Number Bases: <u>Non-Integral</u>:

Say we want to use t. base, $\pi$. T. representation will be a string of integers $< \pi$, $\geq 0$; i.e. 1,2,3.

"32103" means $3 + 0 \pi^1 + 1 \cdot \pi^2 + 2 \cdot \pi^3 + 3 \cdot \pi^4$.

"3.2103" means $3 + 2 \cdot \pi^{-1} + 1 \cdot \pi^{-2} + 0 \cdot \pi^{-3} + 3 \cdot \pi^{-3}$

To convert a no.; N;

.08    Find largest power of $\pi$ $\ni$ $\pi^n \leq N$. If it is =, $N = 1 \cdot \pi^m$.

If $\leq$ ~~find t.~~ find t. largest $a_{n+1}$ integer $\ni$ ~~$N \geq$~~ $N > a_{n+1} \pi^{n-1}$, $a_{n+1}$ is t. first "digit" (msd)

Take t. remainder $N - a_{n+1} \pi^{n-1}$ and get t. next msd via .08. Loop until remainder

.11    is zero or is small enuf to satisfy t. applexn. needed.

~~vaguely foot first less~~

Using t. algm. .08, t. conversion is, I think, always unique, if th. base is not an integer. I'm not sure under what conds. (if any), one is assured of a finite conversion time. e.g. it may take an infinite amt. of time to decide if, say $\pi^3$ is $>$ or $<$ or $=$ to $N$, if $N$ is of arby precision. This will happen infreqly, hvr.

The system is <u>inefficient</u>, hvr. To represent an integer $N$, ~~always~~ usually needs more digits to base $\pi$ than to base 4 (tho. in bo'82 cases, only t. digits 0,1,2,3 are used.

Integers are ~~usually~~ (usually) represented as <u>non-terminating</u> expansions — but not always. ~~e.g. consider~~ ~~$3 \times^2 + 2$~~

e.g   $3 x^2 =$ integer   has a soln. for $x$ that is a non-integer $> 3$.

$3 \cdot 3^2 = 27$, $3 \cdot 4^2 = 48$   so   $3 \cdot x^2 = N$ has a soln betw. 3 & 4 if $N$ is betw 27 & 48.    or even "$3 x =$ integer" can have non int solns for $x > 3$.

One thing I was thinking of using these non-integral bases for was CPI. By using the shortest decrn. length /as proby index # for t. string, we get $2^{-N}$ for probty if base 2 is used.

{ Using other bases is o.k. for number radixes, but how do we use non-integral bases for <u>derns of strings</u>?

E.g. using $1.1_{10}$ as a base, gives us probty values 10% appart.

<u>Negative bases</u> are of some interest. to use base $-3$;

~~$2 \cdot (-3)^2$~~ $2 \cdot (-3)^0 + 1 \cdot (-3)^1 + 0 \cdot (-3)^2 + 2 \cdot (-3)^3 =$

$2 \quad -3 \quad + 0 \quad -54 \quad = -55$

I don't ~~t~~ know if it's more or less efficient than other bases. Using a positive base with a ~~negative~~ sign symbol looks like a Mixed base. Arith may be simpler w. a neg. base (say $-2$). I think there were some recent papers in Computer, C ACM or t. IE3 comp. kens on neg bases.

Number Bases : Negative.

One apparent advantage of neg bases: ~~In addition~~ When adding 2 nos., the carrys do not propage. This is because a carry is always a _borrow_ from the next higher digit! _Woops!_ If one borrows from zero, this _can_ propage!

_Anyway_: This type of arith _may_ have advantages.

~~Adv. Analy??~~
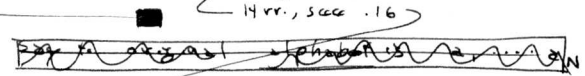
coding w. binary defns! This is a more rigorous, exact method of doing wtk was derbd in

## I i. c II :

~~▓▓▓▓▓▓▓▓~~    I will derb. a method of assigning probys to sequences.

T. codings details, I will not go into, but refer to t. inverse Kraft Thrm.

~~▓▓▓~~ Apriori info. gn:   ① No. of symbols in t. original (raw) corpus. $\equiv N$

.10    ②   Aprip, ~~▓▓~~ $\delta$, of a ~~▓▓▓▓▓~~ punctuation symbol indicating a defn. is to be made.   (— This is to be chosen by t. "coder": I think it is not a critical choice.)

$\blacksquare$   ⌐14 rr., see .16 ⌐

~~▓▓▓▓▓▓▓▓▓▓▓▓▓▓~~ N

$\blacktriangleright$ ~~▓▓~~ ~~▓▓▓~~ :

<div style="float:right; border:1px solid;">⌐SN⌐ t. coding of t. pure bern. seq. is unclear! what pc. is to be assigned to t. various probtys?</div>

.16   →Actually : .10 is probly <u>unnecy</u> : as we shall see. First, to code t. Bern seq. of $N$ symbol alphabet: ($N$ known a pri).

~~▓▓~~ Derb. t. CPM first in <u>prefix</u> code. This is simply a list of the probys of t. $N$ symbols, to adequate accuracy, (actually, ~~all~~ all accuracy levels should be used, and all conceivable sets of pc's for t. alphabet symbols ~~▓▓▓▓▓▓▓▓▓▓~~ — ~~▓▓▓▓~~ These give various CPM derns — ā t. final pc. of t. corpus is t. sum of ⟨ pc's of each CPM times pc of corpus w.r.t. that CPM⟩.

─── To Code t. <u>Bern Seq.</u> in which $N$, t. ~~▓▓▓~~ alphabet size, is known.

$N-1$ ~~▓▓▓~~ probty params, $P_i$, are needed. $(i = 0 \vert N-1)$ $\left(P_N = 1 - \sum_{i=1}^{N-1} P_i\right)$.

To express these, we use $R \equiv 1000$ bit accuracy ( T. result is indp of $R$ for $R >> 1$).

If t. sequence in $k$ bits long, t. accuracy needed in t. $P_i$'s is $\sim \frac{1}{\sqrt{k}}$.

As a result, of the $\not\leq 2^{1000}$ parallel codes for each poss. $P_i$ value.

.25   $\sim \frac{1}{\sqrt{k}}$ of them are about rite — so t. pc of each $P_i$ is $\sim \frac{1}{\sqrt{k}}$

The pc of t. $N-1$ params is $\left(\frac{1}{\sqrt{k}}\right)^{N-1}$, and the pc of t. seq. itself w.r.t. this RPM dern is $\prod_{i=1}^{N} \not\leq P_i^{k P_i} = \left(\prod_{i=1}^{N} P_i^{P_i}\right)^k$ .

.29   T. entire p.c. is $\approx \left(\frac{1}{\sqrt{k}}\right)^{N-1} \cdot \left(\prod_{i=1}^{N} P_i^{P_i}\right)^k$   ⟶ 3.30 for exact expressn.

.30   The $\overset{..}{\approx}$ is because t. p.c. of each $P_i$ gn. is probly not exactly $\frac{1}{\sqrt{k}}$ ) — but has to be figured out .... it is t. <u>acquist</u> pc., obtaind by using all ~~▓▓▓~~ $2^R$ codes for each $P_i$ : so $(2^R)^{N-1}$ difrnt codes.

The pc. is obtaind by some kind of integration, (rather than a summation) of codes (assuming uniform aprip density in $\vec{P}$ space ··· the "Laplace's Rule" assumption).

<span style="float:right">which is equivst. to</span>

.36   The way the coding is done! $N$ and $R$ are given apriori.

The first ~~RAWWW~~ $N-1$ sequences of $R$ bits each, give t. $N-1$ $P_i$ values.

The following seq. is t. code of t. corpus w.r.t. t. ~~▓~~ CPM just derbd.

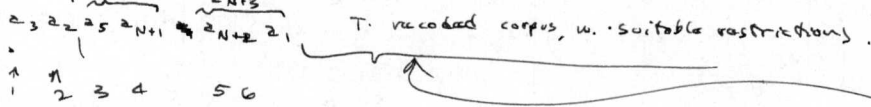We sum over all $2^{(N-1)R}$ CPMs, and we sum over all derns of t.

CODE

corpus w.r.t. each CPM ( There will be an ∞ of codes for t. corpus w·r·t each CPM,

.02     ordinarily .... unless the $p_i$'s are a terminating binary fractions $\Xi$. ( I think ! ) )

---

For t. coding using Binary definitions :.   If we are initially given

N , t. no. of Alphabet symbols $\overline{a_1 .... a_N}$ & D, t. no. of Binary defns used,

we code t. corpus by writing e.g.    ( D = 3 in this example )

.10   $\overbrace{a_3 a_2 a_5}^{a_{2N+1}} \overbrace{a_{N+1}}^{a_{2N+2}} \overbrace{a_{N+2} a_1}^{a_{2N+3}}$,   T. recoded corpus, w. suitable restrictions.
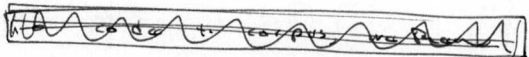
   ↑ ↑    ↑
    1 2 3 4   5 6

In positions 1 & 2, only $a_1$ thru $a_N$ are legal.

In positions 3, 4 , $a_1$ thru $a_{N+1}$ are legal

"   "    5, 6 ; $a_1$ " $a_{N+2}$ " " .

In Most positions after $\Xi$ , $a_1$ thru $a_{N+3}$ are legal.

— except for restrictions: after $a_{N+1}$ has been defined, $a_2$ can't follow $a_3$

"    $a_{N+2}$ " " " $a_{N+1}$ " $a_5$

"    $a_{N+3}$ " " " $a_1$ " " $a_{N+2}$ .

we then rewrite t. corpus using

t. defns. & write this prelimy coded corpus here

Then we code t. string at .10 as a Bern seq. w. N+D sized alphabet

with the restrictions mentioned .

To write t. rest of t. CPM, we first write N+D−1 strings of R bits

each, ( R >> 1 ) representing the R+D−1 probys involved .

The probty of t. corpus w·r·t each such assignment of probes in N+D−1 space

is computed as t. product of t. probys of t. symbols, taking t. restrictions

into account.     e.g. if, at a particular position, $a_3$ occurs,

& $a_7$ is   the only symbol illegal in that position, t. resultant probty used is

$P_{a_3} \dfrac{1}{(1 - P_{a_7})}$ : If $a_7$ & $a_4$ were illegal there, $P_{a_3} \dfrac{1}{(1 - P_{a_7} - P_{a_4})}$ would be used .

I don't know how to go about computing t. accuracies of each of t.

probys ( & t. pF of _ deriving each is ∝ its accuracy )

I'm not sure about how to get t. peak values

of t. probtys themselves.     $\not\equiv$

Hvr., t. discn. of pure Bern coding of 1.30 - 2.02 is to be applied

to obtain the "accuracy" & equivt & t. equivt peak values of the

$p_i$'s .    → see 3.30 - .40 & 4.01 ff. for further discn, hvr.

CODE

There remains t. p.c. of the ~~precise~~ coding parameter, D.
This has to be coded in a prefix code that is inserted before t. rest
of t. code I've derbd.

One interesting way to assign apripd to t. integers $z = 1, 2 \ldots \infty$. $(z = 1/\infty)$

Use $\cancel{\Sigma} P_i = \frac{1}{z}$ (!).  Consider the fact of non-convergence!

Say R is the upper limit of nos. we will consider;

Then t. normzn factor is $\approx \dfrac{1}{\ln R + \gamma}$

             $\underset{\text{Euler's const}}{}$

Since we are interested in ratios of apripd, this / normzn factor is irrelevant!
——— T. ratio of / apripd / D values $D_1$ & $D_2$ is simply $\dfrac{D_1}{D_2}$ (!)

               exact value of this

Actually, I can use any apripd I like for D — usually w. t.
restriction of convergence. But this $\frac{1}{D}$ thing really looks neat!

In general, it is not of much importance, hvr.  We are interested
because we want to know how much t. p.c. of t. corpus is ↑ by
a particular defn.  The main thing that tends to ↓ t. p.c. is
t. factor equal to $\dfrac{1}{\sqrt{R}}$ (1.25).  The factor $\dfrac{D}{D+1}$ is
so close to 1 (even if D is only 2), that it is usually of
~~no importance~~ negligible importance wrt. .

Hvr., if $D = 2$ $\dfrac{D}{D+1} = \dfrac{2}{3}$ : which does mean that one
method of coding will be gn. $1\frac{1}{2}$ times as much wt. as another ———
——— This could be of some import.  Also, & when $D = 1$ (an
impt. case) $\dfrac{D}{D+1} = \dfrac{1}{2}$ , which is a far sized factor !

     I ~~think~~ that much of t. foreg. stuff is relevant to
t. problem of coding w. linear (& non-linear) regressn;
determining how many (& which) regn. coifs. to use.

.30: 1.29 :  **Instead of** $\left(\dfrac{1}{\sqrt{R}}\right)^{N-1}$ , all we want is t. total (wt. in $N-1$ dim space of
t. $P_i$ coifs:

            This is exactly

Here $m_i$ is t. no. of times $\cancel{=}$
     as occurd; $P_N \equiv 1 - \sum\limits_{i=1}^{N-1} P_i$

$$\underset{N\text{-fold}}{\int \int \int \cdots \int} \prod_{i=1}^{N} \left( P_i^{\; m_i} \, dP_i \right)$$

I'm not sure this ← is an N·fold $\int$ — it may be only $N-1$ fold.

~~anyway~~ The integral ~~shown ...~~ involve is well known and rather simple in structure.
The complicated part involve a $\cancel{\&}$ ~~Gamma function~~ — which may be where the
approxn $\frac{1}{\sqrt{R}}$ comes in.

                       Maybe not so easy ··· see
                               4.17 !

Code

For t. corresponding factor in "coding w. binary datas",

we have to integrate $\underline{\underline{not}}$, e.g. $\rightarrow p_1^{17} p_2^5 \cdots p_9^{21} (1-p_1-p_2\cdots -p_9)^4$

but $p_1^{17} p_2^5 \cdots p_9^{21} (1-p_1-p_2\cdots -p_9)^4 \cdot \boxed{\left[ (1-p_2)^4 (1-p_5)^3 (1-p_5-p_3)^5 \right]^{-1}}$

The new factor occurs because of the restrictions in coding at certain positions

$\underline{\underline{Woops!}}$ There $\underline{is}$ a difty! This factor causes divergence! $\int_0^1 \cdots \frac{1}{(1-p_2)} dp_2$ diverges.

It $\underline{may}$ be cancel'd out by the $\approx (1-p_1-p_2\cdots -p_9)^4$ term

Also, since we do have t. constraint $\sum_{i=1}^{N} p_i = 1$, if $p_2$ here got close to 1,

then all of the other $p_i$'s would have to be close to zero!

So: This integral ~~w~~ is $\underline{not}$ so simple to do:

.17 Note also that this integral is more dift than I thot, since

$\sum_{i=1}^{N} p_i = 1$ must hold over t. region of integration! So

we can't just integrate each variable indiply. — This modifys t.

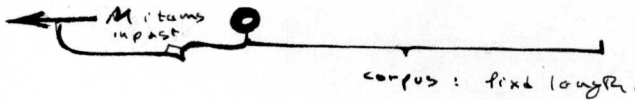ergt. of 3.30 — .40 as well!

Code:

Linear (& Non-linear) Regn:

One ~~[xxxxx]~~ part that was rather messy, was coding t. first # M ~~new~~ items of t. corpus ( M is t. no. of coils to be usd ).

.03     I had previously considered coding each of them optimumly in terms of to previous/ items —— "known" being "coded thus far".    ("known")

A much easier way, would be to consider t. first M items as given free : then code t. rest of t. seq. w.r.t. that info. We use something like

complexity $(x/y)$ ≡ complexity of $x$ given $y$.

~~Then~~ For larger M, we have to consider items further into t. past — if we keep t. corpus itself of constant length.



corpus : fixd length.

M items in past

As we ↑ M, we keep corpus constant, but M extends further into t. past.

---

8583 ! See more recent Notes on this ; e.g. Notes for Harvard talk

→ in 3/83 ! Trouble is each of these m. ~~data pts~~ represents an enormous amt. of "info". Giving M+1 rather than m coils is prostration of price of I would ! which ~~xxxx~~ very much.    ("info")

Actually .03 isn't bad. We start w. 1 ( or 2 or 3 ) ~~xx~~ data pts known. whether it's 1, 2 or 3 ~~xxx~~ will only affect whether we decide to terminate t. search w. 1, 2 or 3 coils (I think) — so it probly does not affect t. final answer — but check on this

→ Keep to ... is !

PSG-discy    stochastic                                              H.C.
# On discovery of ~~stat~~ Models Using / & Entropy estimates ~~is~~ as hill ht.                    1
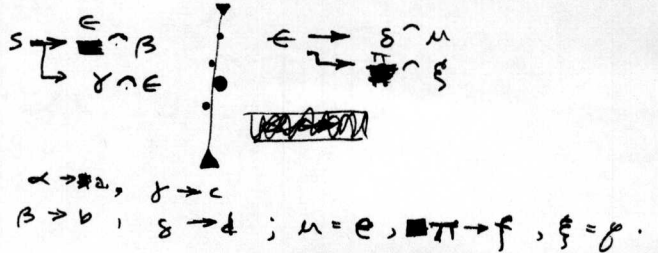
— Mainly on how large a SSZ is needed to "identify" PSG's

— or what ~~xxxxxxxx~~ predn. accuracy is obtainable w. what SSZ.

    Also on t. efficacy of using simpler grammar models for ssz,

(or even SSZ), & then modifying t. simpler grammars into PSG's

(ie. introducing loops or rules giving potentially infinite lays.

    I'm thinking of 3 levels of Grammars: <u>ＧＧ</u>.

    1) Grammar ~~~~ & <u>finite</u> lays:

ie. <u>no</u> loops , but otherwise a psg: <u>stochastic choices</u>.

    2) <u>F</u>inite State lays., w. stochastic choices . . . . . . .

    3) stochastic PSG w. loops being legal.

•···· Big Bang

$$S \to \in \char`\^ \beta ,\ \gamma \cap \in$$
$$\in \to \delta \char`\^ \mu ,\ \pi \char`\^ \xi$$

$\alpha \to a,\ \gamma \to c$
$\beta \to b,\ \delta \to d ;\ \mu = e,\ \pi \to f,\ \xi = g$.

---

In ① if one had a large enof ssz, one could get t. prby of each s in t. lays.

& then try remove Grammars to reduce t. no. of /params. derbing t. lays.    (continuous)

    Hvr. perhaps t. usual case is that t. lays. is ~~to~~ has too many members for this to be practical, so the loopless PSG ~~~~ gives a very <u>large</u> # in t. amt. of info in t. grammar.

---

    Actually, I usually think of / finding a type 1) model ~~easily than~~ going to a type 3) model. — ~~~~ I haven't really considered that a type 2) model ... being able to model an infinite lay., mito be closer to 3) than t) is! (<u>E</u>xcept that the idea of t. "a new method .... 1960" paper was that PSG's are very ~ to PSG's. & this fact mito be used to implement PSG discy ).

    Any type 3) Grammar can be written as a (potentially infinite)

.37    type 1) Grammar!

    type 3)

$\alpha \beta$.

$$S \to \blacksquare \alpha$$
$$\alpha \to \alpha a$$
$$\ \ \ \to b$$

so all as's are : $b a^n$.

PSG-discy

An infinite type 1 [crossed out] PSG for this:

$$S \rightarrow \alpha \quad \Big| \quad \alpha_1 \rightarrow \alpha_2 c \quad \Big| \quad \alpha_2 \rightarrow \alpha_3 c$$
$$\alpha_0 \rightarrow \alpha_1 c \quad \quad \rightarrow b \quad \quad \rightarrow b \quad \cdots \cdots$$
$$\rightarrow b$$

.03

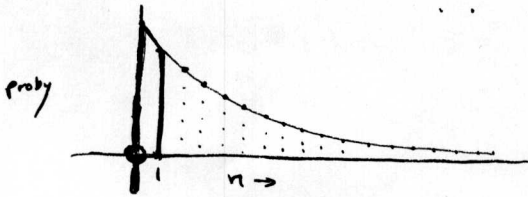Its not clear whether this grammar would have less info than the original corpus. — for any finite corpus, only a finite no. of grammar rules would be needed. Also, if we allow repetition of αs's in the given corpus, the corpus would ordinarily have much more info than the grammar.

e.g. say the proby of $\alpha_i \rightarrow \alpha_{i+1} c$ was $p$; $\alpha_i \rightarrow b$ was $1-p \equiv q$.

So, if $p$ is close to 1, we will have a mean string length of

$\sim \frac{1}{q}$ ⎰ a string of length 1 has proby $q \cdot p^0$
⎱ .... 2 " " $q p^1 = q \cdot p^{2-1}$
.... $n$ " " $q p^{n-1} = \frac{q}{p} \cdot p^n$


proby

If the SSZ is such that we have 1 string of length $n=n_0$, then we have $\sim$ ( [scribbled] )

$\frac{1}{p^1}$ strings of length $n-1$
$\frac{1}{p^2}$ " " " $n-2$ , etc.

so, total string length $= \frac{n}{p^0} + \frac{n-1}{p^1} + \frac{n-2}{p^2} \cdots \cdots \frac{1}{p^{n-1}} + \frac{n-n}{p^n}$.

$= p^{-n}( p + 2p^2 + 3p^3 \cdots \cdots np^{n} )$

$\Sigma = p^{-n} \sum_{i=0}^{n} i p^i = p^{1-n} ( \sum_{i=1}^{n} i p^{i-1} ) = \frac{d}{dp} ( \sum_{i=1}^{n} p^i )$

$\frac{d \frac{\alpha}{\beta}} = \frac{\beta d\alpha - \alpha d\beta}{\beta^2} = d( \alpha \cdot \frac{1}{\beta} ) = \frac{1}{\beta} d\alpha + \alpha d(\frac{1}{\beta})$ ⎰ $\frac{d}{dp} (\frac{1-p^{n+1}}{1-p})$ ⎱

$d \frac{1-p^{i+1}}{1-p} = \frac{(1-p) \times (-(i+1)) p^i - (1-p^{i+1}) \cdot (-1)}{(1-p)^2}$ ⎰ $\frac{\alpha - \alpha d\beta}{\beta^2}$ = $\boxed{\frac{1 - p^n(1+(1-p)n)}{(1-p)^2}}$

[scribbled-out lines]

proby of string of length $n$ is $q p^{n-1}$

so $\Sigma = \frac{1 - p^n(1+(1-p)n)}{p^{n-1}(1-p)^2}$

.40

The idea here, is that    if the longest example string is of length $n$, ($n \gg 1$), then the no. of ~~this~~ choices in t. raw corpus would be very large:    for $p = .9$ & $n = 20$, t. no. of choices made is $\sim 470$ :   (from 2.40).

On t. other hand, t. no. of choices ~~in~~ to create a non-loop PSG would be $\approx 3n$ (see 2.03) which would be $\ll$ t. no. of choices in int. raw corpus.

Hvr. this is __not__ t. point!   Even using a good grammar, t. no. of choices needed to create t. corpus, would be about t. same as t. no. of symbols in it : i.e. (2.40). So the avgt. __would__ have to be very detaild to tell if a loopless PSG ~~did~~ did, indeed, achieve a compression $>$ say a simple __Bernoulli code__!

Actually for t. $g$ lng of 1.37   w.    $s \rightarrow \alpha$  ;  $\alpha \overset{.9}{\rightarrow} \alpha_2$, $\overset{.1}{\rightarrow} b$

t. Bern seq. is not so bad: for z large corpus,     ,,,,,,,,,

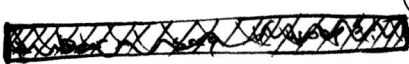$p_\alpha \sim .9$, $p_b \sim .05$, $p_{end \, symbol} \sim .05$

Grammar-Grammar

Hvr., t. __main idea of t. present discu. is:__ If one is gu. z loopd PSG corpus, is it useful to first find a good (loopless PSG) for that (= type 1 gramm) corpus .... then try closing t. loop(s) .... using H.C. z apparent Entropy ($\equiv$ Pcost) ~~as a~~ Gore.

Whether z Bern seq. would give a better Gore is, hvr. (Not relevant) ... because → A Bern seq. is __not__ a type 1 Gramm, since it is ~~for~~ an infl. lang. ( Here we use z stop symbol as one of t. Bern ~~lang~~ alphabet, so we get $>1$ sentence in t. corpus ).

(SN) From a large loopless Gramm of a loopd PSG corpus, one __can__ devise a grammar for this _____ ( potentially infinite) Grammar, & perhaps from this "Grammar Grammar" find an easy way to close t. loop(s).

Hvr., what has always seemd like a u.g. idea, is to find out how to find t. Gramm of a ~~(stochastic)~~ finite State lang. From this use t. correspondence betw. FSG & PSG to devise a method for PSG discy.

There __may__ be known ways to do this for FSG : see Taylor Booth: papers, reviews, book(s).

TM: PSGb3cy:

First make as good as possl. TREE grammar. This is a loop free CF Grammar. For a given corpus size, there will be a bunch of ± "Best" Grammars. From one (or more) of these tree grammars, one tries to xfm them into **looped** grammars.    Here, each Tree grammar is regarded as an object ~~operation that~~ that one wants to xfm into a more "compact" ( more pc ) object — ie. t. loopd grammar.   Once we have t. tree grammar, we don't go back to t. original corpus.

T. only (sort of ) back tracking method usd, is to retain several Tree grammars & try to ~~loop~~ "loopify" each of them, to see which results in t. ~~the~~ best "pc" object.

T. reason I think this is a good approach: "Working backwards" Say one had ~~t. correct loop~~ a good loopd grammar for t. corpus! Then one could break t. loops ~~~~ in various ways & get Tree grammars of/ probably lower pc.   Looking at it this way makes it reasonable that .01 ∯ **WOULD** work.

How many diferent ways can on convert a loopd grammar into a tree Grammar?   Does one always ↓ pc by doing so? Do t. various alternative Tree grammars all have about t. same pcs?

---

consider t. loop grammar language (a∅)ⁿ:

Grammar!
$$S \to a\emptyset$$
~~$S \to a\emptyset$~~  )←

Some Tree grammars
$$S \to \blacksquare S_1 \qquad \to S_1 \, S_2$$
$$S_1 \to a \qquad \to S_1$$
$$S_2 \to S_1 S_1$$
$$S_3 \to S_2 S_2$$
$$S_4 \to S_3 S_3$$

●

or
better
Tree grammar
$$S \to \blacksquare S_3$$
$$S_1 \to a$$
$$\to \Lambda$$
$$S_2 \to S_1 S_1$$
$$S_3 = S_2 S_2$$

$$S_1 = a^{0,1}$$
$$S_2 = a^{0,1,2,3}$$
$$S_3 = a^{0,1,2,3,4,5,6,7}$$

or
$$S \to S_3 S_3$$
$$S_3 \to S_2 S_2$$
$$S_2 \to S_1 S_1$$

$$S_1 \to a$$
$$\to \Lambda$$

How to go from! ← this to this → ?

(Actually, t. Grammars of .35 R(ie .35 m — .90 in general) are not exactly what one gets from $\sum_{TB} 141$ (coding u. deduns.). $Z141$ gives a lang. That's a Bernoulli seq.

Hvr., t methods of $Z14c$ cen be applied to ordinary langs à t. devising of free Gramars for Regx.