then +ₒₚ   so 80                  for addition.
         └─────┘
            +

Same rc for subtracting, mult.                  rc for subtraction is mult by 2 because 2 poss'l calls.
                                                 "       mult "   " " 3        3 "   "
So   1280 × 2   subtraction
     1280 × 3   multiplication

Notation is Rev. Polish.                        So say it got   7, 3, + as input :
It push 7, push 3, +ₒₚ.                          { call(#) 2
So total rc for add would be   4, 4, 80 =   1280
for subtraction  its                            1280 × 2 × 2   2 calls.   or, we may have
for mult   its                                  12 × 3 × 3               (push, push) as a single concept.
                                                                         much ψ in fac !
    3, 4,5, add, mult ·        I don't remember how I got to "general soln ∫"
         5

Perhaps main idea? " If Prog are in comparative notation & a k input function is symbolized, then
create associated a printer."

This is z dep of a dup!
of a dup    & 5

**dynamics associates**

Paul's Grammar w. defns:

.103    st → If cond then op else.

→ call ( number of gram defn)

→ st ; st

→ do st until cond {end do}

→ [ see 1721 for loop inst.]    loop : st : If cond End loop

cond → Num ; +s ; —s ; Xs | # :

op → Push | +op, —op, Xop,

Num → {0, 1}*    (finite (binary string))

.C2    Alphabet = $\Sigma$ = $\Sigma$ (If, then

No! (If, then, else) is one word!

call, do, until, Num, +s, —s, Xs, #

Looks like 14 symbols.

[If, call, Do, Num, +s, —s, Xs, #, Push, +op, —op, Xop, 0, 1]

Push, +op, —op, Xop, 0, 1]    17 symbols.
14 symbols

Perhaps don't need to list "14".

2 concat symbol ? .... prob;
Then I don't think it needs to
be on our op "17" : its "16"

$[0, 1]^*$

15 (until's unnecy)

(No division!) ← Actually no trouble adding division.
÷ 4.2 no problem.

If examines Input list.
Push nm n on puts. # object on input list at
Pointer position, onto stack, increments pointer index.

First two sep. soln. :    ●   &

If (Num) then Push. ( pc = 1/(4·5·4) = 80 from input RS?

This pgm ends up w. number onto
tos.    which is what we wanted.

If num then Push :    "Num" is not a terminal .: illegal.
But still, let's expand our idea of Grammar.
So "Num" is not a terminal, but a language —
( ≡ set of strings).

If +s then +op.    pc = 1/(96) = 80
(96)⁻¹
1/4    1/6    1/9    = 96⁻¹ —1
≡ α₊        80

DEF    If we want our grammar to have defns: we have a special symbol bd between
defns. & at p end of all defns..    : end of all defns is "ed"
between defs ; "bd".

If a grammar begins w. "ed" then there are no defns.
After the last defn, we don't use "bd", we use "ed".

Grammar → Def statement

Def →    st bd Def
→ statement
→ bd

bd
(st bd)ⁿ

Nœs



: Proxy best way to do t. conc. net. ≀ language ≀ vocab ≀ insts ≀
wr.to t.
TSQ: Then make conc. net.! But conc. net is, at first, rather vague:
It is expressive of a vague human mind.
   At this point, we can take several possl. paths!
1) Make conc. net less vague
2)  "    "    "    actually not vague, but exact.
3) If Ø try to derive lang that has enuf "looseness" in it to implement t. conc. net
as ▪ |S-functions!|   Then as time goes on, new concs are discovered ≀ so many of t.
S-funcs become closer to d-funcs!   This looks like a quite different approach
from what I've had in mind in t. past!  — In t. past, I workd on t. conc. net so
make it exact" ≀ then devised a lang/inst. seek to realize it.

      These are 2 (seemingly) essentially diffrnt approaches to
(wng Algebra!  One purely/deterministic, certain, exact. —
The Other very uncertain, fuzzy, probablistical t. wey.

   Either way 's possl. perhaps my Body of them!
   In t. deterministic path, I eventually have to honest learn to work
problems wi S-funcs as solns. (So Tdh unite 'get t. pt. where it
would decide which kinds of problems were likely to have d-solns,
V.s. which needed S-solns

                Making  /Ever
   Also, Remember manage/Conc. net ≀ making TSQ's to go from one pt.
to another (no really starting out "Bottom").
for each pt in t. net, write several Solns. — see which are most easily implemented.

N. PS

(SPAC 26340)

(3) Perhaps make comments on OOPS: (a) how "factor of 2000" is a funct of "73".

(β) How OOPS can't do probilistic induction, but can be easily amodified to do it.
→ Why we couldn't to do it

(γ) How to system can do Boom induction — (its clausib now perf. & prove)

(δ) "Boost" as a kind of Mutation

(ε) "Boost" seems difrnt from AZ141. Just how is it difrnt? consequences? —
in Genzr form is it "Good enuf"? i.e. is it universal? This may be ≡ a very old
Q that I successfully answered (e few times ( ). [we will invoke a 'priori']  ← what is this telling about?
value for p when sz twit, is φ — so "Modified Laplaces rule"

Main problems of TM:

→ 1) Good set of instsll  Good initial language (good set of concs)

  a) Good set of insts, concs.
  b) Good form of Loop (Perth?)

2) Goat way to represent s-functs" p.d. on Sfuncts.
  a) I have a 3 impt ume.
  b) Modify (Parts) (≠ OOPS) to get this.
  c) to Do AZ-type (Born) pc's. — How to do this on second level?
     I. e. P.D. over AZ-type langs.
              2

3) Early Search routines:
  a) Lsrch ( see ref. on NIPS 133 for a way to do search over TSQ for starts. )
  b) On "Continuation" v.s. "Backtrack": when & how much to Backtrack.
→ Design of Langs, arch, so that .26 is feasible.

4) Design of TSQ:
  a) Simpler Algebra: Chose lang (1) & Design such routines (2) so that Press pnts
canbe solved.

M/

5) Context: Used to modify pc's of tokens.

I think it would be desirable to do (4, 1, 3) together: interactively.
As it is, I think I can devise a concept net for 1 vnp. Algebra.
(No sum w. some elements of large < Ss).

(·12) $M_{2R}$ : $P \neq NP$ : undecidable

↓ GOOD!

SC

[SN] On "undatiny" in view of more data. Say I define z upon early in life. Later, t. corpus is → t. defn. is no longer a compressor.

Hvr., this data has been used as components of other definitions, so "undatiny" it would undefine several other defns that have, in view of t. larger (subsequent) corpus be useful compressors. It would be a BIG BACKTRACK

Well, t. thing to do is retain all defns, but update their pc's. There is a "chain rule" involved since "hi'er order" defns have PE's that are dependent on t. pc's of earlier defined upon.

Updating pc's involves retiring all defns but updating their pc's int. temporal order in which t. defns were defined. So one can do this for t. entire set of defns every once in a while: Th. top-level defns pc are updated most frequently

Updating in OOPS: Does it consist only of freezing old solns, to "acknowledgement" of tahoms that have been defined in soln of last problem?

→ Its really not clear in my mind as to just what old solns are accessable (as frozen info) to a new problem-to-be-solved.

17

[SN] [M2R:] It could very well be that t. P v.s. NP problem is undecidable

It does involve behaviour at n=∞ only. i.e. we are concerned w. behaviour for arbitrarily large values of n, only. Seems related to Chaitin's work on t. unsolvability of certain probs in Number Theory.

20

Whether a particular problem is P or not could be undecidable

→ Updating OOPS seems very easy (trivial). Why is Alpha update so hard? Its a matter of (definition!) [Terminology] in Alpha, we spend much time on update, but prediction is always very fast. In OOPS, we spend much time on predn. & there is ≠ ∅ update. The sum of t. update + predn. is what is of interest & they are both big in Alpha & OOPS. Alpha ≡ { Algorithmic Probability (Heuristic Algorithm) }☺

27: [262.12] → One way : Say we have z bunch of density vectors that are all desirable we may simply want to "OR" them together. If there is z subset of these t. vectors that are close together — then OR them. If we can partition the vector space int ///////// — do that and (OR each ) Cluster as a "boostp" vector ( Remember these vectors are BAG representations, so we have to be careful in OR'ing.)

CLUSTERS

on t. "center" of each

If assume P=NP (all probs w/in poly) then t. Q of whether a very prob was for NP could be solved : (e.g. it's always P). So it its impossible to tell if certain prob are for NP, then P could be true. ... then P ⊇ NP

i.e. ideally they are assoc w/ 7. problems that are "close to t. prob one wants to solve

Nier

30

: Report Revision:

1) Look at my letter (from) J: — (7.1.03)   After li by Go thru J's comments.
Several impt items   Also other Letters.   12) see 7/2/03 "things to do" lines 22 ff

2) Table of contents:

3) How to learn from failure.

4) How to design $O^j$ or functions.
2 ways: @ 3 input unc ) 272.24 & forg. discussn.
@ Grammar   to 276.09

5) I said t. long used by OOPs is a kind of
stack (not stochastic!)
long. — How ants me to emphasize

O

"t. long used by t Current OOPs (pilot)   I don't really
implementation is a kind of stack log.   see this!
But in fact, I can only wished comments on t. Paper of 13, not on
discussion future & expected future developments

6) Use of WON instead of elaborate
psych.

[  7) A section on what has been worked out   This is related to   RV: 277.03 –.24
   & what needs to be worked out   → 265.12 ff   ⟹ Essentially RV is
   & what needs to be worked out   12 V. Good   (7) + (8)

20

   8) A new introduction or to Abstract:   See 277.03 –.24
   Explaining what Mech L mg problem was:   + 253.27 –.40
   What were major sub-problems,
   What sols I'm using & why & think(they
   are V.G.

9) Have more detailed discussn of OOPs. in §6
   e.g. $\left(\frac{n_Q}{(1...)}\right)^4 \cdot \left(1 + \frac{1}{n_Q}\right)^3$ @ttack.
   Also give some of the corrections that I wrote
   J. about in my letters to him.

30

10) Also Note OOPs doesn't solve
   s funct! (redu problems — only d funct)

11) Main differences betw Phase 1 & OOPs. ——
a) Lispish v.s. Fortrish langs.
b) Ph. 1 tries to do(lvn) s funct; OOPs does only d funct
implement → c) Ph 1 trys a tsof of essish problems. OOPs staves on hard probs.
philosophy
of research — This may be political, rather than scientific reasons.
d) Ph 1 uses Ganzd contexd to deal w/scaling: How OOPs goals w. scaling is unclear   –726500

00 : 261.40 :  Actually 261.36 differs much from what OOPS did, — what boostq did. Boostq picks a
single ~~the~~ problem-ate rime (plus/background density [uniform]) —— Too one could use 2 or
more "Boostq's", t. result would be of much lower priof.

The previous discussn (see 261.18 on "Dot products") reveals a serious difty in
~ Lsrch: From t. dot products, we will get a lot of ~~smritoon~~ density vectors
that are promising, but similar to eachother. I would like to use them in ‖,
by adding together t. pc's of various ideas occuring (jointly)          t. "WOW" approach
Unfortly Lsrch doesn't do parallel addition. Instead, one way is to choose t. single
best looking density vector. If it doesn't work after a while, chose another
10        t. hy pc density vector that is 'not too hyly correlated w. t. first,
continue trials ~~even~~ for a while. Of unsuccessful, try density vector
.12   ~~t.~~ relatively uncorrelated w. earlier trials, etc. ...          → 264.27

This problem was dealt w. in ordinary Lsrch (using WOW in ‖ updating
system. So we chose "best looking" PST. As we work on it, unsuccessfully
its rank as "best" ~~best~~ decays: As does other PST's that are ~ to it.  "highly correlated"
when (These) one's ~~t~~ get to be no longer #1, we jump to t. new #1 ... which
will be a PST that will tend to be dissimilar ~~to t.~~ from first trial (i. because PST's that
20   were ~ to t. first trial will have their probtys of being best" t.).

In the Won system, the probty of being best will be   $\frac{(0^{th} \text{moment})^2}{1^{st} \text{moment}}$

The moment being of t. ‖P  ⌠∿⌡   function for t. (PST, problem) pair.
                                      →  T

──────  ─────  ─────  ─────  ─────                a soln of t. of Hanoi.
26 : 258.28      Note on $\left(\frac{K}{11}\right)^4$ formula:   the value of pc/ for search

We. ~~are~~ promises trig !

$$\left(\frac{1+k}{7 \cdot h}\right)^3 \left(\frac{1}{7 \cdot h}\right)^4 = \frac{1}{7^3} \cdot (1+\tfrac{1}{h})^3 \cdot \frac{1}{7^4} \cdot \frac{1}{k^4} = \frac{1}{7^{-7}} \cdot k^{-4}$$

30   After trig  $\left(\frac{1+h}{7 \cdot h}\right)^3 \left(\frac{1+k}{12 \cdot h}\right)^7$  = $\frac{1}{7^3} \cdot \frac{1}{12^7} \cdot (1+\tfrac{1}{k})^{10}$        $= \frac{1}{7^7} \cdot \frac{1}{k^4} \cdot (1+\tfrac{1}{k})^3$

                                                                          $7^3 \cdot 12^7 = 1.229 \times 10^{10}$

So   $\frac{w. \text{trig}}{w_0. \text{trig}}$ : $(1+\tfrac{1}{k})^7 \frac{7^4}{12^7} \cdot k^4 = (1+\tfrac{1}{k})^7 \left(\frac{h}{11.053}\right)^4$

So, pc of soln. after ~~boost~~ is ≈ indep of $n_\phi$,
But pc of soln w.o. boost is ∝ $\left(\frac{h\alpha}{11}\right)^4$

7/12/07

NIPS

o: 258.90 : So three problems:

1) What is exact justifin & mechanics of "context" dependant token pc's?

o2 In particular, for various kinds ("levels") of context, just exactly what is t. corpus? → 260.00

2) In Oops: → 1. (Q) about "Boosty". a) Boost looks like a particular kind of "context" — is it? If so it gives a joint pd on successive new tokens — so it is difrnt from any "context" — Can I genz. my "context" to include this "joint pd"?

b) A general Q. about "instructions that modify pc's of tokens". J. says that in general, insts can compute pc of next token". Just how legit. is this? In t. case of boosty, it seems to have some a prior. justifn., but for arby insts that change pc's of tokens .... unclear how reasonable this is. ———— He may have called them "patterns" patterns

I think this was thinking of various density vectors" A density vector is a set of pc's for all of t. tokens. One can combine them in various ways ① linear combination ② multiply 2 vectors together (& perhaps normalize). thus t. "Boost" density vector set was t. only one he included in OOPS. Presumably for various problem domains, one could have different "density vectors". — they give a kind of a priyd for each/domain.

3) There was this Q of how to generate/formalize/derive a universal pd. over s-functions: AZ does it for d functions (in a "manner of speaking"). for univl. d.f. on s.functs. (Solution) The formalism: any 3 input unc defines an s.struct. t. first input derbs t. function. t. next input is for t. inputs, t. third is the "R" (random) input. — As we have normal unc machine derived prcby – types of output. So, random input on input #1 gives a univl. d.f. on s.structs. t. longer t. prgm, t. less pc. Each input is a self limiting input that decides when t. should stop & go onto t. next. input. There is a ref. to a detailed Lynch over such uncs in t. NPS (133,—)

Another way to derb an impt. set of s-functs, is t. Bern. formalism: Simply derb each object & tell what its pc is. To obtain s.functs of this kind! One way is formal langs. Bernoulli & finite-state (= HMM) langs, stochastic CF grammars, stochastic context sensitive grammars .... any other kinds of stoch grammars, Also BBN's (Bayesian Belief Nets) —(257.00)

So (3) doesn't look like any theoretical problem.

As is A2 is 2 units/d.f. over d.facts. **used** **without 2 units/d.f. over 5 thgs/hr.**

Seems to me that I solved this problem with (task pr one). It was lopsided ≈ 2

Scientus Regarded problems as + Prest time. Also much unsuccessful work to be far Prest.

As related to OOPS! It doesn't look like 4.15 ≈ 2 system is hard universal w.r.t. modifying my 15 orig 6

(SIU) Benefit of ≈ 1000 from "OOPS": "Boost"; however, technology depends on

$q_a(z=73)$, no. of tokens & raised w.r.t.

**.30**

$q_n ≈ k$:

$$\alpha = \left(\frac{1+k}{7 \cdot k}\right)^3 , \quad \left(\frac{1+k}{12 \cdot k}\right)^7 = 5.322 \times 10^{-11} \quad \text{for } k = 73$$

$$\beta = \left(\frac{1+k}{7 \cdot k}\right)^3 \quad \left(\frac{1+k}{7 \cdot k}\right)^4 = 4.539 \times 10^{-14}$$

$$\sqrt{z+6} = \frac{\alpha}{\beta} = \left(\frac{1+k}{12k}\right)^7 \cdot (7 \cdot k)^4 = \frac{(1+k)^7}{(12k)^7} \cdot 7^4 \cdot k^4$$

$$z \cdot k, \, k >> 1 \approx \frac{k^7}{12^7 \cdot k^7} \cdot 7^4 \cdot k^4 \approx \frac{7^4}{12^7} \cdot k^4$$

$$\frac{z k}{z} \approx \frac{k^4}{(1200)}$$

$$\text{So} \quad \text{gain} = \frac{k^4}{1200} \quad \text{for } k=11 \quad \text{otherwise} = \frac{k^4}{1200}$$

$$e^{\frac{z}{k}} \approx 1 + \frac{3}{k} \quad e^{\frac{3}{k}} = 14 \frac{3}{k} + \frac{1}{2} \frac{9}{k^2} + \frac{27}{6k^3}$$

$$\left(\frac{k}{11}\right)^4 \text{ should be mult. by } (1+\frac{3}{k})^3 \approx e^{\frac{9}{k}}$$

$$\text{So } z \left(\frac{k}{11}\right)^4 \cdot e^{\frac{9}{k}} \approx \left(\frac{k}{11}\right)^4 \cdot (1+\frac{3}{k})$$

$$= \left(\frac{k}{11}\right)^4 \cdot z \approx 2000 \text{ for } k = 73.$$

Now exactly

$$\boxed{\sqrt{\frac{116}{11.0527}}}^4 \cdot (1+\frac{3}{k})^3 \quad \boxed{2621126}$$

$$(1+\frac{1}{11})^3 = 1.298 \approx 1.3$$

**Perhaps it would be best to think about** .00 ... before considering .05?

In .05 OOPS uses "permitting lists" — they modify pc's of new-tokens to be or-tokens

Boost is like my "context" dependent pc's of newly or-tokens (to b or-tokens).

In the my "context" it ... 12 ... regular regularly Prest & Prest of a certain corpus

[I will dig out just identify compress/user]

thus Boost is a bit distinct from context in t used/cause in Prest Boost fills a gap of

... ... to work together — to create a kind of "Auxgram" (... via a ... of Lattices)

of a previous soln.

N.PS

"Learning English". [(Naive → Real) Physics]

I had Idea of TM learning Algebra, then learn to discuss Algebra in simplified English. To graduate to more complex a/o less "simple grammar" English — Several ways: ① stick w. "Algebra World"

② Get TM to learn Simple "Physics" by learning to play "ping pong". This last can be in graduated degrees of sophistication. To start off, "balls" are points & no air resistance. Bouncing is perfect, There is no net. No "opponent" — only a "reflecting" well. TM is given exact coords of ball at all times. Then difficulties are slowly introduced. Air resistance, Net, balls of >0 diam, eventually, an opponents of gradually increasing skill & speed.

Eventually, I may want it to learn to play a game w. Real opponent — Modes ~ to TM & get TM to model the opponent & realize that the opponent is ~ to itself.

In general it is probly possl. to teach TM to play various "games" that become more & more like R.W. — Including Entities that are like people, So it would be then useful to Get TM to understand English discussion about "T. Games". There may be actual/video games (existing) that TM could learn to play.

A Quite different approach is used in Statistical Machine Xltn (which they do as "Q.A. long". To Some extent "Language understanding" is key to xlt between internal & External (Nat lang). lang. (See D:\PS\ Mach Translation-statistical for ) many papers on this.

~~Probablc~~ ~~Boolean~~ | Bayesian Nets ( Bayesian ~~Belief~~ Nets ) | BBNs

! What ~~Probablx~~ Bayesian ~~nets~~ nets may be!

Say we have discretes, finite alphabets of / variable values: poss.

A ~~module~~ [I,O, module] will have as inputs, a certain variables. For that module, this will induce a pd. on each of its outputs. For each input config. ( n binary inputs → $2^n$ input configs) it will have a pd. ~~associated~~ on its set of output configs ( if there are m binary outputs $2^m$ configs) — So, for each of the $2^n$ input configs, it will have a list of $2^m$ pc's (that sum to 1) for each poss. output config.

≡ Some simpl'ens:

1) only one binary output  2) a single ~~radix~~ radix r output.

3) .07-.08 are a single, radix $2^{m-1}$ output.

4) several outputs, but pd's of them are indep. (equiv. to several radix modules in || ).

inputs:
same as 1, 2, 3 for outputs:

relations of Input to output: 1) most general is in form of a list of probs.

2) I can't think of other common forms

I, O. modules can form a net by connecting ~~the~~ outputs of 1 or more modules to inputs of one or more modules. This usually results in many intermediate & final values of variables being correlated.

So: This gives a way to get fairly complex ~~the~~ ~~discrete~~ s functs for discrete variables.

For continuous Vars, one could also define "Modules" w. defined I-O. characteristics, but it would be more diffc. (e.g. in Discrete case, it is known that a complete (often realizable) descn of a probblstc I/o ~~behavior~~ of a modules is ~~always~~ always poss.)

In the discrete cases, perhaps simplified ~~the~~ modules & nets of modules are being devised so that their behaviors can be easily analysed!
In particular: How to (usefully) constrnt. no. & nature of params defining a module or net, so that learning can occur w. a reasonable ~~eff~~. For most nets, modules, data types, → w. reasonable cc. )
the problem a may be ~~practically~~ ~~practically~~ ousolvable

Three kinds of Probabilistic Induction: Universal distributions and convergence Theorems.

t : 251.40        section **4** :   ~~Incompatability~~ of Universal Distributions.

Abstract:  We will describe three kinds of Probabilistic induction, ~~~~ give ~~state~~ general solutions to each kind (~~and associated~~ with) convergence Theorems to show they give good probability estimates. tend to

The first kind ~~~~ extrapolates a sequence of strings and/or numbers.

The second ~~~~ extrapolates an unordered set of strings and/or numbers.

The Third. extrapolates an unordered set of ~~~~ ordered pairs of elements ~~~~ that may be strings and or numbers.  Given the first part of a pair, to get a probability distribution over the corresponding second part of the pair.   → .25

(1)  ~~For each solution a correct reading an convergence Theorem is given, showing that as sample size grows, the error in probability estimate decreases rapidly.~~

~~These problems. These kinds of induction problems are very general and~~

17 : .29    free  →  The solutions given are very general and/color may ~~perhaps~~ all kinds of induction problems.   Time series prediction, Grammar discovery (for formal or natural languages), curve fitting, The Identification problem/ and the categorization problem, are a few of the kinds of problems amenable to the methods described.  End.

0                                 three kinds of
~~For each of the problems, a universal distribution is given
that solves it, and a corresponding~~

29 : .11    Each of the three kinds of problems is solved using an associated universal distribution.  In each case a corresponding convergence Theorem is given, ~~~~, showing that as sample size grows, the expected error in probability estimate ~~decreases~~ rapidly.   → .17
                     decreases

29

0

so:   : Inserts list. from ~249.12

~249.12   ①   Dimacs 1.22-.28;   1.29-.32   ) do as one e branch 1.22-.32

—————— DM1.tex       DM2.Tex

250.19   ②   Sol §§   25G 27-.37     S99 S99=1.Top

247.30   ③   ( IDSIA Report. )

            BS eq(1) to 5.27     SO2=1.tex

            Appendix B completr.     SO2-2.tex

og : 247.40   If we set $Q_i = \Lambda$ $(i=1...n)$ in equ ( $\overset{B}{\cancel{B}}$ )

o:   It becomes clear that the equation ( $\cancel{B}$ ) for induction on unordered sets is a special case of Operator induction, and that the Convergence Theorem ( B ) holds for unordered eq (A) as well. This also assures convergence of the Operator induction technique of Section 2.1.

     Is there any advantage in using equ ( $\cancel{B}$ ) rather than eq (β) for operator induction?

     eq( $\cancel{B}$ ) exploits regularities in the set $\{[Q_i, A_i]\}_n$. eq (A) exploits regularities in the as a function of $Q_i$.   It includes regularities in the so would seem that we are doing more work than is necessary. set $[A_i]$ — which we discard.   In eq (B) we only find regularities in the functions relating $A_i$ to $B_i$. Superficially, this such regularities may be easier to find than regularities in the more complex object $\{Q_i, A_i\}$. In general, however, the convergence of either of the techniques will depend critically on just what problem is being solved.

The result is that the probability errors for the normalized measure, $P_M'$ can converge much more rapidly than those for the semimeasure, $P_M$.

Gács (Gác ibid) also shows that the corollary corresponding to q. eqs 5 and 6 holds if $P_z(\cdot)$ is an unnormalized semimeasure.

Marcus Hutter ( Hut ) shows that these results hold if we use alphabets with any finite number of symbols.

In the forgoing convergence theorems, the total squared probability differences is used as loss function. The proofs of the theorems also show that the same convergence for the kullback-Lieber loss function (which is greater than or equal to the square loss function).

Hutter (ibid) considers more general loss functions and shows that the use of the universal distribution gives losses that converge rapidly toward the smallest values that they could have.

---

Section 2    Induction on Unordered Sets

Subsection 2.1    The problem and a Solution.

> f.N. section 2.1 follows the discussion of fol 99 pp 256 and 258

Copy fol pg 256.27 - .37

A solution using a universal distribution is obtained by assuming that the data was generated by some unknown stochastic distribution on all possible single finite strings, $\mu(D_n)$

The Universal distribution is a weighted sum of all finitely describable probability measures and semimeasures unordered sets of length on finite strings.

$$P_M([D_n]) = \sum_j \alpha_j \prod_{n=1}^{h} P_j(D_n)$$

(8) eq

$h$ is the number of strings in the set $[D_n]$

$\alpha_j$ is the weight of the $j^{th}$ probability distribution on finite strings

$\alpha_j = 2^{-|a_j|}$, where $a_j$ is the shortest description of $P_j(\cdot)$ and $|a_j|$ is the number of bits in $a_j$

The M index of $P_M$ indicates that the functions $P_j$ are to be described with reference to machine, M. Since M is universal, it can be used to describe any describable function.

The probability assigned by $\mu$ to $[D_n]$ is

$$\mu([D_n]) = \prod^{h} \mu(D_n).$$

(9) eqn

≥8    To start off, we will normalize $P_M$ to create $P_M'$:    /

245.1 ⊢ .15    $P_M'$  ...

maximizes:

~~The particular Method of normalize~~ It is not difficult to show that this method of normalization ~~plus~~ for all ~~ratio for~~ $P_M'(x)/P_M(x)$ ~~Base Map~~ for all x. it will become clear. Later ~~will gives~~ that this condition also ~~gives least expected~~ ~~$P_M'$~~ ~~error in probability estimate~~ a kind of minimal error in prediction. leads us to expect

0    minimal errors in $P_M'$ to have minimal prediction errors.    gives least expected prediction error for $P_M'$.

12 : From ~~~~ Printout ~~2.40~~    Just how accurate are the predictions of $P_M'$?

on    | DIMACS 1.22 - .28 |    but modify by using $M \to P_M' \cdot$ , $P_0 \to P_0'$.

on    $\underset{M}{E} \sum_{m=1}^{\infty} P_M'\big( z_{m+1} = 1 \,\big|\, z_1, z_2 \cdots z_m \big) - \mu\big( z_{m+1} = 1 \big| z_1, z_2 \cdots z_m \big) )^2$    $< -\frac{1}{2} \sum \ln P_0'$    (4)

on    | DIMACS 1.28 - .32 | ~~over 2 ?~~

The truth of eq (4) hinges on the fact that if $\mu$ is a computable probability measure ~~~~ then there exists a positive constant $P_0'$ such that

~~where~~    $\dfrac{P_M'(x)}{\mu(x)} > P_0'$

and that while $P_0'$ will depend on $\mu(\cdot)$ and $P_M'(\cdot)$, it will be independent of x.

Eq (4) can be usefully generalized so that if

• $P_1$ and $P_2$ are any normalized measures on x
- $x(n)$ is a string of length n.
• $\dfrac{P_2(x(n))}{P_1(x(n))} > \alpha(n) > 0$    (5)

where $\alpha(n)$ is a function of $P_1(\cdot)$, ~~~~ $P_2(\cdot)$ and n, but not of x

·0    Then $\underset{P_1}{E} \sum_{m=1}^{n} P_2 - P_1 )^2 < -\frac{1}{2} (n \,\, \alpha(n) )$    (6)

~~The proof of the convergence theorem of eq (4) given in SoL E 78 is for the~~ The Convergence Theorem of eq (4) ~~is true~~ if $P_M'$ is a ~~un~~normalized universal measure Peter Gács ( Gác 97 ) has shown it to be true for the unnormalized semimeasures $P_M$, but the associated convergence constant $-\frac{1}{2} \ln P_0$ is much larger than the corresponding constant, $- \ln P_0'$ for $P_M'$. ~~The difference is~~ The difference between them is

$$-\frac{1}{2} \ln \frac{P_M'}{P_M}$$    / ~~~~ value of the normalization ~~const~~ factor for very large n. $\frac{P_M'}{P_M}$ is the we have selected a normalization technique to make it as large as possible.    ≥250.00

Nips

00: 247.40 :, stuff to l. bottom of 241 ⌠ item 241.36  is perhaps readable: 242.15: review

P1 : 241.36 → We are using a/normalized ~~stuff~~ universal distribution, $P'_M(X)$. Minimal normalization

conditions are $P'_M(\Lambda) = 1$ ; $P'_M(X0) + P'_M(X1) = P'_M(X)$

There are many normal(rezating techniques that satisfy these constraints. (An additional constraint

~~leads to~~ so minimal/expected probability error is that the probability ratios of the normalized and unnormalized

distributions remain the same:

$$P'_M(X0)/P'_M(X1) = P_M(X0)/P_M(X1).$$

This gives the recursion relations:

$$P'_M(X0) = P'_M(X) \frac{P_M(X0)}{P_M(X0) + P_M(X1)} \quad ; \quad P'_M(X1) = P'_M(X) \frac{P_M(X1)}{P_M()}$$

11 $\quad P'_M(X0) = \frac{P'_M(X)}{P_M(X0) + P_M(X1)} \cdot P_M(X0) \qquad P'_M(X1) = \frac{P'_M(X)}{P_M(X0) + P_M(X1)} \cdot P_M(X1).$

15 $\qquad$ with the critical condition $P'_M(\Lambda) = 1$. $\qquad\qquad$ →242.00

To reduce expected positive error probability, we will be using a particular

normalized ~~version~~ of the universal distribution.

.17 $\qquad$ Because certain of the codes, $S_k$ do not result in useful output (i.e.

the computer prints out part of X, but continues to calculate without printing any (thing else.)

the resultant probability distribution is not a true measure, but a ~~semimeasure~~ semimeasure — usually

$P_M(X0) + P_M(X1) \leq P_M(X).$ $\qquad$ 2 italics

~~It can be shown that~~

We will normalize $P_M$ to $P'_M$ so that

23 $\qquad P'_M(X0) + P'_M(X1) = P_M(X).$

The additional constraint, $P'_M(\Lambda) = 1/$ assures us that ~~the probabilities of all~~ the probabilities of all

strings of a given length sum to one.

Later, it will become clear that the larger $P'_M$ is, the less expected error it

has. There are many normalized methods that satisfy these two constraints.

To obtain as large $P'_M$ as possible, we add the constraint.

30 $\qquad$ There are at least 2 ways to use (1) for prediction!

$$P(X0|X) = P_M(X0)/P_M(X) ; \quad P(X0) = P_M(X1)/P_M(X) \qquad (3)$$

$$P(X1|X) = P_M(X1)/(P_M(X0) + P_M(X1)) \quad P(X0) = P_M(X1)/(P_M(X0) + P_M(X1)) \qquad (4)$$

eq(3)

250.40
20: 245.40

· Suppose that [Dₙ] n=1...h is a set of/strings generated by some unknown stochastic device. What is the probability that [xxxxxxxx] our universal distribution [xxxx] assigns to $D_{h+1}$, a new [xxxxx] (possibly) new string?

It is just $P_u([D_n] \cup D_{h+1}) / P_u([D_n])$ .                                (10) eq.

omit ⌐Any by function that can assign a probability to any finite string⌐can also be used to assign [xxxxx] probability to each bit of [xxxx] a string, conditional on the preceding bits of that string.

~~How accurate are these probabilities?~~
~~These probabilities~~ For a suitable ~~[xxxxxxxx]~~ set of strings, [Dₙ] these probabilities ~~can~~ can be very close to those assigned by μ, the true generator of [Dₙ]

In section 3, we will discuss Operator Induction and prove [xxxxxx] an associated convergence theorem [xxxx] that implies a convergence theorem for Induction on unordered sets.

Section 3    Operator Induction.

In the Operator Induction problem, [xxx as described] in the [xxxx] introduction, we are given an unordered set of/ⁿ strings and a number pairs, [Qᵢ, Aᵢ] [xxxxx]    241.36 → [xxx]
                                                                                            243.[x]
Given a new $Q_{n+1}$, what is the probability distribution over all possible Aₙ₊₁?           243.15 ~ [xxx] 15,
                                                                                            243.20
We will ~~try describe~~ two give two solutions. [xxx] The first [In the first, we consider ]243.20 → 243.30
this to be an extrapolation of unordered finite strings, $D_i$ [and] $D_i = Q_i, A_i$.      [xxxx]
                                                                                            [xxxx]
~~[xxxxxxxx]~~                                                                               242.40 → 243.25
eq. 8 is used to obtain a probability distribution on all unordered                          243.25 ~ 243.40
sets of $Q_i, A_i$ pairs and eq(10) gives us a probability distribution over                243.40 → 245.00
[xxxxx] $(Q_{n+1}, A_{xx})$, — i.e. $P(Q_{n+1}, A_{xx})$ for all possible Aᵢ .              245.40 → 247.00

Then $P(\cancel{xx} A_{xx}) = P(Q_{n+1}, \cancel{xx}) / \sum_i P(Q_{n+1}, A_i)$            (11) eq.

Section 3.2
The second solution to the Operator problem is
[p5 of [report] (eq (8) of report) to 15.27 | Then jump to Appendix B. )    This includes "A if x B should   Later 3 or
                                                                            Modify Pts. to "we will show   ← the unclean lines
A23 (Appx B) of report. then through of appx B. Some Modifns.               these".                         I will fill in each.

I can take care of this part myself.

[xx]    Follow Pts in discussion of how Bay induction is special case of
so since our first soln of Op induction problem uses [xxx] induction —
we can be certain that it, too, converges rapidly

251.09
248.00

∞ :244.40: Actually, we really chose cands on basis of max $f_n(x_{(u)})$ — for $x_{(u)}$ being t. curr.sc.

So this $f_n(x_{(u)})$ is perhaps an approxm. of $\sum_u f_n(x_{(u)})$.    ( i.e. $x_{(u)}$ are assume is a "typical" u output.)

or maybe $\sum_u \ln f_n(x_{(u)})$.

.03
> So after all is said & done we select $f_n$'s s. $P(f_n) \cdot f_n(x_{(u)})$ is large. a sum w. $P(f_n)$ wts.
> Simple old Bayes; But we realize that $P(f_n)$ is impot. & it's much better if
.05
> $P(f_n)$ is truely & carefully Updated !

So, indep of all t. conversion stuff, we still use same methods to estimate error

( The ALP does enable us to do away w. t. training set )

10   Y. Cost of eliminating t. Trg set is that we have to consider many more models —

So their a priori s ↓ !  So we may not be getting such a BIG Bayem !

_____

But t. main idea is that all t. careful analysis about accuracy, etc. is not

so important as $\boxed{.03 \to .05}$:  Just to use good standard Bayes w. good updated a priori.

A possl. twin from ALP is to try to uncover a p'p'd "universe!" — That

one can get ideas on how to get t. A p'p'd from the existing

~~Scientific~~ / surveys of t. Domain of inquiry.


20  $\boxed{SN}$  On things that converge "with probability one" : Unforty, this tells no thing about

How fast they converge!  Such theorems are often based on other theorems about

"Converging w. PC = 1"

23    It would be possl. to put many theorems in forms giving rates of convergence,

if one had a large enuf BASIS of Thms on quantitative Convergence Rates.

Unfortely, many theorems in statistics are of "PC=1" character — this is like

in Biology, dis covering  "A influences B" gets a nobel prize, but finding out

just how much "A influences B" is usualy not looked into at all — certainly no

nobel prize....  Quantitative Consideration is usualy regarded as (back work —

" not creative".

    Sn.23 At this work would involve devising a large enuf basis of

terms on Convergence rates.  One would also have to know lots of inequalities".

Maybe a good PHD thesis for someone! — Or maybe a much bigger job!

O: 243.40

In the forgoing ~~Convergence~~ convergence Theorems, we have used the total squared of probability error as a loss function to be minimized →

Hutter has shown that the universal distributions also converge rapidly when more general loss functions are used ~~(Hutt~~ → } Optimality of Universal Bayesian Sequence prediction for General loss and Alphabet

I could expand this — explaining ~~what happened~~. Just what he proved.

Tech Report IDSIA - 02-02 18 Feb 2002. Give web site.

$P_{ij}$ also has part But is different

or :. General loss bounds for Universal ~~Sequence Prediction~~ Sequence prediction

Tech report IDSIA 03-01 10 April 2001

2 useful for both Hutter references

Hutter considers more general loss functions and shows that use of the universal distribution gives ~~a~~ loss functions that converges rapidly toward ~~the~~ their smallest possible value. (Hut )

perhaps a footnote.          perhaps refer to Section "2.1"

Sol 99 256.27 - .37

### Section 2    Induction on Unordered sets

~~In this~~ ~~we will~~ ~~the beginning of this~~ [2.1] section follows the discussion of Sol 99 pp 256 and 258

Suppose we have an unordered set of n finite strings of symbols. · · · · · [2,3 ~~pp~~] refs. PP 256-251 ommit(?)

A solution using a universal distribution is obtained by assuming ~~entire~~ ~~set of~~ ~~probability distributions~~ ~~on~~ ~~strings~~ that the data was generated on all possible finite strings,

$P_j(D_n)$

by some optimum, ~~pn~~, finitely describable probability distribution ~~M~~ on unordered sets of strings, $P_M([D_n])$

$\mu(D_n)$.

$P_M(D_n)$ — limited

The universal distribution is a weighted sum of all finitely describable probability measures and semimeasures. [2.3]

9 [9.4]

$$P_M([D_n]) = \sum_j \alpha_j \prod_{n \geq 1} P_j(D_n)$$

h is the number of strings in the set $[D_n]$

$\alpha_j$ is the ~~weight of the~~ WEIGHT jth probability ~~measure~~ distribution on finite strings.

$\alpha_j = 2^{-|a_j|}$, where $a_j$ is ~~the shortest~~ the shortest description of $P_j(\cdot)$ and $|a_j|$ is the length of that description.

The M index of $P_M$ indicates that the ~~description language~~ functions are to be described with respect to Machine, M.

The probability assigned by $\mu$ to the set of strings, $[D_n]$ is [ Since M is a universal machine, it can describe any finitely can describe describable function.

$$\mu([D_n]) = \prod_{n \geq 1} \mu(D_n).$$

Maybe not "include"

[ Since $P_M$ includes $\mu$ as one of its component probability distributions

$$P_M([D_n]) \geq \text{~~} \quad i.e. = W_\mu \; \mu([D_n]).$$

How ~~~~ $W_\mu$ is the ~~universal~~ weight of $\mu$ in $P_M$.

Since we don't know $\mu$, we don't know $W_\mu$. ]

Nips

:[SN] On radix ≠2 for ALP's Conv. Theorem : For ALP, there is no diffty in formulating t. distribn.
For t. conv. Theorem, showing its true for t. KL distance is easy via G&C's proof. Showing

$\sum (coord)^2 < $ KL dis for radix >2 is difft.

Hmm one could code radix 3 as 00, 01, 10, with 11 never used. We then just have a binary
formulation of ALP, (possibly removing 11, when it occurs & renormalizing). corresponding, 00, 10/01.
I suspect this would give the rite universal dist. Could one use the Conv. theorem for
binary strings to show that these ternary pc's also converge properly?

a b ge are the 2 bits for a pair coding 0,1,2'.         α β δ are the 3 ternary values.
                                                        00 01 10

$\frac{a b}{\text{?}}$  $p(a=0) = p_\alpha + p_\beta$.          $p'_\alpha, p'_\beta$ & $p'_\delta$ are corresponding pc's for $\frac{m}{}$
           $p(a=1) = p_\delta$
           $p(b=0) = p_\alpha + p_\delta$
           $\beta(b=1) = p_\beta$

for t. ternary error      we know, $p_{err} \cong (p_\alpha - p'_\alpha + p_\beta - p'_\beta)^2 + (p_\delta - p'_\delta)^2 + (p_\alpha - p'_\alpha + p_\beta - p'_\beta)^2 + (p_\beta - p'_\beta)^2$

converges:     Q: does $(p_\alpha - p'_\alpha)^2 + (p_\beta - p'_\beta)^2 + (p_\delta - p'_\delta)^2$ converge?

<21  $p_\alpha - p'_\alpha \quad \Delta\alpha$       is  $(\Delta\alpha + \Delta\beta)^2 + (\Delta\alpha + \Delta\delta)^2 \geq (\Delta\alpha)^2$?

$2\Delta\alpha^2 + \Delta\beta^2 + \Delta\delta^2 + 2\Delta\alpha\Delta\beta + 2\Delta\alpha\Delta\delta \geq (\Delta\alpha)^2$         $\Delta\alpha \to \alpha!$
                                                                                                 $-1 \leq \alpha \leq +1$

is  $\alpha^2 + \beta^2 + \delta^2 + 2\alpha\beta + 2\alpha\delta \geq 0.?$

$(\alpha + \beta)^2 + \delta^2 + 2\alpha\delta$        I think it may be any way, but either way its easy!

is  $(\alpha + \beta)^2 + (\alpha + \delta)^2 \overset{\geq}{\underset{\text{?}}{}} \alpha^2$.        No. say $\alpha=1$, $\beta=-1$, $\gamma=-1$

       0    +    0   ... $1^2$.

$(\alpha+\beta)^2 + (\alpha+\gamma)^2$  v.s.  $\alpha^2$

So it looks like it may have to be proved by showing .02 is true: which is difft.
Hutter says he did it but its difft: I haven't checked his proof, but I have shown
it to be very likely for radix 3 (& maybe 4) by Montecarlo trials.

[SN] Say we are using a set of funcs (M) to try to find t. best one to predict our data.
μ t. true generator, is not in (M).  ——  But for each value of n (t. no of bits in t. data we're using)
There is a function $f_n$ & $\frac{f_n(x(n))}{\mu(x(n))}$ ← pc assigning to is max. [←data of length n]

(Even if μ is in (M), it may have a small wt., so that fn in this case will get more wt. than
t. M that is within (M).)

Hutter treats t. case of μ not in (M) in one of his papers (maybe t. 4th pseudo!)
He considers the funct in (M) that has the closest KL distance from t. true μ
(or something like that).
→ In my formulation, (M) will vary w. T, t. CB. T. no value of n may vary (to extent T)
over different funcs. in (M) — causing trouble. Hvr, if n varies betw. cands, — we can't
compare them very well (but we can approximately compare them).

→ T. Davison in 578 § VI on Cover's measure, is easy to read. — T. Reasons/of interest Bravo are easy to understand —

≥ Lampos

∞ : 242 : 90  [SN]  In **Sol 78**  I should Plot Cover's Extension Complexity (differed from my

Norm

Method of Warren by a factor that → ∞ in → n, (but very slowly!)

Since Extension probty is based on a countable no. of continns of the corpus,

à BLP is based on an uncountable no. of continns..... this **may** explain it ....

but look into this !  (Essentially, Cover used a semimeasure, like Gacs version of Convergence thm.

(SN) Perhaps include Gacs of "sequential predn via Bayes." à explanatory
Gene?"

⊃ : 242.40  → The proof

∼̶N̶∼̶E̶P̶(̶Z̶)̶ The proof of eq(2) in Sol 78  is for ∼ 2 ∼

has been

commit

$P_M$ that **is** normalized so that

$P_M^{(0)} + P_M^{(1)} = 1$  and  $P_M(x0) + P_M(x1) = P_M(x)$. ← This is a very general normn's not what I need

The

15 : 241.36  To obtain minimal error in probability, we normalize $P_M$ / using the constants:

to $P_M$

$$P_M'(0) + P_M'(1) = 1 ; \quad P_M'(x0) + P_M'(x1) = P_M'(x) ; \quad P_M'(x1)/P_M'(x0) = P_M(x1)/P_M(x0) . \quad (2)$$

← eq

→ Probly should give normn eqs (those in (LiV) or Sol 78

→ To use $P_M'$ for prediction.

maybe give
$P_M'(x1)$ in terms of
$P_M'(x1), P(x1) P(x0)$

$$P(x1|x) = P_M'(x1)/P_M'(x) \longrightarrow \boxed{242.00} →$$

$\boxed{\#3}$ ← eq

242.00

The convergence Result  (Sol 78 § p42E) assures us that $P_M'$ gives very good predictions.

25 : 242.40  The proof of the convergence theorem of eq (4) is for the normalized

given in (Sol 278)

universal measure $P_M'$. Peter Gacs (Gac 97) has shown

93

that it is also true for the unnormalized semimeasures, $P_M$, but the convergence

constant, $P_0$ is much larger than the corresponding constant $P_0'$ for $P_M'$.

$-\frac{1}{2}$ in

The difference between the two will be 1 the natural log of the normalization factor

smaller

Renormalized measure

for large n .

The result is that the probability error for $P_M'$ converges much more rapidly

than does for the semimeasures, $P_M$.

Gacs also shows (ibid) that the corollary corresponding to eqs 5 and 6 holds

if $P_2(\cdot)$ is an (unnormalized) semimeasure.

Marcus Hutter has shown (Hut ) that these results hold if we use

alphabets with more than 2 symbols.

Tech Report IDSIA 07-01 ∞ 2001
Convergence and Error Bounds
for Universal Prediction of
Non binary sequences .

Sec.
→ 245.00

248.15

00: 243.30

(Sci 78 p 426)

The convergence theorem assures us that this technique gabs good probability values.

The way it works:

[ Dimacs 1.22 - 28 ]   Suppose you had a device, $\mu$ generating ....

Bau   put in

$$\underset{\mu}{E} \; \underset{m=1}{\overset{n}{\Sigma}} \; \big( P_M(x_{m+1}=1 \,|\, x_1, x_2 \cdots x_m) - \mu(x_{m+1}=1 \,|\, x_1, x_2 \cdots x_m) \big)^2 < -\tfrac{1}{2} \ln p_0' $$

very good probability estimates.

(4)
eq

Dimacs 1.29 to 0.32   $p_0$ is the probability that ... to and of $p_1$ ... and of recovery general.

A useful corollary to this theorem:

The truth this theorem is based on the idea that for any

if $M$ is a computable probability measure (or more generally, any finitely describable probability measure) then

$$\frac{P_M(x)}{\mu(x)} > p_0'$$

$p_0'$ is a positive constant whose value depends on the structures of $M$ and of $P_M$, but is independent of $X$.

Equation 4 can be usefully generalized, so that for any two probability distributions $\mu$ & (not necessarily universal or computable), $P_1$ and $P_2$

if $\dfrac{P_M(x(m))}{\mu(x(m))} > \alpha(n)$         $\dfrac{P_2(x(n))}{P_1(x(n))} > \alpha(n)$

$n$ being the number of bits in $X$

Then equation (5) holds with error bound $-\tfrac{1}{2} \ln \alpha(n)$.

(perhaps use $P_1$ and $P_2$ instead of $P_M$ & $\mu$ — and rewrite eq(2) )   (5) eq

use Balls

Equation 4 can be usefully generalized so that if $P_1$ is any normalized probability measure, and $P_2$ is any probability measure or semi-measure, such that $X(n)$ is any string of length $n$

$$\frac{P_2(x(n))}{P_1(x(n))} > \alpha(n) > 0.$$

(7) eq

where $\alpha(n)$ is a function of $P_2(\cdot)$ and $P_1(\cdot)$ and $n$, but not of $X$ ———

Then $\underset{P_1}{E} \Sigma$

$< -\tfrac{1}{2} \ln \alpha(n).$   (8) eq

NIPS                          Possibly ommit ~~section~~.

10:238.40 :   T. ~~Tony~~ 239.00-40 can be ~~a~~ ~~long~~ abstract or short introduction.

At. t. and of t. Intro.

Section 1 deals with ~~the~~ sequencial prediction and its universal distribution.
We discuss generalizations of
~~Generalizations of~~ this problem ~~and as embedded for~~ that are amenable to the same formal solution.
~~A convergence theorem for the normalized universal distribution for sequence pred~~
This is followed by a convergence theorem for the normalized distribution and some
more recent generalizations of it.

Section 2 ~~deyendes~~ ~~deals and~~ deals with / extrapolation of ~~a set~~ a set of unordered
strings and/or numbers, and gives an associated convergence Theorem.

Section 3 deals with Operator induction, and gives the associated convergence Theorem.

Section 4 ~~Discusses the incompatibility~~ Discusses the ~~its~~ incompatibility
of the universal distributions and ~~its~~ (limited) ~~bearing~~ bearing.
~~its lack of~~ on practical implementation.
~~We also~~ ~~We give a quantitative estimate of the frequency with which~~
~~the incompt~~        We discuss the "halting problem" and the frequency with which it
has a bearing on the computability of the universal distribution.

Section 1     Sequencial prediction

The universal distribution for sequencial prediction is a probability distribution on strings
that is obtained by assuming the strings were the output of a universal machine with random
input. We will at first consider only universal Turing machines with / unidirectional input and output
tapes and an infinite bidirectional work tape. ~~Though~~ Most of these conditions can be
relaxed using more general kinds of machines and input output mechanisms.
~~To obtain the probability of a finite string x~~      How do we ~~use the~~ this distribution / can express the / probability
We will ~~rugged~~ the probability of a finite string x / to be the sum of the probabilities of a particular / finite string, x?
of all finite or infinite strings that have x as prefix. Let $[S_x]$ be the set of
all binary programs for our machine, M such that $M(S_x)$ gives an output with x
as prefix. To prevent double counting we have the additional constraint on the set
$[S_x]$: If we drop the last bit of the string $S_k$, the resultant program will not have output
with x as prefix. With this condition the probability of x becomes the sum of
the probabilities of ~~all~~ all of its programs:

$$P_M(x) = \sum_k 2^{-|S_k|}$$

$|S_k|$ is the number of bits in $S_k$ and $2^{-|S_k|}$ is the probability of an input that has
$S_k$ as prefix.                                                    → 243.36
                                  ~~normalized $P_M(x)$~~

To use this distribution for prediction:

$$P_M(x1|x) = P_M(x1) / (P_M(x0) + P_M(x1))$$

$x0$ ~~*being*~~ the binary string x followed by the bit, $\phi$.

This gives us the probability of a 1 following the string x.

GREEDY
↓

≈0.237.⁜ : We will have several choices of dstns to make. Rather than chose & sing to "Best"
one (& one H. max pc), chose dstns w. ʷᵗˢ ≥ their pc's. ← ( NOT so GREEDY)

___

Try to design my System so it is ergodic w. not too small pc of going into all possi.
states — essentially, this means that it can't get "stuck". On & other hand, we want Y.
System to spend almost all of its time in V.G. states.

Hvr, one idea of Universality is that it has all ∞ of states & that their p's can
be varied. On & other hand, in say an U.D. if no state get pc=0 !

___

In Gary Wolff's abov system he would periodically reparse it. If some bad dstns were
made, (perhaps), reparsing would make them unimportant.

It may be that in Eng English, that parsing rules are as only pc dependent —
which means that usually/likely ᵘˢⁱⁿᵍ our parse will be adequate.

On & other hand, in more general systems in which I will be using (grammatical) Induction,
this will not be true, ≥ & ideas of 237.3 off may be best.

___

Anyway, it would seem that there are 2 essentially diffrnt Approaches here:

1) 237.3 off : to try various dstns, pc's act. sequentially branching various parse trials —
Doing many trials for entire system (parse/dstns) ···· predict based on
wtd. combination of the best ones.

2) Do various dstns but repeatedly reparse after every (or every few) dstns,

2) would seem better if it is likely that there are only a few V.G. parses a/o (set of dstns).

NLPS

GFG Discovery

**: 237.40:** Say $zb$ occured 10 times & t. history to pars.it would do $z \to of ab \to C$ & leave 7 as $ab$. So we randomly chose $ab \to Cw$. per $\frac{3}{10}$ to obtain & parse we then increment the corpus ($15 \to 16$, say) and look for new doctus. If none are found & me corpus size to search: loop on till new doctus. is found, then randomly prune as before.

**.06** We run thru t entire corpus this way making random choices. Ideally, we would go thru t entire corpus this way, many times, & use a wtd sum of t resultant Grammars, for predn. ~~Alternatively~~ Wwat kind of Dist. we'd get on.wts of t. Grammars/parses is unclear. Will it be a few agrmt. or a broadspread agrmt? Will we have a Gaussian Dist. on ~~wts?~~ In wt. ? (most likely).

**.11** An alternative methodology more room, (perhaps): Do corpus length 20, Select out best 100 grammars & continue with them to 30 symbol corpus: return best 100 grammars & continue etc.

Another way (in t. spirit of .11) Go thru t entire corpus as .00 → .06. Next go thru t. entire corpus again (random choices), but as soon as (i.e. avars), we get too far behind ~~the best~~ nor any point in t best run thus far we abort x. trial, & **backtrack** (don't merely go back to ~~zero corpus size~~), back to a point at which "things aren't so bad". We work on a particular branch/trace for a certain amt of time & if we don't seem to be making enuf progress per unit time, **then** we may go back to t. root (corpus size = 0) & start over.

~~For~~ certain kinds of stochastic processes ($\overset{one}{\text{ergodic}}$ kind, perhaps), if we go thru this branching process, we will almost always end up in t. hy.prc region. Whether t. process of 237.30 ff is of that kind ··· I don't know. In ergodic processes we are ~~almost~~ always at a finite distance (prc - wise) from ~~~~ any other state, so for a long run, the worl/distance between runs, will be $\overset{\text{re}}{\text{"Negligible"}}$.

**.29** [ It is conceivable that in certain kinds of processes, one could branch into regions of state space **.29** [ that were quite diffrnt Statisticly from others.

**30** While a long English text may be Ergotic, my method of analysing it may not be. My impressn is **31** that decisions made early in t. process cannot be "undone" — so if one makes a bad mistake in defining something, certain symbols will be permanently recored from t corpus, so one couldn't ever use them "properly" later ← Tho maybe not so extreme since normally one does not completely "parse away" any particular symbol — one only reduces its probability (which may be bad enuf)

A way to ~~~~ perhaps avoid .31 ff. At each pt. in t. process,

**CFG** (& maybe **CDG**) **Discovery**

[CFG Discovery] IMPT !

.30 FF looks VERY GOOD! ⟵ Perhaps Great Breakthru!

DO:

: This is not a new idea, but I had sort of forgotten it — I think I had some of these ideas in Sol 56 (& in md. mt machine)

From corpus, using methods of NB 141 (latest methods: old 214 (didn't work).

& define zgms from combine zgms to get ngms, etc.

Next get ngm sets: Some ways: b & t. set of all b's ⟹ bt is likely to be used ngm. or ac ⟹ set of all c's ⟹ ac is used ngm.

Other ways to define ngm sets from ngms. various environments: perhaps context dependence.

After we have ngms & ngm sets, we can try making new ngm sets from old by

① concat ② Boolean combinations (vscally "AND"). If we use any other comb rules we may end up w. a non cfg (but that's ok)

■ Then test recursive rules like $A = AB$ ; $A = BA$, $A = BAC$, etc.

We will test such rules $\iff$ we see an instance of it: i.e. $a_i$ and $a_i b_j$ are common ngms and so we try $A = AB$ ... $a_i$ in an ngm set $A$; $b_j$ in ngm set $B$.

If may be that rules & ngms of this kind are the only efficient way to discover CFG's, ⟵ LIKELY!

It seems clear that w. adequate SSZ, longer and longer ngms will be found, & they will be expressed by grammar rules, & the heuristics for trying grammar rules will be "best t. rule works for 1 (or a small SSZ).

We can "tune" this heuristic by balancing SSZ in heur. w. cost of trying t. heur, & "Yeald" of t. heur.

OK: Back to 237:30 :40 ! In town: Starbucks — abstract of Sol 89.

[5] Non CFGs: I had stared clear of doing $pc = 1.0$ parses. T. result was many possible 11 parses. — lots of work to keep track.

I could try 100% replacements of $ab \to C$ say. Maybe not always right, but easy to go to "next level"

Another possy is less greedy parsing. Do 2 levels at a time — so I have some "look ahead" for opportunities at "next level",

Alternatively, after a new doc has been used, select a set of random parses & see which are best compressed. (Actually all parses would have to be investigated) Actually, .30 gives t. true picture! Each parse is a highly branching stochastic (rooted!) — So we just pick the parses at random & continue along each one. We trace thru as many as we can ... obtaining many grammars ... then we pick a set of t. grammars of hyest pc.

.32 should also be used in t. more complex parsings resulting from use t. recursive grammar rules discovered using .00 – .48 (+ .05 – .26)

This .30 AF looks VERY GOOD! T. way it works! We do t. first 45 symbols of t. corpus! try to find good ngm zgms. For each zgm we have a parsing decision!

⟶ 238.00

∞:236.40 Problems in ~~inductive inference~~ probabilistic induction are of two general kinds: In the first, we are given a linearly ordered sequence of symbols to ~~extrapolate~~. A very general solution to this problem is based on the universal probability distribution, and there has been much work done in ~~approximately it~~ finding good approximations to it [references]. It has been ~~proved~~ shown that for long sequences, the errors in probability estimates converge rapidly toward zero.

~~Though the second problem is~~ In the second kind of problem, we want to extrapolate ~~on~~ a (unordered) *finite* ~~sequence of strings~~ finite strings and/or numbers. ~~Though ~~ A universal distribution has ~~been derived~~ that solves this problem (Sol 59), we will give a convergence theorem that shows ~~that~~ it to give small errors as the number of examples increases — just as with sequential prediction.

~~In~~ In Operator induction we have an unordered sequence of ordered (pairs of elements) ($Q_i$, $A_i$) (that may be strings and/or numbers). Given a new $Q_i$, to obtain the probability distribution over possible $A_i$'s. The $Q_i$ can be ~~questions~~ questions in some formal or natural language, the $Q$'s can be associated A answers. The $A$'s can be inputs to some unknown stochastic device and the $A$'s can be outputs (The Identification problem). The $Q$'s can be descriptions of Mushrooms, the $A$'s can ~~be~~ tell if they are edible or poisonous (The categorization problem).

The $Q$'s can be number and the $A$'s can be exact or noisy values of some unknown function of those numbers (The curve fitting problem).

We will give two solutions to this problem based on universal distributions, and give associated convergence theorems that support their precision in prediction.

~~All~~ All universal distributions are of necessity, incomputable. We will ~~discuss~~ show how the practical application of these distributions is usually not affected by their incomputability.

The incomputability ~~is~~ of the universal distributions is associated with the unsolvability of the "halting problem". We will show that ~~it~~ for ~~practically~~ most prediction problems ~~the~~ The occurrence of a ~~halting problem~~ becomes very unlikely ~~Often a finite amount of data~~ as the size of our data set increases ~~decreasing frequency~~ that rapidly decreases with sample size.

(While this fact is of theoretical interest it is very rarely of importance in practical prediction.

We will show that the frequency of occurrence of halting problems, ~~rapidly~~ approaches zero rapidly, ~~as the size of our data set grows~~ with the growth of the data set.

N.B.

30 : (SN) & ES pco? conv. The reason *they* mean that for ~~certain~~ M ⊆ CPM, ~~~~ one can
get ~~an~~ arbly small/~~knowable~~ ~~by~~ doing ~~so~~n; But ~~probably~~ for the Rubick kind of "computability"
it is not practical computability: e.g. ~~if~~ if we limited our models to prim rec. funcs —
**They** would give "computable" ALP, but not ~~practically~~ ~~computable~~ ALP.

If μ ≡ Born seq. in p̄≡.5 This should give ~~interesting~~ statements about halting probly
of an arby ∪mc!  To simulate "μ", just copy input onto output — so it should have
z fairly short desa. T. desa $c_{du}$ be 1 bit in a ∪mc as so! for any input w.
prefix φ, rest of input is copyd to output. for prefix 1, machine acts like ~~computer~~
∪10, ∪mc; for rest of ~~~~ input.

10 [SN]  It may be that my/~~proof~~ of QA induction in "Appendix B" of Report is wrong! [PP 23-25]
[conv.] [(space)]
T. idea of proof:  that+. seq A₁sA₂s···Ans ⟵ (s is marker betw A's.)
constitutes z seq. to which the [convergence w. radius] "corollary" applies ———  is correct; but K ~~~~ has 2 parts:

1) dervn of t. operator M (which is indep of ~~~~ [Qᵢ] )
2) derm of the ordered ~~substit~~ sequence, [Qᵢ] : ~~[ord]~~ ~~This Fact~~
The ~~~~ dervn long [Pvt-2], grows with n.  ~~~~ is relevant to t. last # of Appendix B or p.25.

Unclear as to whether this is a Bug or not!  It looks like a Bug: i.e. To derv [Aᵢ] ~~[orig]~~ ,
t. dervn of M is inadequate.  But is it legl to assume the [Qₐ]ᵗʰ is "free"?
.19  T. Q is can we obtain germzn of # T? can! # that includes "free info"?  [STS]
20  ~~~~ Well, ~~even~~ as the [Qᵢ] sequence changes, the length of the M that can get a particular
∏ P(Aᵢ|Qᵢ) for that sequence, changes.

Whether t. {Qᵢ} are "free" or not, is irrelevant, here, th. Q is "Does t. corrolory of (.19) apply?"
Well, not knowing [Qᵢ], T. U.P.D. guarantees a M that gets a dervn length of the
[Aᵢ], which is *some* constant.  How t. constant will grow as [Qᵢ] grows, is unclear ———

If t. Qᵢ's remain in "same Domain", its possl. that t. const will not ~~grow~~ much (14 std!).

~~~~ A similar problem occurs when we use the ∏ P'(Qₜ, Aᵢ) form of induction.  (235.02)
If we have z certain pⁱ ~~that~~ works well w. t. known {Qᵢ, Aᵢ} set, — but if we
change t. character of our Q's much, the pⁱ will change — z it will probably
30  need a much longer ~~desa~~ dervn.  So it would be well to give both conv. thms z perhaps
explain t. deficiencies of **each of** them.  Maximum U.P.D.  [Based on associated UPDs.]

Qu. : Back to 235. 30-.40  Intrdn: maybe start w. abstract of Sol 99?  [Mention UPD's.]
Probablistic induction/~~of~~ of two kinds; .... [probably any] [can!]  [1,3,4,5,6,9,10]. Two general
solutions to the second kind have been described. Sol 99.  ~~Whether~~ Though a convergence
theorem has been proved for sequenced induction, Shown

00:  So t. Q's, just how to present this?

It  P(A|Q) is normed already, /  P̶ ∑ⱼ ∑ᵢ ∏ ∑ P(Aᵢ|Qⱼ)  to unit total

Q2   0 Promise ⌉ Make
              ∑ᵢ ∏ Pʲ(Aᵢ, Qⱼ).    It is not clear that these methods Make t. some thing —
i.e. That they get t. some soln! We have conv. Proms for both, but differ't things converge to
differ't constants.

What I can do,  Give t. "conv. Proms for Bags: Then say it will be provided a special case
of Op induction    a theorem we will prove later.

for Op induction:  2 ways to state it.

10 2→ (SN) can t. f. production of t. latest QATM be regarded as Encyc. induction of BAG.
       T. encyc is One "Given" i t         can use into encyc.
       In t. latest QATM, 1.      Encyc is t. set of Q's.. (Y encyc can be op into in addition
       to t. set of Q's. — except that t. A's are somewhat Tied to t. associated Q's!

       So: We will present two somewhat different solutions
for Operator induction. This kind of induction is very common, and comes in
many guises.   We will discuss the relationships between the two solutions.

20  (SN on) CIAP:  A ██████ young, concious machine is like a Baby Tiger; cute as hell and
lots of fun to play with — but as it grows up, it becomes increasingly
unpredictable  unsurgato, and  ████████  eventually, [ (very) Deadly.
                                                        very dangerous
                                                        very life threatening ]

For Operator induction, we will present two different solutions to
this very important problem.

→ Give the conv. theorem for Bag induction, then say that proof will be a corollary
of another Convergence theorem for a more general kind of induction.
                    review  Sol 99 +                        post give conv. theorem in Sol 78 form,
In introduction: /         seql. induction, / bag induction discussed in ...in Sol 99.

30  2 methods of Bag induction given.
    Conv. theorems for Seq induction are well known. & various geniuses .. ... by Gacs, Hutter.
    We will discuss Here we discuss Op. induction.  Give 2 models & our convergence theorems.
    — and in assoc conv. theorem for Bag reduction            ("T. halting problem is usually not a problem":
         Then → expansion of this : Incompatibility due to unslv. of Halting prob. ( Open problem. If seq. tester is a.w.
                                                                                    CE To; will p(ν) converge to CE To?)
         For Brief abstract; we will give conv. theorems for Bag induction:
    Then 2 solutions to Of induction problem i.e. 2 assoc. conv. theorems / or expo conv. theorems,
    Discussion of applicability of these results to practical induction &
    Some theorems on the in frequency of incompatibility
                                            unlikely ways &.

Nips

β 20:47

154 lbs w.
particulated ; no shoes
5'11"

30:    : What I want in Introdu. : Dern of Problem, review of previous work.

[ Here could be section on seql. preda! ~~Give Go~~ Manbin & Fano. Theorem & Gacs, Hutter's ~~earlier~~ extensions of it.

Then Derb Sol 99 results. Tell what they left out (conv. Thm).

That t. soln. to f. O p^order induction problem suggested, here will give that soln. explicitly, & give a definite soln, perhaps more practicable soln as well as t. assoc. conv. Theorem.

There for t. most part t. forgoing results are all proved. Recorded and all conv. d.h.'s are incomputable

bearing (i.l.ch of beann)

10    we will discuss of t. least incomp results on practical proba.

Maybe ( ~~executes~~ ) ~~discusses has some interesting results~~
In all cases t. incomp 2 30.30 is a nice intro do t. ∑ P(x) result.

Later, I may want to discuss difrnce bctw P(A, Q) & p(A|Q)

P(A,Q) = P(Q) · p(A|Q) : But how is ρ(A,Q) related to P(A|Q)?

$$P(A) = \sum_Q P(A, Q) \; ; \; P(Q) = \sum_A P(A, Q)$$

P(A|Q) = α P(A, Q₀) perhaps normalized

so maybe $\dfrac{P(A, Q_0)}{\int dA \, P(A, Q_0)} = P(A|Q_0)$

Thus, Using t. input some target P(A|Q),

we have to normalize wrt A anyway! — Tho in some cases, (perhaps)

20    it would already be normed.

So small black copy of 60/64 papers.

So it would seem that one could do BOTH 2 ways:

→ ①  indirectly by d.f. on (Qᵢ, Aⱼ) pairs.  $P(A_i|Q_i) = P(A_i|Q_i) / \sum_j (A_i^j | Q_i)$

→ ②  Directly by conditional d.f.  $O^d(A_i|Q_i)$ t.o usually we have to be $O^d(A_i|Q_i) = O^d(A_i|Q_i)/\sum_j O^d(A_i|Q_i)$

Though t. Operator induction as can be looked upon as a kind of BAG induction problem,

we will treat it as a special kind of problem — having its own convergence theorem.

Ordinary BAG induction can be treated as a special case of Operator induction

We will present 2 kinds of solns to the Op induction problem.

30    Sometimes P(Q, A) is easier to model.
         "         p(A|Q) "    "    "    "

33    Conjecture: for every random no. [x] (say on (0,1)) there is an interval ε ∋ all nos. x±ε  on interval

33    have t. cost > some integer, R.  For small enuf ε, R is is a simple function of ε.

This is related to interpolating continuous fc functions

SOS    Perhaps show 32 - .33 for rational nos. in x±ε.

We can consider, say π, a rational nos. close to π. On interval π±ε, all rationals must have denominator > f(ε). (t. f. is some increasing function of ⅟ε.

(can soib⁰ Prob⁰ f(ε) is difrnt for π±ε v.s. √2±ε.

DO: 2317A0. [SN] Actually, Gac's proof has a same *kind of* corrollary as So(TS) TS:

I *think* Gac's proof *also* *proves* S~(TS TS, *and* it is dependent only on $\frac{P(KX)}{P(KT)} > e^{-k}$

P can be any 2 semi measure (it includes computable measures)

Gac's uses it *(Reing probly ratio inequality)*    I *think* t. only property of $\neq$ $P_M$ or $R_{MT}$

In fact (P331 LIV97) coroll z.1 does

06    just P(xt)    ($\leq$ h(s) $> e^{-h(t)} p(x)$)  ;  $\sum_{t=1}^{n} s_i \leq k(n)$

μ & ρ  can be any 2 semi measures satisfying (.06L).

Any normalized semi measure is a measure, so it works *for* normed measures —
except that −k(n) is *smaller* for normalized semi measures.

I *could* just write a shorter paper giving Bog's Operator induction & *assorted*
proue convr. Theorem for Op.ind. & say t. other proof is corollary.    So perhaps do that *first*,
then if I have time, I can add in stuff on sequencial prodn. & "section 5" on
P(us) convergence.

_____

Dern of Beg problem & Soln

*probablistic*

~~Eusstea~~ *in abstract*  We *have* described two kinds of induction :  In *the first we have a*
linearly ordered sequence of symbols that must be extrapolated.  In the second
we extrapolated an unordered set of finite strings.  ~~This Muldon~~ ~~There has been~~
~~much successful work published on the first.~~  For the first type of problem,
there is a well known solution *using A 2* ~~involving~~ universal probability distributions ~~&~~
~~there is an~~  ~~Bertram~~ associated convergence Theorem that assures us that the
probability estimates made will converge rapidly to the correct values as the ~~set~~ length
of the symbol string increases.
~~[illegible scribbled out line]~~  $\neq$ For the second type of problem
we presented a ~~different~~ kind of universal distribution ~~that~~ that would solve
the problem.  The present paper gives a convergence Theorem for
that solution.

We ~~[scribbled]~~ also ~~dean~~ discuss ~~[scribbled]~~ Operator induction ~~[scribbled]~~ .
~~effects~~ associated ~~t~~ universal distribution, and its convergence Theorem.

_____

31    [SN]  On use of sequencial trials ~~&~~ (w. back tracking) for prodn.

In sequancial trials, we try to adapt a pgm by adding onto its end.
A better way mite be to try "*General Modifn* of t. pgm." ≈ "Mutation": how efficient this
is, is on clear!  (what I meant in 31/ was that t. trial funchns ~~[crossed]~~ illus'd f. same
*prefix* function. J's "Oops" does this rather nicely: knowarys (never tried as the profit probr,
are able to ~~request~~ additional bits if they like).

Naps

> :SN It would seem that to discrete & continuous univl. d.f.'s would be very similar.

.00 The discrete d.f. assigns a pc to each finite string.

.01

.02 + contin. " " " " " infinite (or infinite+finite) string.

If we use .01 for sequential predn, we get Cover's "extension complexity".

Th. BIG Difrnce seems to be in implementation.

If we use .02 for predicting we can use a UIO device so we can tell when a code is

appropriate (whether + corpus is a "prefix" of it).

In .01 we often don't use a UIO machine to describe finite objects" — The objects

can have their bits written in any order. I got into this difty when I tried to use "extension

complexity" to get probs using that "dimensionality" estimator pgm for SM preda.

For any string, it gave that the string was obtained by a linear process (?≈) (The BDS pgm
Brock, Dechert, Scheinkman)

(I don't know where the pgm is, but 7.9.93 is a date in notes) : I think pgm was able to assign pc to

any string, — but I wasn't able to get it to do proper "sequential predn" — in which I'd

get pc of past objects as I moved along m+object. It may have been that pc was a

vector "global" function of + finite string. — Tho its possl that now (10 yrs later) I

could devise a sequential pgm.

A stronger pt in SM case above, was that the prediction pc's were not very strong

for + kind of predn I was doing! (I.e. single seq. predn). The S did get some

good looking "yield estimates"! i.e. if + actual corpus was a factor k less probable than

+ "default corpus" then one should be able to get a brokerage free yield of $\frac{1}{k}$ for that

sequence — I think I got some attractive yields!

> [ If we forced the discrete d.f. to use a UIO machine to simulate them, it would
work O.K.
> It may be possl to express + BDS pgm in an incremente(vay", so t.
> pc of X was found, than + pc of X1 was found using a modifnc of + value
> for X & modifns of parts of + "trace" of + calcn of + pc of X
> We want + calcn of X1 to be a minimal addition (in pc) t + calcn of X.

?0

6/14/03

For §2 (toys) use notation of Sol (99) 2 kinds of prob ind.
& Stickclose to Sol 78 notation, terminology.

231
233
272
358

## Section 1: Sequencial Prediction ①

We ~~are given~~ ⌜X(n)⌞ a data sequence of n binary symbols ( 0 and 1), and ~~wish~~ are required to find the probability distribution over the next symbol.

We assume that the ⌜sequence X(n) was generated by~~ the ~~⌟ stochastic source, P(·), that is able to assign a ~~conditional~~ probability to ~~the~~ any finite string.

■ P(·) is ~~a~~ sequencially normalized probability in the sense that

$$(P(X0) + P(X1))/P(X) = 1 \quad \text{and} \quad P(0) + P(1) = 1 \tag{1}$$

We ~~also~~ assume that P has a finite description." Later we will ~~express~~ explain more exactly, just what this means.

~~The universal distribution for sequencial prediction is obtained by assuming that the given data was obtained as the output of~~

The Universal distribution ▨ for sequencial prediction of binary strings,

■ is ■ the probability distribution on

① footnote: ~~The sequential prediction of section 1 follows the notation~~

The ~~discussion~~ explanations and notation of section 1 follow the discussion of Sol (78).

■ ~~The universal distribution~~ normalized Universal Distribution is discussed in

(Lev 97) pp 272-~~288~~ 314 ~~discusses~~ discusses sequencial prediction using the ~~universal distribution~~ semimeasure, and has some ~~(this provides)~~ discussion of the ~~normalized~~ universal distribution.

pp 315 - 334 discusses inductive reasoning ~~and~~ ~~and use of the~~ normalized universal distribution. Section 5 of the present paper ~~explains~~ ~~why our particular kind of normalization was used and why it gives better predictions~~ of ~~Section 1?~~

~~than the unnormalized semimeasure.~~ gives some recent results that illuminate the use of the normalized universal distribution for ~~inductive~~ prediction

Actually z ~~would~~ df has become obsolete for at most countable alphabet.

The Sequencial predn. w stay of an arby sequ of objects, has not been dealt w.

Perhaps as a corrolary (theorem) of Q Findexction?

In paper "2 kinds of prob ind."
I have one.w. on how below stationary sources. This would be OK for v. finite long strings being "adachied" so ternary A/phabet.

Claims re: QATM
① uses Univ df. — a best inductive key

ⓑ Incremental How induction based on past data

②) The algorithm transfor for learning so far wrapper can completely replace soft if just driven by data

③) can do OSC MPL code

→ predictions using

2
to predictions using.

:  • one or more ~~whole~~ members of a population of Genetic Algorithm trials.  Mutation
corresponds to | n=1 ; / Crossover corresponds to n=2.
Generalized

~~If the data elements are from languages.~~

Section 3 deals ∧with the third kind of problem; operator induction. We are given an unordered
set of n ~~object~~ pairs of elements, $Q_i$, $A_i$.  We are then given the n+1 th Q ~~object~~ element,
$Q_{n+1}$ and are asked to ~~that~~ give a probability distribution ~~given~~ for the n+1 th ~~element~~ element,
$A_{n+1}$.  ~~██~~ We assume that the data was generated by some unknown, but finitely describable
stochastic algorithm that assigns ≤2 probability. to each $A_i$ based on its corresponding $Q_i$.

  •  The Q's ~~input~~ and A's might be questions and answers in a formal or natural language.
~~████████████████~~  ———————(but not Rorschach)

  •  The Q's might be the inputs and the A's the outputs of some unknown stochastic device—
in which case operator induction is close to "System identification".

  •  ~~the~~ Q's might be a list of characteristics of a type of Mushroom and the corresponding
$A_i$ would tell if it were poisonous or not.  ~~Induction of this sort corresponds to~~
~~████████████~~ This is ∧usually regarded as a ~~████~~ classification problem.

                                                    for
~~If~~ each of the three prediction problems we will show that the associated
universal distribution converges rapidly to ~~the~~ generating distribution — i.e. its probability
assignments to the predicted element ~~are~~ become very close to those of the generating
distribution as the number of data elements, n increases.

While this would seem to give an ideal, almost ~~perfect~~ solution to most
induction problems, the road is not very direct since all universal distributions
are formally incomputable.  Section 4 will ~~to~~ show that ∧then incomputability of the
universal distributions does not ~~with~~ limit their use in ~~predictions~~ obtaining accurate
predictions.

Section 5 will ~~to the~~ give some ~~recent~~ results on the ∧very rapid rate of convergence
of the ∧undefined "output of universal distributions. |

The incomputability of universal distributions is associated with our inability to know,
in general, whether a particular universal machine will eventually ~~██~~ stop ~~████~~
( Turing's "Halting problem").  Section 5 gives some new results showing that
as the sample size, n grows, the probability of a universal distribution being undefined
because of the halting problem, approaches ~~zero~~ very rapidly.

Nips

'0:

: In Sec. for # 229.175: The Q, A pairs can be questions and answers ~~(not necessarily correct)~~ (not necessarily correct), or the inputs and outputs of some unknown stochastic operator (the "identification problem") or ~~might be~~ Q might be a list of characteristics of mushrooms and A ~~could~~ be ~~the~~ ~~a label~~ "poisonous or edible" (the categorization problem). Q ~~a problem~~ description and A ~~a better~~ algorithm that successfully solves the problem. (They must be input Q = problem ~~A better way!~~ Just list the possys on separate lines w/ balls "•". A = solution?)

~~Intro~~ Give section nos. in which each kind of prob is treated ( §'s 1, 2, 3)

.07    Int. abstract (intro do tell about discussion of incomputability & how to univ. d. f.'s are actually approximated.   How much of this I will want to write, is unclear, at present.

It is possl. that I may not want to spend so much time on .07 #.

For discussion of (Multidim. General.⌐ sequential prediction. I could ~~name~~ Cana. "prefix set" — or just a "prefix". Define "M strings" (Multidim strings) —— M strings need not have ~~a~~ data on all ~~of~~ its "cells". We can have large finite alphabet ~~or~~ ~~three~~ countable alphabet

$X_1$ is a prefix of $X_2$, if : If a cell of $X_1$ ~~is assigned a~~ certain ~~a~~ elefabel, then t. corresponding cell of $X_2$ ~~it~~ has t. same symbol. $X_2$ may have ~~additional~~ symbols in cells which $X_1$ has no symbols.

Use a UMC w. no "stop" state, so no structure for output cells.

~~unwanted section~~ No point to use any thing but a unidirectional √random/input tape.

~~See re write of the first par.~~ 229.02: Title: Three Universal distributions and Their Convergence Theorems

(Introduction)   We will describe three kinds of induction problems, and will associate with each, a kind of universal probability distribution that can be used to solve it. Each such distribution has a convergence theorem that ~~gives always~~ quantifies the precision of these induction techniques.

In section I ~~deals with~~ The first kind of ~~induction~~ problem. sequence extrapolation. We are given a data sequence of n ~~objects~~ elements — they may be strings and/or numbers. The problem is to give a probability distribution on ~~the possible~~ the unknown, n+1ᵗʰ ~~object~~ element. ~~the sequence~~ We assume the data was generated by ~~a~~ an unknown stochastic source that is describable by a finite sequence of symbols. The sequence ~~about be:~~ ~~Bell "type" of p(x)~~ might Be daily values of a ~~financial~~ Stock or Bond ~~instrument~~ security, a set of meteorological parameters for a sequence of days at ~~a~~ New York City, ~~a~~ pictures of a particular person ~~taken~~ yearly, ~~the values of an unknown mathematical function~~

Section 2 deals with the second kind of problem. ~~em~~ We are given an unordered set of n data ~~objects~~ elements and are required to find a probability distribution over the ~~next~~ n+1ᵗʰ ~~unknown~~ element. Again ~~Here~~ we assume the data to be generated by an unknown & but finitely describable stochastic source. The set of objects might be: Called "Bag Induction" or "Grammatical Induction" We are given an variously called "Bag Induction" or "Grammatical Induction"

• A set of acceptable sentences in some ~~~~ formal or Natural language.

N ips

00 | : [SN] Q: how to divide up into better. Intro & rest of report. e.g. In deem in
sequence extrapoln., how much is in "Intro" v.s. how much in "Section 1"?

02 | Restart: We will describe three kinds of induction problems — how
they are solved using Universal distributions and give theorems concerning the accuracy
of those methods.

The first kind of induction problem is sequential extrapolation. We are given
data $x_n$ (sequence) of $n$ objects (R ose are strings and/or numbers), and we must give a probability
distribution on the unknown $n+1^{th}$ object. We assume the sequence was generated

10 | by some unknown stochastic source that is describt
by a finite sequence of symbols.          as data em.

15 | ⟶ see .20
In the second kind of problem, we are given an unordered set of objects,
and we must give a probability distribution over the $n+1^{th}$ possible object

.135 | We assume the data to be generated & an unknown stochastic source of finite description.
"Operator induction"      data        (.25)

14 | In the third kind of problem, we are given an unordered set of $n$
($Q_i, A_i$)object pairs, we are now given the $n+1^{th}$ Q object and we must
give a probability distribution for the $n+1^{th}$ "A"$_{n+1}$ object.

.175 | We assume the data was generated by some unknown stochastic algorithm,
that assigns a probability or accessible A, based on its associated Q object.

20 | (Insert for .115) The sequence of objects can be the daily values of
a stock, or a set of meteorological parameters, or pictures of a
particular person, taken yearly, or the values of some unknown or describable function
for successive integral values of its argument.

25 | Insert for .135   The set of objects might be a set of acceptable
sentences in some natural or formal language, a set of pictures of people — not good enough; usually non-recursive(able) Brain.
⊗→ known to be criminals, } I wrote list of Gramm d. very prob "BAGs", w. only
pos. + ve data.     [ Actually Gramm induction is v. rare: QA is much commoner.→ what about Mutations, crossovers? for GA? SSZ=1, SSZ=2.

[Good idea!] When I'm writing & I get stuck at a point for a while, write down just what
the problem is, & continue on. Keep at all times a list of unsolved "problems" &
work on one when I have an idea about it. I think this will save time in writing

30 | (& in other prob solving as well! working out clearly what the problem is, often solves it!

[SN] Re: GA: If I were able to get a set of cmds w. P.d. over expected G,
and time needed to create a test, we could do more intelligent GA! Also do it w.
SSZ = 2 or 3 or more → "many". We should try Cmds in $\frac{\Delta G}{\Delta T}$ order.
(woops! This assumes "linearized G", ... so we need to Linearize G "first"!)

[SN] in .10 & .14 perhaps mention "self delimiting" & practical ideas.} Not yet! I'm just
during the problem, not the soln!

Nips

Maybe abstract 20.??

~~Intro.~~ These given by Ra

We describe three kinds of prediction problems; ~~each~~ with each is associated a universal probability distribution that enables the solutions of problems of that type. ~~Associated with~~ Each of the distributions ~~has~~ a convergence theorem that limits the difference between the \universal distribution and ~~its~~ ~~given~~ probability distribution that generated the data set.

All universal distributions must be incomputable. This incomputability ~~arises~~ ~~caused by~~ ~~the without~~ arises ~~when~~ ~~various inputs to universal~~ because we cannot know ~~that~~ whether certain inputs to universal machines will ~~cause~~ the machine to eventually stop or not .... (The unsolvability of the Halting problem). We will show that as the sequences of data grow larger, the probability of ~~its~~ not halting converges very rapidly toward zero.

~~It is fortunate that~~ ~~Fortunately,~~ The incompatibility of universal distributions ~~usually~~ ~~has~~ rarely inhibits their ~~on their~~ ~~Little effect on their~~ use in practical prediction rarely ~~does this inhibition~~ ~~Seriously~~ use for practical prediction.

~~VII~~

Introduction!
We will describe three kinds of universal probability distributions ~~used for prediction.~~ and their uses for prediction.

The first, which is dealt with in Section 1, ~~is~~ has been called ( LiV '78) "The Continuous universal distribution" It is the probability distribution induced on the finite and infinite output strings of a universal computer (e.g. a universal Turing machine) having random input strings.

~~If~~ If $P_M(X)$ is the universal probability induced on finite binary string $X$ using ~~the~~ machine $M$ as reference ~~the machine,~~ then the conditional probability ~~that~~ (knowing $X$) that $X$ will be followed by ~~the~~ the symbol 1, is

$$P(\# X1 \mid X) = \frac{P_M(X1)}{P_M(X1) + P_M(X0)}$$

$X1$ being the concatenation of $X$ and 1.

If ~~we~~ we have a finite string $X$ ~~n~~ n bits long, that has been generated by a stochastic algorithm having a finite description, then the stochastic algorithm will ~~associate~~ assign a probability ~~of each~~ for each bit of $X$. Similarly the universal distribution will give a probability for each bit of $X$. The first convergence theorem tells us that the expected value (with respect to the generating algorithm) of the sum of the squares of the differences in bit probability for the stochastic v.s. the universal distribution will be bounded by a constant. ~~Since this~~ Since this is finite, independent of $n$, the length of $X$, it is clear that ~~the~~ this difference must decrease more rapidly than $\frac{1}{n}$.

6/8/03
N↓S

Conv. Thms Paper: .00

6 : 2.02.40 :   Title : Convergence Theorems for Universal ⎧ Induction techniques ⎫
                                                            ⎨ Induction models      ⎬
                                                            ⎩ Probability Distributions ⎭

>1   $\tilde{\approx} SN$   Is t. Operator Induction problem Covered by "BAG induction"?

At first Glance: "Yes":   In Bag induction, we ~~are~~ are given part of an object
( Bcs "part" can be ≡ Λ(null)) — to get P.D. over rest of object.
The $Q_i$ $A_i$'s can be regarded as parts of t. object $z_i$. ( $z_i = Q_i, A_i$).

.08   $z_i$ ⎧(usually will)⎫ have ≡ punctuation, say comma, betw. $Q_i$ & $A_i$.
           ⎩ can ⎭

Furthermore, sequence predn. (+ continuous distribn) can be regarded as a special case of
bag induction, in which one or more of t. objects are allowed to be enlarged infinitely.

Data can be whole or (~~usually for~~ always for infinite objects), part(s) of objects.

The OP induction $QATM$ ➡ is a special case of BAG induction, it does
SEEM diffrnt :   T. solution I'm using ⎡for QA⎤ seems to take a diffrnt form
                                          ⎣     ⎦
from ~~that~~ suggested by BAG Induction.  i.e. we feed in ≡ $Q_{i+1}$
& we directly get a d.d. over $A_i$'s w. hypoth $pc_i$ first.

In Bag induction, we usually get a D.P. over t. set of "next object"'s
T. Diffrnce is n't very much, hvr, because in QATM we often have
part of a $A_i$ to complete — & probly in BAG induction, given part of a
a BAG element, we can ~~not~~ ask for P.D. on t. rest of t. ~~t.~~ element.

This "completion" can give a soln to t. sequential predn. Problem.

There may be some differences in QA v.s. BAG:  In QA, we are
given t. Q's:  In BAG, we have to code t. Q's as well as A's:

27   ➡ In QATM t. raven? no info into Q's.  In BAG, Q's hovers much in t.o.s A's.

I think I thought about .27 a lot before doing present QATM Model.
Superficially, it would seem that in QATM corpus, Q's & A are symmetrical: that one has as much "info"
about t. Q's as about t. A's.  Anyway t. discussn of .01 —.08 , on how QA is exactly covrd
by BAG induction — seems Rite. !  Hvr in QA induction, T. TM would not have any idea
as to what to do if given a null Question, or partial Question.  BAG TM would be able to deal
w. that, hvr.


Verification & Reasons
Supported by Wallace.
A unified Theory of prediction.

(Lsearch)

2) It uses the universal distribution (or — more exactly, approximately it) to do induction. As I have mentioned, this is a very good, very general method of induction.

3) To search for models and to search for good, relevant problem solving each needs, we started in what I call Phase I by using Levin's universal search algorithm. For the additional information state that we start out with, this is close to the most efficient search technique possible.

In later stages of training [which I call Phase 2] it acquires different kinds of information and uses a different search technique that is better than Lsearch, for that kind of information matched to that kind of information. I expect it to be much better than Lsearch.

4) At all stages of its learning, the searches are guided by a conditional probability distribution. The rapid success of the search is much dependant on how good that probability distribution is. The system executes the modification, the "updating" of this guiding probability distribution as one of its normal induction problems. Its a real "Bootstraps" process, so that the method of improving the guiding probability distribution, improves as the distribution itself improves.

5) The system is able to relate any problem domain to any other problem domain. If the problems in the two domains have common elements, — structural, analogies or whatever, the system is able to use this commonality for induction. It is an idealized kind of "transfer learning"

2 1 ; 2.24.39 ! insert ! In inversion problems,

"My main Goal ···· B' 5.36 - .38 :
Some early papers I wrote in '86 and '89 I gave a preliminary design for system designed to learn to solve difficult problems. The present talk is a progress report on that system.
To pique your interest, I will list a few of its more exciting properties.

( "insert"

Here's how it works: B' 6.34 - 7.04.
Then discussion of a few details of the system.

Print up the talk up to here / give talk see how much time I have left.

Try to get entire talk printed up in big type.
Check on slides :
The main first heading retyping is Seminar
Also a 15 lines from A .05 - .22

For example, we have to design a car having certain specifications. We want the cost of the car to be small, but if ~~the design is completed after~~ it takes too long to complete the design, its value ~~to~~ the ~~car~~ cooperation will decrease in a known way, because of emerging in our days ~~the value~~ of competitors.

For a long time I was unable to obtain a useful solution to the same varying

G problem, but ~~about~~ a few ~~payments are~~ / I found ~~obtained set~~ what looks like a useful solution.

~~At the types of 1988 and 89, I and won't~~

Induction problems can be of at least three kinds:

1. The extrapolation of an ordered / sequence of symbols and/or numbers: ~Time series~ extrapolation.

2. The extrapolation of an unordered set of finite objects — Restore strings and/or numbers.

3. The extrapolation of ~~████████~~ a set of unordered pairs of objects.

SN ~~███~~ 3 ~~███ is~~ clearly a case of 2; 2 is clearly a case of 3: If stated in this way, it may be sure the conv. theorems will say what I want them to say. —

But the conv. theorem for 3 is significantly different from that for 2 / The Q's in 3 ~~can~~ ~~into the coding cost in a quite~~ don't enter into the coding cost — the they would if (3) is regarded as a special case of (2).

2 is ~~for~~ Though (3) is a special case of (2) I have found it more useful to consider (2) to be a special case of (3).

→ 'My main goal (B'5.36–.38) story B.S.

~~The approach I will describe~~ It is clear that we have to get an enormous amount of information ~~it~~ into a machine before it will be very intelligent. One approach, by Lenat ~~██~~ is to carefully program "common sense" into the machine, so it can eventually learn by reading books and/or ~~and perhaps~~ the internet. My own approach has been ~~to teach the machine~~ how to do inductive inference, ~~such give it a sequence of problems of increasing~~ ~~difficulty~~ more like the way I believe humans ~~acquire~~ acquire ~~their~~ knowledge — ~~by~~ by / inductive inference on a suitably designed training sequence.

There is a trade off between ~~the amount of information in the apriori distribution and~~ ~~the information in the training sequence.~~ The case w.

~~Finally~~ To motivate ~~the~~ interest in ~~the rest of the talk and perhaps to~~ people to actually get ~~the~~ ~~to~~ read the report, I will list a few claims for the system.

1) It is able to / learn to solve both inversion problems and time limited optimization problems. These problem types cover almost all, if not all problems in sciences and engineering.

~~Inversion problems are dr~~

Time limited optimization } see B'6.26 & 6.28. The P and NP problems of computational complexity theory. and — include all kinds of machine inference, surface reconstruction. All problems

5.27.03
Nips
222 (crossed out)
can't find 222
does it exist?

Well defined problems:

Discuss Ind probs & Oz prob: Well defined probs are Eqv. probs.

20 : 221.40 : Clips list: 213. 28R tf: (3 times) 1), 213,
Oz and 2 given, we don't know we have best soln, but we have no. That tells us how we could we are.

4) Use V.g. transfer Irny.

List kinds of induction problems it can solve;
looks like way humans work.
is very general prob. solver.

Seq. general time series

Bromdirision: Languages (Natural & artificial), clustering.

Q.A. induction: Operator induction, System identification, Categorization.

it several closeup and data are given it because 2 QA:

Another task would be to key Prot Sol 86, 89, 94 reported earlier work on this system —

This talk is a progress report — since Penn.

The system I will describe is one that I first wrote about in 1956 ···
It was called " An inductive inference Machine".
Its success was limited by my incomplete understanding of the induction inference process.

After 1960. I mainly worked on how to understand and apply the universal distribution to practical induction.

In 1986 and 1989 I wrote papers describing my progress in this area.
It had evolved into system to learn to solve very general problems. In those papers I was concerned only with two kinds of problems: INVersion problems and time limited optimization problems.

Not well written →

In version problems are " well defined problems" the P and the NP problems of computational complexity theory. —
They include solving equations, constraint satisfaction problems, symbolic integration, theorem proving and traveling salesman problems.

Time limited optimization problems include all kinds of inductive inference problems, surface reconstruction, devising scientific theories/ research planning. Designing an automobile in 6 months, satisfying certain requirements and having minimum cost is an example.

30

Most practical problems in science and engineering are time limited optimization problems. They are of several types: In general we are give a function G(·) that maps strings and/or numbers into real numbers. The problem is to find an argument X, for G(·) such that G(x) is maximized — and we only have known time T to find the x.

There is often G(x) takes a long time to evaluate, so we want to use few carefully selected trials for X.

• There is no time limits but G is a known/decreasing function of computation time T, as well as of X.

**Nips**     Continuos (in a sense) discussn of "Ti measuring of $P(c_i) > 0$" (00)     <238 fixd / 222

0 :220.40: 1   Superficially, it would seem that any $\mu$ likely to be a model of a real world phenomenon, would be an ordinary semi measure and be **Normal**. (i.e. not needing norman).

Does this make it a **CPM**?

A quick reply would be that particular p.d.'s are then $P(c) \le 0$: but still, Approximations to point rec. functions could be good models of processes that take a long time to generate (in < RW). — i.e. lots of c.c. — Not nearly ∞, but large enough to beyond to computing capacity we have. Could be in Geology & Economics (?ESM).

**SN** On the "hyper chaoticness" of a gas of molecules: I found that the uncertainty in one momentum component of a molecule increased by ~ 100 to 1000 factor for each collision.

Hvr, how does this face up w. **Louville's Theorem**, that keeps the h. volume of a region of state-space constant, as the system evolves? Maybe related to the manifold containing the system evolution into an extremely "fractalish" shape!

(218.00-.07): Introduction. or B' 1.02-.05 z' just → Cuso
- a) B' 1.20-.29 — what is 6. UPD & hours. Tied to predict.
20 : 218.40 » (A) 1.05-.32 ~ Karpool +19: Connectivity.
- b) B' 3.37-4.29 Incompatible
- c) B' 4.31-5.33 Subjectivity
- d) B' 6.22-7.04 Roughly how TM works.

219.18R ## "CLAIMS! Also : kinds of problems solved. : also kinds of doctoral problems solvd.

Maybe "claims" betw ____ c&d ....
- e) (C) — 6.04-11.42   how TM works : Details ( Explns about "Phase 1" & "Phase 2"
- f) OOPS

Claims can serve out w. claims for UofD; Plea Claims for TM, ___   commit to prior claims of B'—put them in here.

Don d) B' 6.22-7.04, ruffly how TM works, then some details from (C) from OOPS

Another poss is to leave all claims about UPD where they are now (3.28 ff in B'), then have special section for claims for TM — later.

I expect this model to be developed into a rather complete prob solving system for just about all kinds of well defined problems.

In respect to '" I have derived evolution of E. & prpd in a person as he solves problems in life — an ongoing updating of E. prpd ". The system will cast this as a user's model for its cognitive development ... active existence

NIPS

(Augmented) PST update .00

(Continued
Normal'zn .20)

: "PST update": This includes both updating of G. distribus of old, ≠PST's & t'
generation of new PST's.

I had a good idea on how to do this, but I've forgotten it & lost references.

"Formal"
I can perhaps write out a kind of "Theoretical" soln* :

Given all post (probi, pstj, T, (Fail/solve))) quads, to extrapolate to

PST (probk, pstj, T, Solve.) = probed or $P(T_{solve} | pst_j, prob_k)$
⌐→ a specific problem     ranging over all problems, (like "A" in QA prob'm solns)

context 217.00
ps update 220.00

192.00 –
194.22

Just as in solns to QA probs, I want to know pstj "values" that over particularly good:
i.e. zeroth moment of "T.d.f.
──────────────────── = max ( ≈ prob of soln/Time of solu ) = max.
(first n " " )²

So .09 is "what I'm looking for".

Next of importance: Say trainer inserts a bunch of PST's into t. system for "consideration".
How should TM best use this info? I assumed that TM should factor "t. set of PST's,
& make a grammar for them — (that Trainer could help by (partially perhaps) prefactoring
t. PST set.)

A possl. way for this to work: Trainer gives TM & a small set of (perhaps factored) PST's.
TM tries them on its set of problems, so it conpar data for (.00). The factoring
of t. set of PST's is used to help get solns (or trials) for ▓▓▓▓ (.06–.07)

.0 : 215.40    On Normalization:

{ What L..um probcy should! While any Universal semi measure > c · (any other finitely derivable
semi measure), some semi measures have norm const. that → ∞ as n → ∞,
so their Normd forms would not be dominated by a univl semi measure.
For there to be no mult. dominant univl. measure, this means that every univl. semi measure
must have at least 1 seq. for which t. norm const → ∞ w.  n.
Also, different these seqs must be dfferent. for difrnt universal measures.

──────────────────────────────────

On t. other hand, if we assume that t. seq. being predicted is very long (tho we only know
t. be ≤ 1000 bits or so) then M is unlikely to have a norm const. that diverges,
since then $\sum p_i^2(z_i)$ diverges & probcy $\sum p_i(z)$ also diverges — which means it
is very unlikely that this seq. did not terminate! In fact, $\prod(1-p(\omega_i))$ is t. probty that
this sequence did continue: — & it would (say) go "rapidly"(?) → 0 & it t. norm constant

$$\left(\prod(1-p(z_i))\right)^{-1} → \infty. \; \text{G}$$

So, it is likely that M does not have an infinite norm·n constant & ∴ that /& M is
normalized
must dominated by any universal semimeasure. If M is not finitely normalizable, it is likely that
ani
it would have stopped (if ▓▓▓ we have a lng data seq).

? ( What if M is a normed semimeasure w. an infinite norm·n. const.? — then normality
Universal measure ... ill ....... (?) Much ................ ... ≠ G..!

(right margin, rotated) probability prob·c·ty so measure x·es of no significance!  not  probability  is a  $\prod(1-f(z_i))$ is a probability  on t. other hand, note that

NIPS

0: 218.40: Perhaps better approach to writing this talk/paper:

Write it up roughly, w.o. necessarily good "turns of phrase". Then mark th.
parts that need rewriting, & work on them. This way, I will quickly have
a usable talk, & I can slowly make corrections if needed.

→ Start out by talking about th. this talk is th. T. Univ. D.G.'s Mach. Srvey — (much longer) } at my website
213.21 — .27 is menu outline:           +   "A Progr. report"  . . .            HHP:X www.ida.com/vmjal/.

Other things written: Ⓐ ① NIPS talk
                      Ⓑ ② kol. talk                    Ⓑ' kol paper
                      Ⓒ ③ summary R. Holloway,

Present plan:         ▮▮▮     Intro: Accuracy, incompatability, subjectivity
                      ▮▮

Intro: **218.00 – 07**

Discussn. of what is UPD?:  B 1.15 — 30    or   B' 1.20 — 29.

Accuracy   T. conv. from & its recent genesis; A 1.05 — .32  | so good
                                                              | B' isn't generally
                                                              | rather vague

    Incompatability: B' 3.37 — 4.29 ▮

    Subjectivity: B' 4.31 — 5.33

▮▬▮▬▮▬▮  ▮  ▮▬▮  ▮  ▮▬▮  ▮

        In the next part, I want an introduction that will motivate their
listening to the some of t. tech details that will follow — & also motivate
their reading of t. Report.

        Software will be a connecting discussn. "Rosetat." betw.
B' ▮ 5.33 and B' 6.22 → 7.04  which is t/ dern of (lets non-technical)  TM.
I will probably want to include stuff from Ⓒ summary. in this latter section.
                                          → 6.04 — 11.42  (entire summary)

        Write "Summer", I & assumed that almost all of audience heard
the ▮▮▮ "kol Lecture", the previous day.

        So perhaps to "Rosetat.", use mostly t. kol lecture B' 6.22 — 7.04
Perhaps modified to emphasize t. idea of " t. evolution of t. aping.(  .28 R ff)
        Actually, the discussn of this  B' (6.22 → 7.04) is quite short & might,
▮▮▮ in its entirety, be used in t. talk — perhaps augmented by t.

"Claims" of .28 R ff.
    Be sure to mention "Well Defined Problems"
    Expand "claims" by giving Examples.  ← This "CLAIMS" ▮ section
is a very impt. part of t. Lecture.
        So spend a fair amt of time on "Claims".
        Also discuss "OOPS" — a partial realization of part of Phase 1 (if its only determinist. "produ")

---

If we make a code for a
corpus and that code is not
"random" then we can recode
it and compress it.

Ⓠ if a code is not t. shortest
code for a corpus, must it
be non-random? — must
it have regs(s)?

If this is true, then coding &
recoding will always work!
Counterexample: Tcode is
Ⓐ followed by B.
A is a true min code of corpus.
When A is executed, t. corpus
is printed & machine stops.
A is Random. B is any other
random seq. So code for
A B is random & is a code for
t. corpus — that is not at
all minimal

.28 R  Uses best model of induction
       that we have.
② Uses best model of human
   problem solving that we have.
③ It is a fairly General problsolver.
   & At first I thought
   there were some kinds of probs
   to which it could not be
   applied — but as time went on
   I found ways to apply these methods
   to their solns or approx. solns

(Spec
(213.27)):  Title: T. univl. D.P. & M.L.

I'm going to talk about two things: First, the Universal Probability Distribution.
and some of its properties.

Then I will ~~~~ tell how to ~~apply the universal distribution~~ use this
distribution ~~~~ for a very general ~~~~~ ~~~~~~~~~~~, very powerful
kind of machine learning.

~~There~~ Three critical properties of the Universal distribution are

1) Its accuracy
2) Incomputability
3) Subjectivity.

~~The~~ When I first discovered the distribution, the main question was whether
it worked at all. There were various heuristic discussions, and approximate
applications that suggested it ~~~~~ would give good results

The universe

~~~~~ Expo: (special): While the incomputability of the Universal Distribution usually has no ~~direct~~
~~implications in~~ relevance to practical induction, it does have something to say about
Science as a continuing adventure: ~~~~~~

If "Many Scientists ..."

~~~~~~~~~~~~~~~~ We can never be sure that if we spent two more minutes
searching for a better theory, we ~~would~~ would not find a much, much better one then this
~~~~ best we have now.  This isn't at all bad — it means that nature will ~~~~
For those of us of an adventurous state of mind/ ~~~~ ~~~~ continue to be a never ending source of joy in discovery — That the picture of
Scientific ~~~~~~~~~~ invention will never end !

On "BACKTRACKING"    I have been assuming that backtrack can always be
characterized by t. "No. of problems Back" once Goes.   Not so! A better way would be to
"selectively backtrack", to violate the sequencial order in which t. TSQ has been given.
For early TM, this is not such a good idea ! T. TSQ ordering contains much useful
heuristic info.   For a more mature TM, t. TSQ is not being so carefully
constructed (& may be largely selected by t. TM itself) — in which cases,
"Selective bktrkng" is in order

List ways that
ALPHA is
better than other
systems.

Re: "Boost"
2 scaling factors:
1) In mature machine,
pc's of tokens will be about
rare so "factor of 1000".
1') much smaller
(maybe <1)
2) Immature systems,
pc of "Boost" ~~~~
instructions much
lower due to no. of
successful, frozen
~~~~.

5.22.03

**NIPS**

> 0: [SN] If we restricted our UPD to prim. rec. funcs, would it be "computable"? If so, & so if
> we had a convergence thrm, we could falsify it w. an "intractable funct" —
> would such a funct be prim. rec? — Superficially, it would seem so / i.e. T. generation of such a function
> would seem to be very "straightforward". Hvr., it does involve deciding of a certain real no. is > or ≤ ½ — which can be

> 03 [SN] Context: "Disc. of context" in (138.00 (also 133.02, & 40.16 -.32) may be adequate start.
> Context is used in Phase 1 (perhaps only Phase 1)-type of reasoning— i.e. t. zipnpd is for "similar
> problems in t. past " — not for optzn.

> 06     It is possib. that in th. searching which context is introduced, I "theoretically" used t. same
> zipnpd as ever, but because of t. ... search techniques effectively impose an additional
> bias (of an uncertain nature) on t. resultant zipnpd!

>     A remembrance of my ideas (more recent) ideas about "context":

> 1)   A somewhat General idea of context: We are doing trial codes to derive a corpus!
> T. pc of a token ... can be indep. of context. (This would be a Bern. model for
> t. set of tokens) — Or it could depend on all of t. tokens thus far, or on some of
> them a/o on a larger context — t. "nature of t. problem being solved" (... a very general
>                          ranging from
> dscrn of t. problem (e.g. Chemistry") to a narrow "Verbatum" dscrn of t. current
>     or "o z problem"
> problem.

>     More generally, this context could ... (help) control not only t. "next token",
> but t. pc d.f. on t. entire sequence of tokens to follow.

>     [SN] If Σ p(U̅ᵢ) < ∞ then p(U̅ᵢ) is normalizable! Is this at all a useful
>     way to look at p(U̅ᵢ)? T. nomzd p(U̅ᵢ) would be t. relative pc's. Σ p(U̅ᵢ) would be
>     t. total no. of "hits".

> (continued): IN A2(41), t. only context used was freqs of tokens, ... on t. previous present code.
>                                                                      partial
> — Also, we could define new tokens as we added along. More generally we could
> (.251) augment t. A2 (anyway) by including any kinds of "Legit" regys. Real naive
>                                            regy
> found ( "Legit means t. reduced total pc).

>     It may be that context is only relevant to ... devising induction codes i.e.
>                                                           conditional
> {Short codes for a corps} — Not relevant to devising / p.t. over PST's (for a particular
> problem)?
>                 (predictive)
>     Each new kind of regy found wrt. "corpus up to now" may need a modifctn of t. language
> used to derb. t. corpus. So we have a "Meta long." that tells us how we are allowed to do this.

> **Is this a "Tₘ - like" idea** (i.e. t. metalong)?

>     From t. foregoing (limited) analysis: Contextual regys are always a kind of "re coding": a finding
> of regys in a code. One then attempts re-re-coding a var(?) coding.
> Hvr., ... t. argts. of .00ff should be written more carefully, so that a person not familiar
> w. t. ... proposed system, could understand — (i.e. Me of Nov+ 1yr!)
>     NEEDS REVIEW!

MIPS    <u>Rev</u>

Also Note "<u>STACK</u>" of <u>ID 611.00</u> ff

ID 612.00 on "Anytime" Problems,
190.09 : "Things To Do" (Re: TM).

0:    :  A list of some impt. ideas I may have lost!    Previous list : 192.00    of G distribution.

1)  ▨▨▨ Conds under which { WON (OR) Net present ≡ L sech. :  One coup passed atim! Give
What functional form  G̅(·) must take.. What simpler equreten.    $, 1, 2 moments as touchout opep.

2) Discussn of "Context" : Just what TN2 ; its Geners, i how it is to be use for
"Updating." What is R. Updating G.c i how does "Context" info help opt z it?
133.02 is ref to Context, also { 139.00 } ← (This is a good idea !) , 149.16 — .32 : (Note 217.03)

3) I got a good understanding of t. G̅ares for discrete i continuous prodn i just how narrowly oft. wid R
continuous d.f.  was equivalent to $ pc density of soln. ( If t. choice muslvd was so tipaoh of f.d.t.) usually
192.20 — 194.25 is early discussn. It was pt on t. stack ! ID 611.38 (#4) || Perhaps adequate soln: 186.16 — .35    tree

4th
4) • Anthropic principal : 213.24R ; 214.00 — 215.10    187.20 — 23  || [186.30 — 35] → mem arg't
Anthropic •    ( 192.20
— 194.25 )

5) 190.09 / THINGS TO DO (Re: TM) /    I scanned from 216 back to 190. ... continue scanning !
Make index of impt id'ces .

Anthropic Principals.                    www.cursor.org

00: 214.40 : Actually, 214.32 & maybe the key to the "Antropic Principal." That the low pc. of any particular kind of world, — be it epistemological ("initial" zprip) or the constants of the universe — that this low pd is not important. Any state of the Universe is of very low pc. The only interesting things are not pc's but ratios of pc's (214.20) — which are used for prediction & for assignments of pc's to theories (or possib. "causes").

So it would seem that the Anthropic principal is not of interest: it does not give predus. — it is not "falsifiable". That this is true for Epist. as well as Physical Anthropic Principals

08: 214.26 : "Using the appr. hypoth to Evaluate itself" : It May be that there is No Alternative !

09
0
      So the "Anth principal" answers the Q "why is the pc ... of t. universe so small?" : This Q should not be asked, since ratios of pc's not pc's are of interest.

15 : 211.40 : A simple proof of ... that $\Sigma U_i$ & $\Sigma U_i^2$ :

1)    $-\ln(1-U) > U - \frac{U^2}{2}$    by 211.19 - 22

2)    $-\Sigma \ln(1-U_i)$    $< -\ln P\mu$  by 211.07

.19

3) from (1 & 2) $\Sigma_i \left( U - \frac{U^2}{2} \right) < -\ln P\mu$.    Well! not such a breakthru! I still have

0

to know that $\Sigma \frac{U^2}{2} < -\ln P\mu$ in order to get useful info out of .19.

.19 says   $\Sigma U - \Sigma \frac{U^2}{2} < -\ln P\mu$ ; but from that alone, $\Sigma u$ & $\Sigma \frac{U^2}{2}$ could be arbly large. If I knew $\Sigma \frac{U^2}{2}$ may $< \frac{1}{2} \Sigma U$ that would be enuf to prove $\Sigma \frac{U^2}{2}$ & $\Sigma U$ are both $< -\ln P\mu$ or $\left( < -2(\ln P\mu) \right)$ — but I do seem to need more info than just .19 (& the fact that $\Sigma U_i$ converges so that .. $\Sigma U^2$ converges)

that   $\Sigma \frac{U^2}{2} < -\ln P\mu$ is sufficient !

If I could show that $\sqrt{\Sigma \frac{U^2}{2}} < -\ln P\mu$ (No E) then this implies (by .19) that $\Sigma U_i < -\ln P\mu$ — a much stronger result — But it seems unlikely to be true. e.g. say t. sequence to be extrapolated was a random seq., & u was some simple deterministic seq. ... like $1^{(cos)}$. In this case, I would guess that P(U_i) would be ... constant — is indep of i & neither $\Sigma P(U_i)$ or $\Sigma^2 P(U_i)$ would converge. In this case, howp, $\Sigma P(U_i)$ has to converge — it always converges. So P(U_i) must ↓ w. i', it can't be constant.

Drop this for t. while !                                    → 220.20

β 18:42

space

▷ of 213.37 (R) "Anthropic principal": We ask "Why is 113 + 217 → 330?" An explanation might show how we then predict that this method of to obtain this sum from a table of addition and carry rule." using tables + carry rule will work on other addition problems like "43 + 91", etc.

So: ■ In Math, the tautology, "An explanation" can be a recoding into a more useful form that may reduce cc of comp'n

In general, having an "Explanation" in Math, does help solve future problems — i Math problems always involve "cc".

So, I may need to Extend my idea of what an "explain" is. T. stuff on other of 213.24ff ff & .24.00ff may be very impt.

Hmm! T. Anthropic principal in Physics, → suggests a principal in Epistemology (unite aprip).

This linkage is a kind of "Productive" being since they are very similar "Structurally" (.:. "Analogous"). On t. other hand, if both principals are "content free" it is of not much interest that they are "linked".

in .00 t. explanation of "113 + 217 → 330" suggested a general scheme of doing addition — T. only reason it wasn't "predictive" was that all addition (i all true Math) is tautology.

Drop this for a while. I can mention t. Anthropic Principal at DIMACS-talk, but I don't ● need to go into its meaningfulness (R. I may figure out something by then).

18    5·20·03 } Re: Anth Princ: (In Physics): It could be used to justify any theory, no matter how unlikely? (Including t. "God Hypothesis"). What we are interested in, is not

▷  absolute probabilities but relative probs of various hypotheses. → 2·15·09
         Chemistry,
. In particular, at an earlier state of Physics, chemistry, we would have an enormous list of essentially "Ad Hoch" constants that described t. Universe — T. P.E. of that set would be very low — much

23  lower than that of t. present small set of universal constants.

24  [    Re (20)ᵃ (rel pc's): We predict on t. basis of rel. aprobs of Models:
         Hvr. in t. case of our evaluating any Model of "aprip", we are using t. hypoth to
26  [   evaluate itself. (! ?) ← [This Seems] ↳ Hume's objection to usual Method of induction. ————→ 2·15·08
                So it may be that .18 ff say that t. Anth. Principal / is not needed, i doesn't
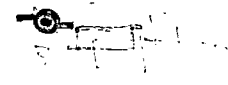                                                          in Physics
     help any Phys. Hvr. .24 - .26 is bothersome! t. application of this line of reasoning
0  →  to t. aprip. Model itself!

         (.20)ᵃ is related to t. "Unlikely ness of Biochemical Evolution ever producing
32   "MAN" "or any similar creature. A MAIN Counter Argt.: Low PC means nothing. Ratios of PC's are what's important!
                                                               ∠.32R
         If ~~~~~~~~~~~~~~~ .18 - .23 is correct (note d'fty of .24 - .26)
     Then this may solve t. Main Problem of Epistemology: "Where does Aprip come from".
     Except, perhaps for .24·26 & .26R in particular.
         A vague suggestion for UPD: That we are born w. this aprip "built in" § parts of it is "things that work well in
     (as are likely to work well in future, if ... t. models have t. aprip." If we have this aprip for models "built into us"
                                                                                         ↳ ← 2·15·00

**NIPS**

(Spac 23)

> 0 (2.12.23): In discussing e need for a priori info in order to learn to walk, talk, and construct speech:

— Also mention people are d priori able to breath air (fish can't do this) — eat certain kinds of food. (mainly milk). They have a set of sense organs that are a priori known to be used in the env. in which they will live. Other creatures have senses that are more or less sensitive in different modalities (color, intensities, types of odors ···).

That a priori for learning is ▇▇▇▇ part of the design of the organism for a certain expected environment — just as is having facilities to breath air to metabolize certain kinds of molecules, to be able to perceive certain sense modalities with certain sensitivities. Sense in U.V., I.R.; electric fields, magnetic fields, ···
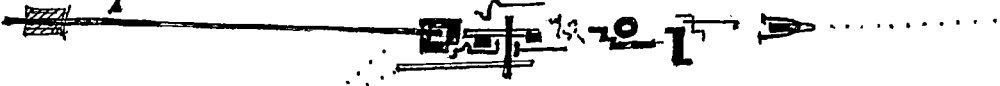
o   Inherited info is called *instinct*: Our a priori ability to transduce certain things very easily is "instinct".

.11   Perhaps mention: for many Statisticians, subjectivity is to be avoided. He wants something that the Statistical Community Agrees on. / Using subjective info may involve extracting it from his client — a difficult and imprecise operation. → while it may be agreed that properly used subjective info can help establish inquiry

I in my own case, I am both the statistician and the client — so this .18   problem does not occur.



> 0   So: Just do the London talk w. folg mods:

21   1) Shorter explann of "incomputability"   write: $P_M^T \rightarrow P_M \rightarrow \omega$
      Give exposition that I now have on website (of London talk)

23   2) In e discu. of subjectivity: Go into it more! 212.15 — .33; 213.00 —.10 (maybe 213.11-18)

24      So how impl. evoln of a priori is: That s intr to how this in my System perhaps Anthropic Principle.

     3) To Section on → Also 212.06 — how QA covers sequential, Bdd, categorization   Ground of discovery
        Alpha: 210.14 —.18  ; Ency. modal'dentification etc.                        .24

27   So I can just use t. web version of London talk as a guide. i. do. 21727
                                                         ↳ 218.00

──── (right column box) ────
.24 | Main Problem w. Anth. Princ.: It doesn't seem to be predictive / "falsifyable".

Is M&Roush's "Predictive" or "Falsifyable"?

M&R is always ("Trivial" Tautologic) if correct.

Concep: "M&R statements are Analytic".

Note would explain M&Roush only because what it is meant to "explain" needs no "explanation".

.37R                              → 214.00

37

NIPS

Backtracking Approach! .34  (needs Work!).

00: 210.18 : [SN]. Perhaps main reason I don't like writing R's is that expousium is such an "Ill defined Problem"!

A possl proni Discuss UPD! "incompatibility"

Subjectivity! How aprop evol of animals, mem, (individual mm evolves over a Life)

Anthropic Principal of Epistenology

Big Q is! How is updating done?  ( How QA this (like Google)

.06 — 2. Phases!  ① + QA problem.  This covers ⓐ sequencial predi ⓑ Beg induction (longest)

ⓒ categorizing Identifn problem, surface reconstructn.  ⓓ Engr Implementations

0  Use of Levels for induction.  ⓔ can it do "clustering"?

So: explain how it does Levels & updating.

Then explain Phase 2.

Perhaps do explain about "incomputability" (also irrelevance of semi-computability)

So start out w. London Lecture started! but more sophisticated!

15  More emphasis on importance of "updating" : One can do induction, prob. solving w/o it, but no real progress

One can get Good results w/o updating, by coding large/diverse corpus, but R's would probby Cost far too much.
Perhaps

So what I may do, is give a picture of how I mg evolves in Man — how it must
20  evolve in Man (& in Life) — (the "Pig Ivng.)

One can start out w. much less aprio info than a child — (but it must be some
minimum of info on how to run) Then if it's small, we learn to put more in to TS Q,
( which is usually good, but it means more try. time & more TS Q has to be
written — ... There is always this tradeoff betw. amount of apri info &
info in to TS Q.

And A way to look at "evolution of Aprpd": One can start w. aprip & no data
("neonato" / tabula rasa.) Then do ALP on entire large corpus leading to present problem / enormous srch.
or start w. do small problem, update aprpd, do next problem, etc. I save speeds out
30  time searching, we will get same results as before. — we lose srch time, not as good a soln,
but acceptable srch time. ( Perhaps Mention "back breaking" to make evoln. of
aprip "More realistic. If one stores many possl/codes, "backtracking" is faster

spec 213.00

34  [SN] BACKTRACKING: HOW TO !  Say I have 3 codes for corps up to n  $\alpha_{1,2,3}$
of lengths $\frac{\ell_1}{2}, \frac{\alpha_2}{2+3}, \frac{d_3}{2+4}$! For coding out to n+1, try continuations of all 3 $\alpha$, but
give relative wts $1, 2^{-3}, 2^{-4}$
[This can perhaps be geared to backtracking for Peerback. ← Not obvious ! Maybe false!
Needs work ! Huri looks like its worth spending time on !

N $\psi s$

$$\boxed{E_\mu \Sigma p(u_i) < -2\ln p_\mu}$$   Bounded $E \Sigma f(u_i)$

By .15–.19;

$\infty : 208.40'$

Similarly; $\ln(1-u) < -u + \frac{u^2}{2}$

$\Sigma -\ln(1-u) > \Sigma u_i + \frac{u^2}{2}$
$\Sigma u_i - \Sigma \frac{u^2}{2} < \Sigma -\ln(1-u)$

$\frac{p_\mu}{\prod(1-u_i)} < 1$   normzn const.

$\frac{1}{\prod(1-u_i)} < \frac{1}{p_\mu}$   $\prod(1-u_i) > p_\mu$   $\therefore -\ln(1-u) > u - \frac{u^2}{2}$

$\Sigma \ln(1-u_i^2) > \ln p_\mu$;   $\Sigma -\ln(1-u_i^2) < -\ln p_\mu$

$\Sigma u_i - \frac{u^2}{2} < \Sigma -\ln(1-u_i) < -\ln p_\mu$

so   $\Sigma u_i - \frac{u^2}{2} < -\ln p_\mu$   $\therefore E\Sigma u_i - E\frac{u^2}{2} < -\ln p_\mu$

$\geq 0$   $E \Sigma \frac{u_i^2}{2} < -\ln p_\mu$   $\underline{E\Sigma u_i < -2\ln p_\mu}$

.09   So   Theorem:   $\boxed{E_\mu \Sigma u_i < 2\ln p_\mu.}$   What intuitive significance is this?   $\sum_{i=1}^{n} \frac{1}{i} \sim \frac{\ln n + \gamma}{-\ln\frac{1}{n} + \gamma}$ ← Euler's const.

0   old Thm.   $E_\mu \Sigma \frac{u_i^2}{2} < -2\ln p_\mu.$   m ) matter to Prob V. at $\approx 1$.   if $p_\mu \approx \frac{1}{n}$, this would do it!

so $E\Sigma u_i$ & $E\Sigma u_i^2$ both have same upper bound of $-2\ln p_\mu$.

These 2 Thms can be readily generalized to 'p_\mu being a function of n, the sequence length

.15   To show $(\infty)$   $\ln(1-u) < -u + \frac{u^2}{2}$:   $f(u) \equiv \ln(1-u) + u - \frac{u^2}{2}$

$f(u) = 0$   $f'(u)\big| = \frac{-1}{1-u} + 1 - u = \frac{-1 + (1-u)^2}{1-u} = \frac{-2u + u^2}{1-u} = \frac{u(u-2)}{1-u} = f'(u)$

$f' = 0$ at $u = 0$.   Consider $0 < u < 1$   $f' > 0$.   $\therefore$ for $0 < u < 1$, $f(u) > 0$,

.19
20   since also $f(0) = 0$, $f'(0) = 0$ & $f'(u) > 0$ in that interval. (i.e. it can only ↑), it must be $> 0$ on that interval

22   Summarized:   $f(u) \equiv \ln(1-u) + u - \frac{u^2}{2}$ : thus,

out. interval $0 < u < 1$:  ① $f(0) = f'(0) = 1$;   $f'(u) = \frac{u(u-2)}{1-u}$   which is $> 0$ on that interval from .19

Note: if $\boxed{E_\mu \Sigma u_i < -2\ln p_\mu}$ then $E_\mu \Sigma u_i^2 < -2\ln p_\mu$ since $u_i^2 < u_i$.

So Prob 3 is major result.

→ I should check on Prob 3, hmm, the juggling of inequalities is a bit tricky!

This is a re derivation.

→ T. beginning at $\infty$:   $\frac{p_\mu}{\prod(1-u_i)} < 1 \to \prod(1-u_i) > p_\mu$ : $\Sigma \ln(1-u_i) > \ln p_\mu$ : $-\Sigma\ln(1-u_i) < -\ln p_\mu$. ✓   ← normzn const   α

50   Then via .15 ff.   $\ln(1-u) < -u + \frac{u^2}{2}$;   $-\ln(1-u) > u - \frac{u^2}{2}$ :   $u - \frac{u^2}{2} < -\ln(1-u)$; $\Sigma u_i - \Sigma \frac{u_i^2}{2} < -\Sigma\ln(1-u_i)$   β

from α, β   $\Sigma u_i - \Sigma \frac{u_i^2}{2} < -\Sigma\ln(1-u_i) < -\ln p_\mu$

so   $\Sigma u_i - \Sigma \frac{u_i^2}{2} < -\ln p_\mu.$

$\therefore E_\mu \Sigma u_i - E_\mu \frac{u_i^2}{2} < -\ln p_\mu.$

$E_\mu \Sigma \frac{u_i^2}{2} < -\ln p_\mu$   i.e.   $E_\mu \Sigma u_i^2 < -2\ln p_\mu$

so   $\boxed{E_\mu \Sigma u_i < -2\ln p_\mu}$

3 disn facts to give derivation:

① $E_\mu \Sigma \frac{u_i^2}{2} < -\ln p_\mu$

② $\prod \frac{1}{(1-u_i)} \cdot p_\mu < 1$

③ $u - \frac{u^2}{2} < -\ln(1-u)$ (if $0 < u < 1$)

conclusion: $E_\mu \Sigma u_i < 2\ln p_\mu.$

→ $\boxed{215.15}$

202.20: Notes for Paper on Convergence Params.
for Chris Wallace "Festschrift."

4 Mor
2 afternoon
5 counb!
late aft.

∴ **On t NIPS talk**: The idea is to Motivate people to read t. paper!
T. abstract should do much of this.

What I may do is start by showing that INV & OZ & predn. (CQA) problems cover
just about all ~~~~~. Perhaps Give a few interesting/general examples.

Consider what we mean by "well defined problem".

　　1) ▨ Usual inversion problems
　　　　a) Inversion problems w. Badness threshold (how far from exact soln can one be?)
　　e.g. solving complicated ~~~~~ Eqns.
　　2) ▨ Open: List several types. Maybe time varying, Noisy, Expensive Generator.

Perhaps show how a clever TM could help one "well define" a problem (perhaps later — Letter to Ivon on this)

1 Book/day
= 1 h books/msys
3 k /decade
(15 k in 50 yrs)

(compare to G. Report)
So I want to show that this is a v.g. way to do A.I.

14　Argts: 1) ALP is a good, very general way to do induction
　　2) Just about all well defined probs are INV or OZ (OZ covers ≥ 05)
　　3) T. Main problem is "updating" & I think I have a good general way to do it.
　　4) 2 phases of updating: Phase 1 simple  Phase 2 Advanced.
16　　5) Use of (such) Pi a (Better) WON.　　——→ (212.00)

─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─

20 [SN]　In Prediction of [reals]: We want to do Searchover /continuous & discrete params at same time.
21　(not "discrete, then optimize" — then new discrete, then optimize.)
Contrast speed at [.21], with discrete trials & associated random trials
for contin params. In this coding, we will have a "Goodness of fit" criterion
for how good t. code is, so we don't have to hit t. "A" value exactly.

.25　There is some Q about whether this system will work, because
.26　it does not do ∑ of pc's of codes; it just tries to find v.g. individual codes
　　There may be some special way I could devise to deal w. t.
.25 - .26 effect.
　　(On a previous look at something like this:) T. trials are not in pc order — so
30 → not so good for t. search for induction. I don't remember whether I found a good way to deal w. this diffy.
(even for discrete predn.)
This is a more general problem that occurs when we do "lossy coding" w. "correction".
We do not know, apriori, how big t. "correction" part of t. code will be. — So this is usually not such
a good method to do "Induction"

Nips ......

On t. Superiority of ALP: Reply to "Further Experimental Evidence against utility
of Occam's Razor" Journ. of AI research)

$$\int = p(w_i) \times p(c(corpus|A)_i \quad (1996) \quad 397-417$$    in PS folder

6/13/02

**Q:** Suppose $PEM_1$ gets more wt. for a certain corpus than $PEM_2$; but empirically,
$PEM_2$ gets better preds? Well, that means that $PEM_2$ gets better p. h. y. on
p. of corpus than $PEM_1$, but t. apriors are much inverted. It could be that
aprior's for $PEM_1$ à $PEM_2$ are way off (i.e. very bad computers (avg.) or it could
be that $PEM_1$ is very a.H. à and actually will not continue to give pred'n.

**or,** My system could be **wrong**!!

Most likely that .02-.03 is true, that 2 p. h. y. s were way off. If $PEM_1$ is
really better à longer corpus will fix wts properly.
When .03-.05 (à above) is true, one can usually tell, by how t. 2 p.h.y.s
were created.

**So :** If corpus is **short** enuf, main wt. is from apri, à can be **way off**, if apri info is poor!
Say t. "true μ" has a very long dern. The only way one could justify very it, would be to
have a long corpus. With a short corpus, one would not even consider that **complex μ**.
One **might** do it by using/many continuous params! in which case "a.H." ness could occur,
à be apparent. — One would tend to select **wrong** μ — A.H. **correct** for that corpus.

It might be good to give examples of corpi w. continuous params, in
which one uses a μ that was too large for t. corpus à Say t. "true μ" was too
(i.e. "too complex")
**large** to be dscvrd by t. corpus!

.19-.20 could easily be t. cause of .00-.01

_____

T. foregoing is related to a paper that I have (in PS I guess) on some guys that tried
using MDL on various corpi à found it worked well on some, but poorly on others.

_____

If μ has a long dern, than a /**too short** corpus, will always give more wt. to params w.
Shorter derns, that will not do as well as M. — But M is essentially "undiscoverable". If one does
μ) One **must** have reasons for choosing μ, à if one does, it's a prior is very hy.     Choose

_____

Re: t. paper at top of page! If they had (wo. seeing new data) modified their corpora algms.
in ways that they that would improve it — either on t. basis of logical reasoning &/o/their
coherence in t. past, — then they were giving apriors of their model (≡ it "complexity").
If they had no reason at all to use these models, then there are /**very** /many other very bad ones
they should have tried.

If they feel that modifying à model increases its "complexity" they were thinking of
complexity" in à purely syntactic sense. This is **Bad**. The ...
Occam's razor is rite only if t. lang. has been modified so that Occam's /found mt. post,
razor was correct in t. past. This is done by giving useful prediction rules, shortcodes.
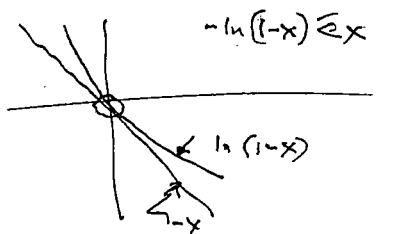If Occam's razor doesn't work, it's usually because t. apriod or language was not properly "updated"

Attempts to find upper bnd for $\sum_i^2 p_c(u_i)$
— we know it converges to $< \infty$.

$\dfrac{\text{note}}{\ln 2} \equiv k'$ chol cost $\ln 2$

: $\boxed{SN}$ : $\prod \dfrac{1}{1-p_c(u)} = \prod \dfrac{1}{1-x_i} < \dfrac{1}{\boxed{p_c(u)}}$
$\equiv z$

$-\sum = \ln(1-x_i) < -\ln z$

$z \ll 1, -\ln z > 0$

$-\sum \ln(1-x_i) \lessgtr k'$  $\left( \equiv -\ln \dfrac{1}{p_c(u)} \right)$

$\le \sum x_i$

$\ln(1-x_i) = -x_i + \dfrac{x_i^2}{2} - \dfrac{x_i^3}{3} \cdots$

$> -x_i$

$\sum -\ln(1-x_i) \le \sum x_i$ $\left( i.e. \quad \ln(1-x) \ge -x \right.$
$\left. -\ln(1-x) \le x \right)$

$-\ln(1-x) \le x$

So while $a(\infty)$ falls vs. $\sum x_i$ converges, it gives no upper bnd for $\sum x_i$ .

$\sum x_i \ge -\sum \ln(1-x_i) > k'$



$-\ln(1-x_i) < x_i$

$\prod \dfrac{1}{1-x_i} < \dfrac{1}{p_c(u)}$

$\sum -\ln(1-x_i) < -\ln p_c u = k' > 0$

$\ln(1-x) > -x_i$     $-\ln(1-x) < x$     so $\boxed{\sum -\ln(1-x_i) < \sum x_i}$

hvn, we do however expected upper bnd for $\sum \dfrac{x_i^2}{2}$ which is $k'$, from $p_c(u)^2$ conv. Bren.

So $\sum x_i - \dfrac{\sum x_i^2}{2} + \dfrac{\sum x_i^3}{3} \cdots < K' \sim \dfrac{k \cos t}{\ln p} \sim 1.4k \cos t$.

$\dfrac{\sum x_i^2}{2} < k'$ or $2k'$ or whatever $\quad < $ for $z$ seq. of $n$ data pts. from plus $z \cdots$ I don't know if this is $k'$ or $2k'$.

So $\sum x_i \underset{\approx}{<} 2k'$

$\dfrac{1}{3} + \dfrac{1}{5} + \dfrac{1}{7} \cdots \dfrac{1}{2n+1}$

$\to \dfrac{1}{2} \ln n$.

If $\{x_i\}$ is a positive, non-$\uparrow$ seq. and $\sum x_i^2$ has known upper bnd.

$\sum x_i < k'$, conver. by

" $\sum x_i$ is known $< \infty$ .

Conver. by any pty. $\ge 1$   if $\sum$

Up to .16 I hadn't been able to get a bound on $\sum^2 p_c(u) \cdots$ only that it converges,

$-\ln(1-x) = x - \dfrac{x^2}{2} + \dfrac{x^3}{3} \quad < x$

is this $\ge x - \dfrac{x^2}{2}$ ?  i.e. $-\ln(1-x) > x - kx^2$ ?

No! for $x$ betw. 0 & 1 the poles of $\ln(1-x)$ will dominate any $x^2$ (or finite power of x).

a parabola thru 1,0, having matching first & second derivs. of $\ln x$.



Drop this! It's an interesting puzzle, but not all critical to TOM !

$\to 2k \cdot 0\beta$

Why $\frac{x}{x+y}$, $\frac{y}{x+y}$ is BEST Normzn possl. .10 → .19

(Convergence of $\sum p_i(u)$ .22)

$\sum_i \ln(1-q_i) = \Rightarrow -\epsilon_i + \frac{\epsilon_i^2}{2}$

if $\sum \epsilon_i$ conv. then $\frac{1}{2}\sum \epsilon_i^2$ conv. Hur....

00: 206.40! ▓ higher include t. no. zero.

4) Monotonic $\beta$ funct of s.c. is s.c. But $\sin(s.c.)$ is s.c. only if $(L=B)$ is all on an $\beta$ section of $\sin(\cdot)$. $(s.c.)^2$ is always s.c.

5) $s.c. - s.c. \geq \frac{s.c.}{s.c.}$ are usually (almost always, always) yes s.c.

∴ T. result is that s.c. uos. are not very useful (for arithmetic! — so sometimes one can use R's property in proofs

→ 6) "enumerable" is not a good way to relate s.c. uos., because when one says a certain set of uos. is "enumerable" its not clear whether one means cardinality or computability!

7) Say $M(x,i)$ is a semi measure, let $\vec{F}(a,b)$ be a "normzn", if $A+B=1$.

$\vec{F} = F_1, F_2 = f_1(a,b), f_2(a,b) = A, B$.

Then $F_1(\overbrace{M(x_i=0)}^{A_i})$, $F_2(\overbrace{M(x_i=1)}^{B_i})$ will be a seq. of real parts.

D·16·04 [I think this ∧ is backwards! $M$ is already normzd; we want to normalize $M$, t. univl.d.f. → Nvr. this doesn't affect the argument

for each $i$, there will be $\rlap{▨▨▨▨}$ a maximum value of $\alpha_i \Rightarrow A_i \geq d_i M(x_i=0)$ and $B_i \geq d_i M(x_i=1)$

The normalization function that gives t. largest $\alpha_i$ $\rlap{▨▨}$ will be t. usual normalization —

i.e.: $f_1(x,y) = \frac{x}{x+y}$; $f_2(x,y) = \frac{y}{x+y}$.

8) If t. theorem of .15 is used, then t. $E \sum (c_i)^2$ before t. normzd and $\rlap{▨▨▨}$ (t. normed measure that generated t. data) will be minimum.

i.e. t. ratio $\frac{M_{sol}}{M}$ will be as large as possl. — using any other normzn will ↓ this ratio. → See 4TM 313.10-20 for more "proof".

9) Usually, for a long seq., t. probty of "u" will be very small; $\sum (p_i(u))^2$ converges is bounded.

→ 10) If $\mu$, t. generating dist., is normzd i.e. a "Measure", then t. normalization contour

for M cont $\overbrace{b, \rho}^{} > \frac{1}{p_i(u)}$, because this would make $p_i + p_0 > 1$.

(t. normzn constant is $\frac{p_i(u), p_i(0+1) = 1}{p_i(0) + p_i(1)} = \frac{1}{1 - p(u)}$ for a single term)

N.B. t. $p_i(\cdot)$ are all Conditional Probs

t. product is $\prod_i \frac{1}{1-p(u)}$ which converges ⟺ $\sum p_i(u)$ converges.

so, a fortiori, $\sum p_i(u)$ must converge.

i.e. is bounded ... so $\prod (1-p_i(u)) > 0$.

we also have then that $E \sum p_i(u)^2 \leq \ln \frac{1}{p(u)}$

GO OVER (10) carefully! .22-.23

t. reason .26 must converge, is that if $M$ is a measure, t. normalization constant for it must be $\leq \frac{1}{p(u)}$ : see .22-.23

[ At first glance, it would seem that t. convergence of $\sum p_i(u)$ could not be related to $\mu$ —

[ But it is! $\mu$ generates t. sequence that decides under what conditions a u should occur!

More exactly $\prod_i (1-p_i(u)) \geq \frac{1}{p(u)}$ = max size of Normzn Constant.

NIPS

On Unimportance of semi-computability. 33 ff
Significance of P(U)>0 (.00 - .30)

: [5N] What is signif. of P(U) >0 ? : for Sun or HR predn. it would be meaningless.
Because there is no info (we know of!), n+ data to imply "end of seq".
For/daily weather data also, I suspect that "U" would be not acceptable.
There is a tendency to interpret "U" as "end of the Universe". I think it isn't ½ a good interpretn.
A few pp ago I noted that ( I looked back to 194 (4.24.03) and didn't find it ):

If I had somehow 2 kinds of info about a t.sq. ① I knew last 100 bits ② 
I knew say. Had no "U" for last 100 bits. Then I should give very low pc for U 
being next bit.

The U symbol can be considered as being one ∋ symbol alphabet for Bernl. predn.
Using Lap's rule, the U always has ≠ 0 probly of occuring. Modifying Lap's rule (because
we have aux info) we get very small pc for the 2nd symbol. ∴ Normally in statistics,
if an event has never occured before, we usually/may want to give it some possy of occurance.

( particularly if it is an "alphabet" "symbol"( ! ) —

Normally ALP will give, for short seqos, a fairly large pc for U,
but then P(U) → 0 as maybe 1/n (or converges) as for many/data seqos.

Laplace gave(?) example of (365×10⁴)⁻¹ as pc of sun not rising tomorrow,
because, he that earth was ~ 10⁹ yrs old — the sequence has only 10⁶ bits.
This corresponds to "aux info". I can't say, "U" means the seqo will stop
at that pt, not that the Mkt will close for a days open the next day.
This (occurs during war, plagues + conditions of Natl. emergency.
— The kinds of info that ALP is using as a basis for the "U" prediction,
is mainly, that the seq. is short. T. fact that P(U) usually converges
faster than 1/n means that aux info about past lengths of seq.,
can be used to suggest, corroborate low expected pc for U.

T. way we do this — by renormzn — seems a bit A.H., but it's a reasonable
way to do it. ... ▨▨▨▨ The suggestion that we consider p(U) >0 as

a real possy. would have very poor corroboration in t. past.
E.g. If we start w. 10 yo/daily/SM data: At first we get frequent occurances of large
p(U)'s. They are always wrong. As t. seq. continues, p(U) ↓ rapidly to
say ~ 1/3600 ~ .0003. ~(.03%). (its .003% if we consider 100 yrs of SM data)

So: All of t. foreg. suggests that S. normzn should be used in usual prediction.
If we have a (meaningful)/useful/interpretn. of "U" so that P(U) >0 is a useful, reasonable
pc, we should use Lap's normzn.

▨▨ On: [203.00 to here]: 1) First that s.c. v.s. non-s.c.(but "knowable have a limit") — s.c. is it
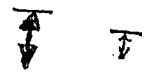signifly much better, if one is looking for an approxn. of t. number.
2) s.c. nos. are characterized by a) a seq. of monotonic(incr.(none) vals. b) an upper bnd c) t. limit of a).
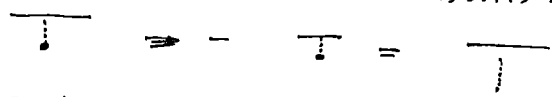the quantity B−L will be the smallest amt. of error one could know in L.
3) sum of s.c. nos are s.c. products are s.c. ...  usually A t. ranges (B₁≠L₁) · (B₂,L₂)

While Mathematics is
"T. Handmaiden of t. Sciences", she certainly
has a life outside that
pit ( occupation ).
on the other hand, she
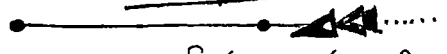really is employ during
by t. sciences.
"No all work and no
plays she".

∞:

: If we take t. distance (or radius) of ≥ sc nos., t. say has a limit, but is not a monotone seq. The "limit" is known,

The resultant/non-monotone seq. has a limit: It has a. known upper bnd, but no known lower bound. From a practical pt. of view, the s.c. value of a semi measure is not knowable in a useful way. Say we want t. semi measure of (. seq. 0110!110. 8 bits. There are 256 possb seqs. of length 8. # finding as many codes as possb for all seqs. of length 8, will narrow t. upper bnd. on t. pc's of all of t. seqs. If we include in P on codes that stop before they print 8 bits or codes known to loop before 8 bits off, we have a corresponds smaller upper bnd for rc of say 8 bit seq. If Δ is t. total pc of all codes that might be a code for an 8 bit seq, but don't ever converge à are of "unknowable" convergence, then Δ is t. uncertainty in t. actual pc of any code of an 8 bit seq. If its true pc is p then we can at best know its betw p ≥ p+Δ.

My guess is that usually Δ ≥ p (à more likely ≥ Δ >> p) so Δ gives no useful bound on t. pc.



So t. moral is, that for actual practical approaches, t. known limit is too far from t. known pts to be of much use in getting accurate estimates. In most cases, à for either semi computable or more "totally incomputable" numbers, one will usually look at t. apparent rate of convergence à extrapolate to ∞. —— if one wants that info. More likely, for practical statistics, it's not necy to know how close one is to t. limit (≡ ALP). Even if one knew one had $P_M$ exactly, we would still have to estimate how accurate it was, in t. "usual ways" (cross validation, with or w.o. training set).

19.——

23

T. subset of "U" cases that are "undecidable" gives a low lower bnd on t. precision of any estimate of a semi computable no. My impression is that normally one's estimate of this "undecidable"; à lower à upper bounds on what could be—are too poor to be of use in knowing à estimating the semi computable no. of interest.

In getting t. "non enumerable" nos. that are differences or ratios of s.c. nos., We really don't have much defence in our ability à estimate accuracy of correct approxn. In this area, t. estimated seq. of accuracy is not monotonic, à t. bounds on it are too far away to be of utility. —— We estimate error by apparent rate of convergence -- - ``

But anyway, t. remarks of .19—.23 applys — i.e. this is usually not a simple problem in ALP.

So: Motivation for Normized ALP: If M is a cpm it will give more accuracy to t. estimate of $P_0, P_1$. Also true if M is "approximatable" by a cpm. I'm not sure about cases in which M is simply a sc. semimeasure.

N/es

8) Levin's analysis (4 & {20.2 #}) is relevant to ⓐ normzing ⓑ semicomputable M.

$(200.21) \leftarrow ?$

Must reason switch as a nonzer constant "

wierd

I think it is only of interest if the normz constant $\to \infty$.

9) T. "natural" way of normzn gives max ↑ of "normzn factor" ∴ it gives least MSE (error)?

Say $P_0$ & $P_1$ are PE's so $P_0 + P_1$ maybe $< 1$. $(P_0; P_1 \geq 0)$.

we for normzn, we want $P_0, P_1 \to x_0, x_1$; $x_0 + x_1 = 1$, ②

say $\alpha$ is t largest no. ∋ $x_0 \geq \alpha P_0$ & $x_1 \geq \alpha P_1$, and ($\alpha$ is a max mult factor of $P(u)$.) Then T₀ $x_0, x_1$ mapping that gives max $\alpha$ is $x_0 = \frac{P_0}{P_0 + P_1}$; $x_1 = \frac{P_1}{P_0 + P_1}$

This is t. Basic Normzn, indep of (8) (Levin's analysis).

10) A possi. interpretat of Levin: That even using my normzn, There is no normzd UPD, that is which $\geq$ all other normzd UPD's (within a constant factor). (w.c.f.)

If all normzn constants were bounded, -11-12 would be false, ∴ there must be sequs's times for which t. normzn. const" $\to \infty$. (This is a necy, but not sufficient, cond for be. to be wise about this. If we assume M is a CPM, then Any Normzd UPD is [w.c.f] better than any other within constant factor

11) If we knew t. seq. has been running for $10^6$ bits, (but we only have last 100 bits for data) then we should normz our predictions to $1 - 2^{-6}$ rather than to 1

20 §N 3 d.fs:

1) from $-\infty$ to $+\infty$ : $\propto e^{\frac{(x-\mu)^2}{2\sigma^2}}$     Gaussn.

2) from 0 to $\infty$     $\propto x^n e^{-mx}$     Gamma

3) from 0 to 1     $\propto p^n (1-p)^m$     Beta (?)

These 3 cases would seem to cover most situations!

12) For discussions related to this stuff see $\left(\frac{L-V}{97}\right)$ : pp 259 - 305.

13) Re: semi computable nos: a "s.c. no. ("enumerable") no. is defined by an incr. increasing (non-↓) seq. of reals ∋ an upper bnd. (or an upper bnd + an algn. to generate such a seq.)

One can tell 2 s.c. nos ∋ Eat t. s.c. is closed under, a) addition b) Mult by positive constants, or by another s.c. that is known to be $\geq 0$. c) any monotonic (non-↓) operation.

$(sc)^2 = sc$; $\sqrt{sc} = sc$ if $sc$ is known to be $\geq 0$. Subtraction & division/s betw. 2 sc's are never poss.

Any way, t. point is, one can't do very much w. s.c. nos. (usually one can't compare them — Pro

A regular real is a s.c. no. & often it can be compared to a s.c. nc.

If z is t. true / value (limit) of a s.c. no. & b is it's upper bnd, then b - z is t. lower bound for uncertainty in t. value t. no.

50 : On "Superiority" of Normed Univl pd. = Universal Prob Measure.

I had this idea that Levin proved that there was an Universal Measure that multiplicatively dominated all other Universal measures.

[ L.V. '97 (second edition) (p 300)(p 301) 4.5.4: Possibly §4.7 pp 307-314 (history) is relevant.

[ Seems to contain no mention of Levin's proof.

The BONE of Contention: ~~Defines that~~ ( unival / semimeasures & measures both have ∈ Σ errᵃ

converge to < = constant — here, the constants differ by a normzn factor for t. unive. measures.

Th. Q is: How large can this factor be? — Can it be infinite? I think it can't be infinite —

Woops! The normzn "constant" is not + a single number for t. entire Normed PD! —

0 Its usually different for every bit in its domain!

Say P_M actually was ≠ αM for t. say, of interest, X. Then t. normzn constant would ~~be about~~

1/α az → ∞, but ~~be~~ be 1/α. for all ~~other~~ bits of X.

So 1/α would be upper bad on normzn constant — so + t. Normed Unit d.f.

~~would~~ could be 1/α times as good as t. semimeasure: ( So ~~+∈~~ Σ errᵃ could be

.18 ◆————————■—————1/2 P(x) semimeasure! ←── N α .18 The "Best" I'm considering, gives Normzn constant 1,

1) In general, this normzn constant has to be bounded by 1/α. | equivalent to (==1 for M)—
                                                                | this gives Σ errᵃ = φ!
2) I am considering only production of EPM's (computable prob measures). | So while normzn ↓ agree, we have no
                                                                          | useful idea abt HowMuch.
3) re: ① My impressions that usually t. normzn cnce will be << 1/α, since the normzn
  ○ constant only fills out the "U" probabilities.

2) ᴺᐧᴮ* 4) What Levin probably showed was that if you considered all possl. normzn methods, there was no
  one that was ~~within~~ "better" (with in a constant factor) than any other. This may be false!
  It implies that some normzn constants must be unbounded. ! If so, the SZ&S corollary would
  say they couldn't be useable probability distribos — i.e., SZ&S wouldn't apply to them.

5) It may be that t. arguments mainly about predicting non cpm's.

→ 6) Actually, only t. usual method of normzn has a normzn "constant", i.e. t. pc's above
  ▦ for θ ≈ 1 ▨ are mult by t. same "constant". (that t. constant ↑ (usually) as t. seq. length ↑.

  [SN] If a monotonic funct. A common ~~inverse~~ funct.
  Well, A⁻¹ is n't defined for many of its args. Hvr, if, in "defining" t. inverse of t. Ack funct.
  one lists t. args for which it is defined, then computing t. inverse is trivially prim. rec.
  To find A⁻¹(k) ~~simply that A's args for which~~ simply list the # / args in order of size.
  t. n'th one will be K; so A⁻¹(k) = n.

7) In view of .18 K : ≤ errᵃ ≤ b but we don't know how much; (φ is our bad! ①).
  If we compute t. normzn "constants" we will know how big it is & how much Σ errᵃ has ↓.

∞ : 3 interesting cases of trying; that seem ~.

1) Mult÷ division needs to take much longer time than Addition Subtraction: This helps TM to discover "laws of Alg", since they can shorten time needed to do many Arith operations.

2) Discovery of concept of "option" as needed for certain prob solns as being faster (shorter) than simply trying to copy examples of "option" that were given in the TSQ.

3) From examples of very simple stories in algebra, TM invents humans w. very simple properties. ( A, B, C, : Alice Bob, Charles, Dick and Jones···ect) later TM sees stories about humans in R_w. → makes models.

A great breakthrough in compression occurs when T realizes "it" can use same / or analogous models for both Domains

A poss. case in which they are "~": Restriction of a TM (i.e. HW or for calcul. access to Mems, APL; to proof S's, etc) or info in TS.. can make possi. for certain concs to be useful that would not otherwise be useful. This fits examples [1] & [2]: I'm not sure about whether Example 3 fits.

4) TSQ's: Carefully designed at first, then less care as time goes on.

Energy: where volume.

Also perhaps that: form is productive (not too bytes of compert however to be lower, that its likely few [likely very small...] [it is normal].

20₀ Paper on "Convergence Thms": Some things I can include: GOOD!
>> Gauss of UPD each have their Conv. Term! So Energy & Gauss's Great. Conv. Terms. Usually the proofs will be about the same.
Include Multidimensi, missing data, Multidirectional {maps, some bag elements have finite lengths, others semiinfinite, other z or (grau) doubinfinite. → also E Σ P(ω) < z info
Include: Conv. of Σ(P(ω))². → Also all corollaries corresponding to Sol 78 T3 corollary corresponding to sol 78 T3 corollary,
Maclaurin C. Wallace correspondence as implying gauss of bags & Gauss.
to Multidim, Multiinfinite corpi.

Monthly Town Meas discovery supporting this Corollary

Give each Conv. Term in detail: Give proof (library: if not obvious or if not given elsewhere). tell where proof is.) Explain what kinds of problems it applies to: Perhaps give examples.

30 Perhaps discuss signif. & lack of signif. of incompatibility of UPD(s).

6/8/03 A poss. way to write paper: write /introduction: But mainly, write up a minimal paper to start. Then add in more material from .20 → 30. Then rewrite introduction;
finally: write Abstract.

I seemed to have promised to get paper to Dave in final form by "Jun 2005" Not clear when in June; try June 30. 22 days = 3 wks! I better get started in more detail!
→ In intro, Give Motivation, importance of conv. Terms.

/NIPS

Let me explain: Many years ago in ancient Greece, the Pythagoreans discovered that $\sqrt{2}$ could not be expressed as the ratio of two integers. It took the mathematical community many centuries to get a good understanding of this problem, but well before that time, approximations were made and used. None of the approximations were actually $\sqrt{2}$, but they got arbitrarily close.

In the case of the UPD, we can make a sequence of approximations and just as for $\sqrt{2}$, the approximations will eventually get arbitrarily close to the Platonic Ideal ... the true UPD. The difference in the two situations is that for $\sqrt{2}$, for each approximation, we have a good upper bound on how large the error is. On the other hand, for our approximations to UPD, we cannot ever know a useful upper bound on how much we deviate from the ideal UPD.

Fortunately, for almost all applications, we don't need this information. What we usually want to know, is not how close our approximation is to the ideal, but rather, how accurate is our approximation for prediction. If we have a reasonable sample size we can estimate this accuracy by cross validation, but often we can do better. If the approximation we use is entirely a priori (devised before the data was known), we can use all of the data for testing, since none is needed for training the model.

Though the incomputability of the UPD is usually not relevant to problems in practical prediction, it is of much interest in the Philosophy of Science.

Many scientists are repeatedly disturbed by the need to revise their understanding of their sciences. They look forward to a "Final Theory" that will put an end to all revisions. However, the incomputability of the UPD assures us that this cannot ever happen. With any amount of data and finite computing resources, we can never be certain that we've found the ~~Best, the~~ Final Theory.

Some of us are not at all disturbed by this state of affairs, but find it instead to be a neverending source of joy in discovery. Newton ~~said that his discoveries were built~~ on the shoulders of giants. It is often necessary to ~~progress by~~ climbing over the ashes of dead heros.