

語彙量予測法について

松岡 潤, 絹川博之, 筒井健嗣, 尾本 健
(日立製作所 システム開発研究所)

1. はじめに

電子計算機によって、自然語を処理する場合に、大抵の場合において、標準となる言葉や文字列の集合を、予め計算機システムの内部に蓄えておく(これを機械辞書と呼ぶ)、入力される処理対象自然語と、この機械辞書のデータとの比較を行ない、その結果に基づいて処理が行なわれる。このようなシステムの性能を評価するには、機械辞書に含めているデータが、一般の自然語全体に対してどれだけの部分を覆っているか、またデータを新たに追加したとき、覆率がどれだけ増加するか、が推定できることが必要となる。

自然語の中から、機械辞書に収容する予定の語を抽出・選択するには、いわゆる語彙調査が必要である。語類(言葉や漢字など)の出現率に関する調査は、国立国語研究所をはじめ、多くの機関や、人々によって今までに行なわれて来ているが、^{(1)~(4)}数万~数千万の言葉や文字を扱うため、膨大な作業量を必要とする。それ故、個々の調査対象語類に対し、語彙調査に先立って、それに必要となる作業量を推定できることが望まれる。

語彙調査結果に対する要求覆率(自然語の中から任意に1語を採って来たとき、その語が調査結果内にすでに入っている確率を「覆率」と呼ぶ)を与えられたとき、それを満足するための語彙の必要調査量を推定する方法についての検討を、本稿は述べるものである。

2. 語類の出現率の近似式

2.1 用語の説明

以下に、本稿で用いる用語について、予め説明する。

語類: 自然語の中で、着目している単位語の類別を表わす総称である。たとえば、日本人の姓に着目して考えるとき、対象語類は日本人の姓であるといい、ある文献中の漢字に着目するとき、対象語類は漢字であるという。

語: 1つの語類を構成する単位要素をいう。例えば、日本人の姓を対象語類とすると、「鈴木」は1つの語であり、ある文献中の漢字を対象語類とすると、文字「木」は1つの語である。

覆率: 1つの語類に属する語の集団 α を考え、 α の覆率とは、対象語類に属する任意の1つの語を採ったとき、その語が α に入っている確率をいう。

非覆率: 語の集団 α を考え、 $1 - (\text{覆率})$ を α の非覆率という。

2.2 出現率の近似⁽⁶⁾

1つの語類において、個々の語による出現率の変化の様子の例を、図1.に示す。同図は、文献(4)による雑誌の語彙についてのもので、同文献の表1をグラフにプロットしたものである。横軸は、出現率の高い語から順に1, 2, ... とつけた語の順位番号 n を表わし、縦軸は、順位1から順位 n までの語の出現率の和(すなわち、累積出現率)を表わしている。

この曲線は

$$y = 1 - e^{-ax}$$

(1)

と似た形をしているので、

$$Y = 1 - y \quad (2)$$

とおき、 (x, Y) を片対数グラフにプロットしてみたところ、図2. のようで、直線とはならない。しかし、図2. は横倒しの拋物線に近い形をしているから、直線化するために、さらに両軸の対数を取りプロットする。すなわち、横軸には $\log n$ をとり、たて軸には、 $\log\{1 - \log(1 - y)\}$ をとる。これを行なったものが、図3. である。同図においては上記のことを能率的に行なうため、横軸は対数目盛であり、たて軸は、原点から $\log\{1 - \log(1 - y)\}$

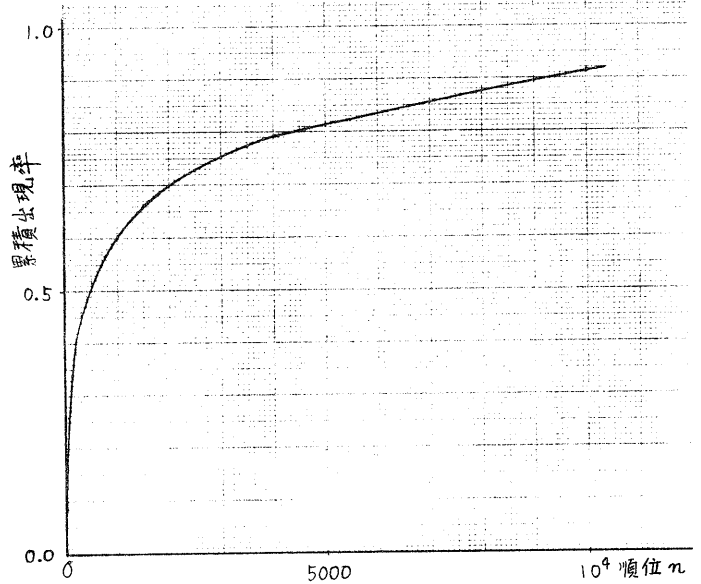


図1. 雑誌の語彙の累積出現率

の折に $y \times 100 (\%)$ と目盛ったものを用いており、一種のワイブル確率紙である。これをCFグラフと呼ぶことにする。

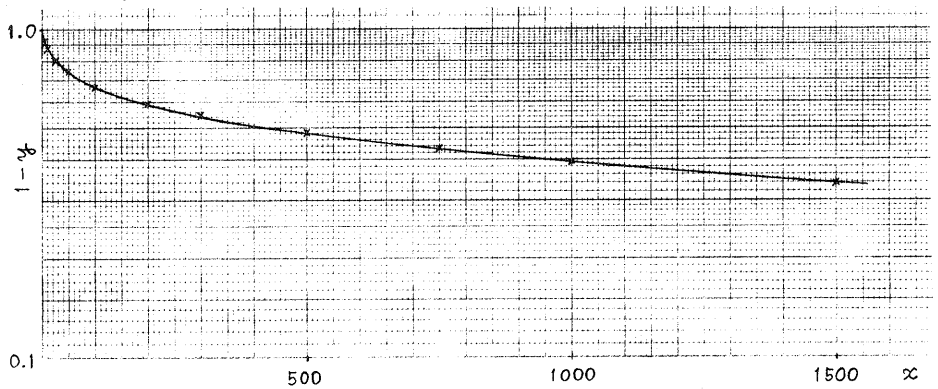


図2. 雑誌の語彙における累積出現率の補数の変化

この図で

は、曲線はほぼ直線に近い2つの部分と、それを滑らかに結ぶ曲線部分から成っている。このことから、この曲線は次の式で近似できることがわかる。

$$\begin{aligned} f(n) &= \phi_1(n) \cdot \phi_2(n) & (a) \\ \phi_i(n) &= 1 - \exp(-p_i n^{q_i}) & (i=1, 2) & (b) \end{aligned} \quad (3)$$

ここに、 n : 語の順位番号

$f(n)$: 順位1から n までの語の出現率の和

p_i, q_i ($i=1, 2$) : 定数

式(3)では、 $f(n)$ は $\phi_1(n)$ と $\phi_2(n)$ とに別して対称である。そこで次のように ϕ_1 と ϕ_2 とを区別することとする。CFグラフは、 n の値が約50以上の部分では、ほぼ直線に近い。したがって $\phi_1(n)$ または $\phi_2(n)$ のどちらか一方のもの(仮にこれをAとする)は0.99以上となり、他方のもの(仮にBとする)が $f(n)$ の値を殆んど決める状態となっている。このBを $\phi_1(n)$ と命名し、Aの方を $\phi_2(n)$ と命名することとする。

筆者らの調査した会話の語彙、および他の発表されている語彙についても、同様にCFグラフを示すと、図4.、図5. のようになる。これらの図のデータの出典

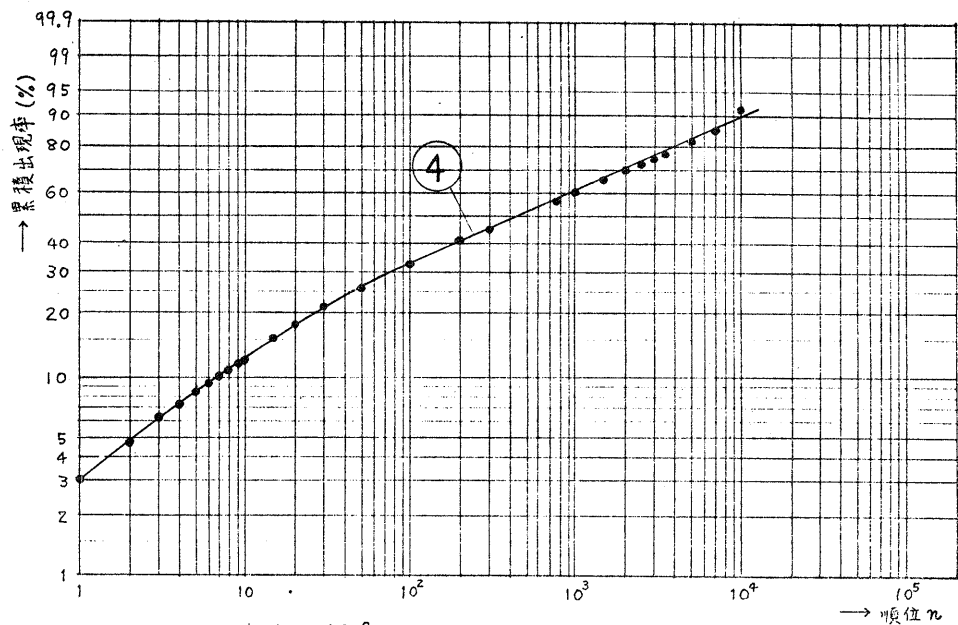


図3. 雑誌の語彙のCFグラフ

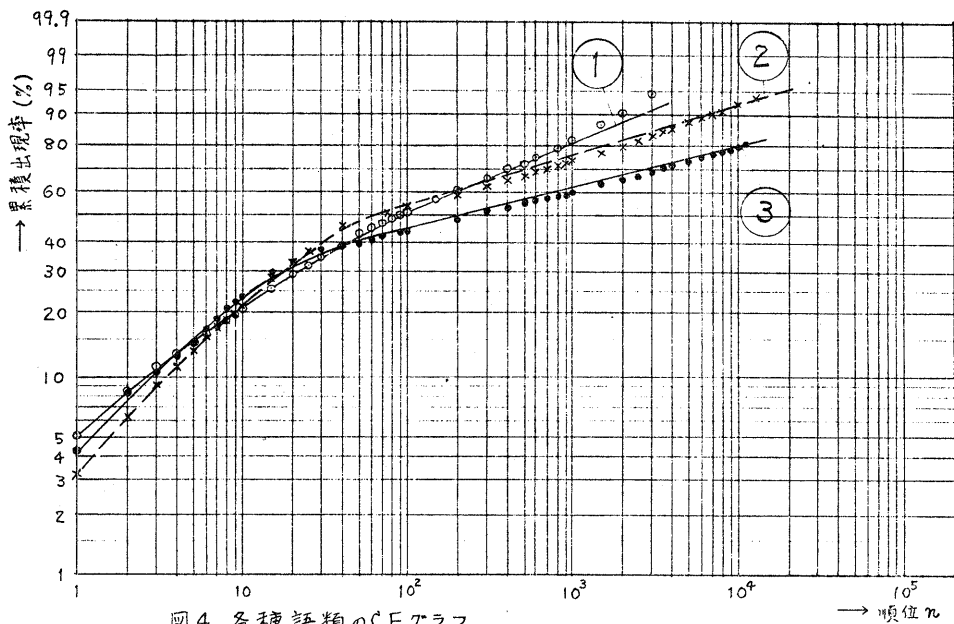


図4. 各種語類のCFグラフ

等は表1.に示すとおりである。図で、線によつて示してあるものは、各語類について、定数 p_i, q_i ($i = 1, 2$) を表1.のよつた定め、式(3)によつて計算された値である。これらの図から、 $f(n)$ が1に近い部分を除き、式(3)で近似できることがわかる。相対誤差は5%以下である。

定数を定める方法は、次のようにする。まず、実測値グラフの中央部分の点が直線に近い部分に着目し、これを近似する直線を求める。これにより $\phi_i(n)$ が定まる。次に n の小さい部分の適当な m 個の実測値から

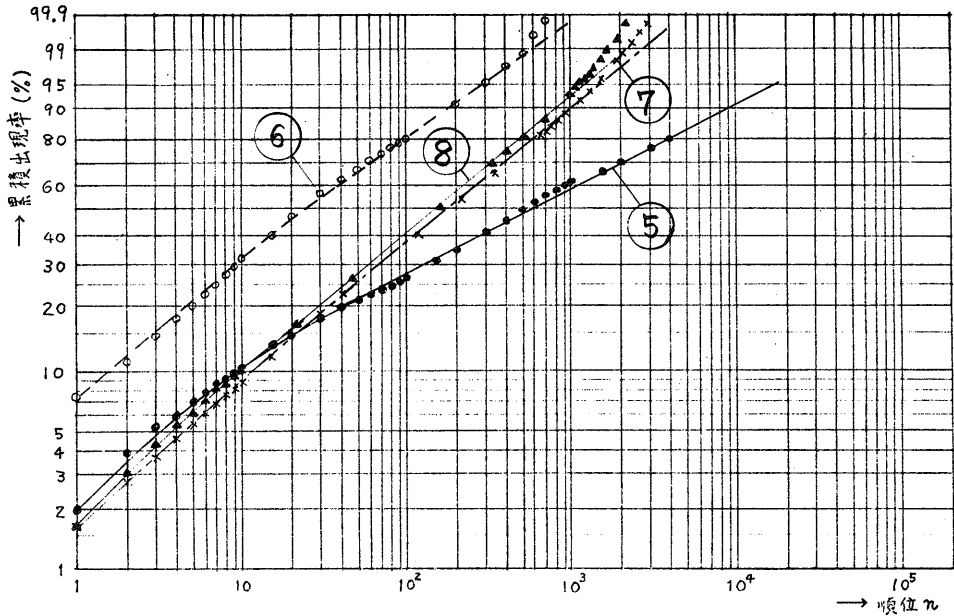


図5. 各種語類のCFグラフ

表1. 各種語類の出現率調査データと近似式の定数

項目		語類名	(1) 会話の語彙	(2) 新聞の語彙 (短単位)	(3) 新聞の語彙 (長単位)	(4) 雑誌の語彙	(5) 日本人の姓	(6) 姓における漢字	(7) 雑誌の漢字	(8) 新聞の漢字
実測データ $f_x(n)$	(a) 調査対象	お運記録	3大新聞 (S41年分)	同左	雑誌90種 (S31年分)	都道府県教職員名簿	文部9)	雑誌90種 (S31年分)	3大新聞 (S41年分)	
	(b) 報告文献 (参考文献番号)		7)	同左	4)	9)	10)	5)	8)	
	(c) 調査機関	日立研究所 シテム研	国立 用語研	同左	同左	佐久間英 田中原仁	国立 国語研	同左		
	(d) 標本語数	39588	940533	679342	532795	約35万	約1.2万**	280094	630313	
	(e) 異なり語数	5118	47805	101081	約4万	約10万*	911	3328	2879	
	(f) グラフ番号	①	②	③	④	⑤	⑥	⑦	⑧	
近似式 $f_x(n)$	(g)	p_1	0.114	0.217	0.221	0.069	0.040	0.096	0.018	0.021
	(h)	q_1	0.397	0.266	0.210	0.378	0.452	0.625	0.706	0.699
	(i)	p_2	0.647	0.179	0.241	0.593	0.715	1.658	2.229	1.574
	(j)	q_2	0.469	0.801	0.798	0.445	0.706	0.194	0.215	0.658
備考			* 約35万の中での異なり姓ではなく、日本全国でのもの ** 原文とないで筆名が推定							

$$h_j = \frac{f_x(n_j)}{\phi_1(n_j)} \quad (j = 1, 2, \dots, m) \quad (4)$$

を求め、この m 個の点 (n_j, h_j) について、これらを近似する直線 $\phi_2(n)$ を求める。ただし、 $f_x(n)$ は実測の累積出現率である。この方法は最良近似ではないが、簡単であり、実用上の精度が損なわれない。表1. に示した各定数は、この方法で、相対誤差を最小とするような最小二乗法で求めたものである。

$f(n)$ が1に近い部分では、実測値グラフは n の増大とともに急激に上昇し、式(3)からはすれてゆく。このことを仮に「グラフの急上昇」と呼ぶこととし、次節で考察する。

2.3 グラフの急上昇の理由¹⁷⁾ および異なり語数

有限の語の集合である、1つの標本をとれば、その中の異なり語の数も有限である。異なり語の数をいま k とすれば、標本の世界では、当然 $f(k) = 1$ であり、CF グラフ用紙¹⁷⁾では、その点は上方無限遠に位置するから、グラフの急上昇は当然ということができる。しかし、その近辺におけるグラフの状況を、次のように考察することができる。

先づ、語類の母集団を、無限要素をもつものと考え、そこにおいて式(3)が成り立っているものと仮定する。語類の出現確率密度関数 $g(n)$ は

$$g(n) = f(n) - f(n-1) \quad (5)$$

と定義される。

いま、この母集団からランダムに、大きさ N の標本 S を抽出したとする。この標本の中に、順位 n の語が ν 個含まれている確率を $P(N, n, \nu)$ と書くことにすれば

$$P(N, n, \nu) = \binom{N}{\nu} \{g(n)\}^{\nu} \{1 - g(n)\}^{N-\nu} \quad (6)$$

のような2項分布に従うものと考えられる。この分布に従うバラツキによって、母集団で順位 n であつた語が、標本 S では順位が n となるとは限らない。 S での順位番号の期待値を $n^*(N, n)$ と書くことにすれば

$$n^*(N, n) = n - \sum_{n'=1}^{n-1} P(N, n', 0) \quad (7)$$

と書ける。これは、語 n が標本に現われたとすれば、母集団において順位 1 から $n-1$ までの語で標本 S に現われなかったものの分だけ、番号がくり上るからである。母集団での n の近辺の順位の語が、バラツキによって相互に順位が入り乱れる現象も起るが、それらは多数回の試行ではバランスし、期待値としては式(7)のようになることがわかる。

式(6)を用いると、式(7)は

$$n^*(N, n) = 1 + \sum_{n'=1}^{n-1} [1 - \{1 - g(n')\}^N] \quad (n > 2) \quad (7)$$

と書ける。

さらに、以上のことから、異なり語数の期待値長 $k(N)$ は

$$k(N) = \lim_{n \rightarrow \infty} n^*(N, n) \quad (8)$$

によって表わされることがわかる。

$n^*(N, n)$ 、および $k(N)$ の数値計算に当っては次のようにした。

(a) 出現率関数 $g(n)$ の計算に当っては、 n の小さい所では定義式(5)によって計算し、 n の大きな所では、次式によつた。

$$\left. \begin{aligned} g(n) &= \frac{df(n)}{dn} = (A_1 + A_2) \frac{f(n)}{n} & (a) \\ A_i &= \frac{g_i t_i}{\exp(t_i) \sinh(t_i)} & (i=1, 2) & (b) \\ t_i &= \frac{p_i n^{q_i}}{2} & (i=1, 2) & (c) \end{aligned} \right\} (9)$$

(b) 式(8)の極限を計算機で求めるに当っては、次のようにした。

$$k(N) \doteq 1 + \sum_{n'=1}^{n_r} [1 - \{1 - g(n')\}^N] + N \cdot \{1 - f(n_r)\} \quad (10)$$

n_r とは $Ng(n_r) \leq 0.13$ となるようにとることにより、相対誤差を1%以下とした。

数値計算の結果の例を図6.に示す。対象語類は前出の会話の語彙である。図で、筆者らの調査した2つの標本サイズの場合についての実測結果と、上記の方法に

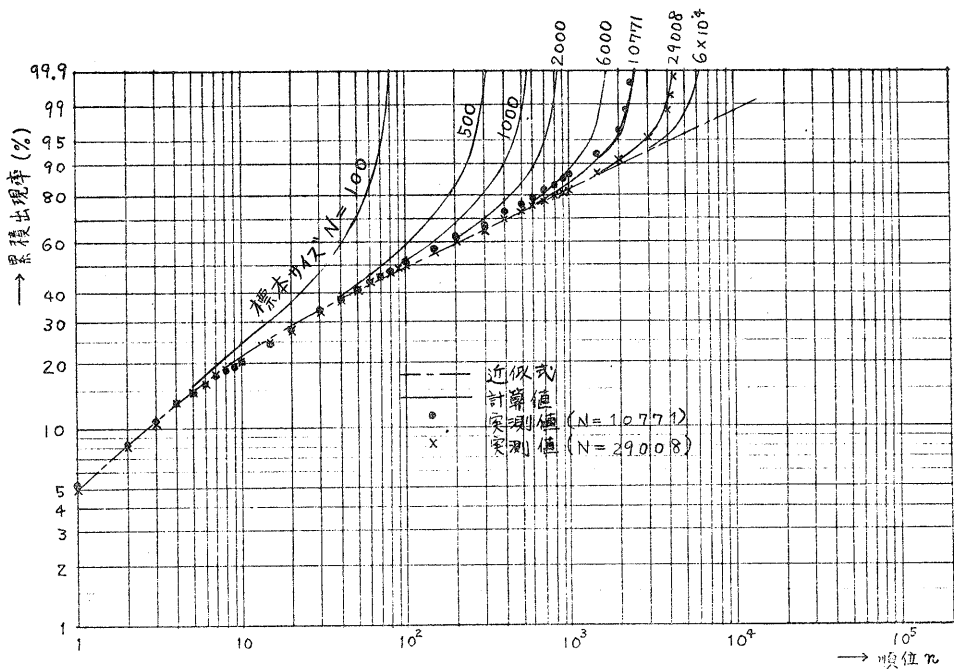


図6. 会話の語彙の標本サイズとCFグラフ

よる計算結果の対照をも示した。また、他の例について、報告されている標本サイズ N から異なる語数 $k(N)$ を、式(10)によって計算したものと、報告されている $k(N)$

表2. 異なる語数 $k(N)$ の実測値と計算値との比較

項目	会話の語彙			新聞の語彙 (長単位)	雑誌の語彙	
	標本の大きさ N	10771	29008	39588	679342	532795
$k(N)$	実測値 k_1	2342	4391	5112	101081	約40000
	計算値 k_2	2460	4232	4951	104860	42811
相対誤差 (%) $(k_2 - k_1) / k_1$		5.0	-3.6	-3.3	3.7	7.0

の予測値との比較を、表2.に示した。

以上のことから、「グラフの急上昇」は母集団から有限の標本を抽出したことによって起るものであり、無限要素をもつ母集団では n の大きい時まで式(3)によって表わされるものと考えられる。

3. 語の大量集団の覆率の推定

いま、調査対象母集団から、ランダムに大きさ N の標本を抽出したとし、この標本の覆率、期待値を $S^*(N)$ と書くこととする。母集団における順位 n の語が標本に出現しない確率は、式(6)による $P(N, n, 0)$ であるから

$$S^*(N) = 1 - \sum_{n=1}^{\infty} P(N, n, 0) \cdot g(n) \quad (10)$$

と書ける。再び式(6)と

$$1 = \sum_{n=1}^{\infty} g(n)$$

とを用いると

$$S^*(N) = \sum_{n=1}^{\infty} g(n) [1 - \{1 - g(n)\}^N] \quad (11)$$

とも書ける。この和は次のようにして近似的に求まることかである。

$$S^*(N) = \sum_{n=1}^{m_1} g(n) [1 - \{1 - g(n)\}^N] + \sum_{n=m_1+1}^{\infty} g(n) [1 - \{1 - g(n)\}^N] \quad (11')$$

であるが、 m_1 が充分大きければ、

$$0 < g(n) \ll 1 \quad (12)$$

であるので、式(11)右辺第2項は
第2項 $\equiv N \sum_{n=n_t+1}^{\infty} \{g(n)\}^2$ (13)

である。いま

$$I \equiv \sum_{n=n_t+1}^{\infty} \{g(n)\}^2 \quad (14)$$

$$\eta(n) \equiv -\frac{p \cdot q_1}{2 n^{1-2p}} \exp(-2 p_1 n^{2p}) \quad (15)$$

とおくと

$$0 < I < -\eta(n_t) \quad (16)$$

であることが証明できる。このことにより

$$S^*(N) \equiv \sum_{n=1}^{n_t} g(n) [1 - \{1 - g(n)\}^N] - N \cdot \eta(n_t) \quad (17)$$

によつて $S^*(N)$ は計算でき、その誤差 ε は

$$|\varepsilon| < -N \cdot \eta(n_t) \quad (18)$$

である。

この式によつて計算した例を
図7. に示す。ただし、この図に
おいては、覆率の期待値 $S^*(N)$ の
代りに、非覆率の期待値 $S(N)$ を
左2軸にとつて示してある。ま
たお曲線上には標本サイズ N に対
する異なる語彙の期待値 $k(N)$ を
目盛った。この形式のグラフを
NKS グラフと呼ぶことにする。
この図の利用法を次の例に示す。

(例) 宛名のカナ漢字変換シ
ステムのような、姓を扱う情報
システムで、処理対象データ
の99%以上を処理可能とした

という場合。図において $S(N) = 1 - 0.99 = 0.01$ の所を見ると、

$$N \approx 1.7 \times 10^6$$

$$k(N) \approx 4.5 \times 10^4$$

が得られる。それ故、170万個の姓を集めて来て調査し、そこに現われるはず
の約4万5千個の異なる姓全部を機械辞書に入れる必要がある。(例終り)
以上のことにより、語彙の大量調査により得られる語集団の覆率の推定を行なう
には、次のようにすればよいと言える。

(a) 1万~2万程度の標本を対象とする予備調査を行なう。

(b) 調査結果をCFグラフ化する。

(c) 近似式(3)の係数を求める。

(d) 定まった近似式により、適当にとつた幾つかの N に対し、異なる語彙 $k(N)$ を
計算する。(式(8))

(e) 適当にとつた幾つかの N に対し、覆率 $S^*(N)$ を計算する。(式(11))

(f) 上記(d), (e)の結果をNKSグラフに描く。

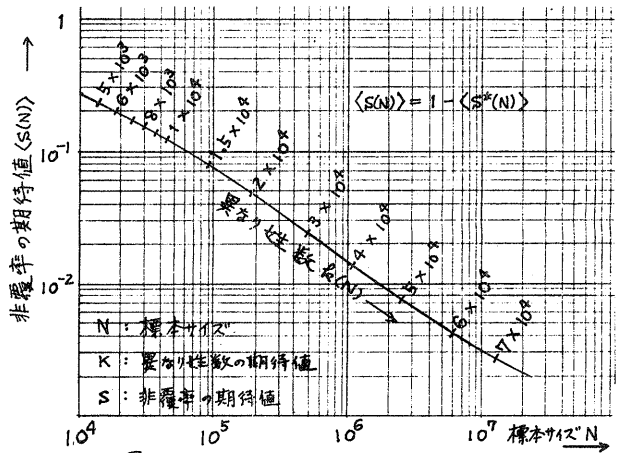


図7. 日本人の姓のNKSグラフ

なお、予備調査での標本の大きさを、あまり小さくすると、係数 p_1, q_1 として読んだ値を得ることになるので注意を要する。1万～2万は経験的な値である。

4. 最適な語集合とランダムな語集合

母集団における各語の出現率順位番号 n がわかっていると、各語を $\alpha(n)$ と表わすことにする。任意の正の整数 m を与えられたとき、語の集合 $\{\alpha(n) \mid n=1, 2, \dots, m\}$ を大きさ m の最適な語集合と呼ぶことにする。自然語処理システムの機械辞書を作るには、特に他の条件がない限り、最適な語集合を辞書に格納することが、能率のよいシステムのための必要条件となる。それは辞書用の記憶場所の大きさを一定とすると、最適な語集合が覆率の期待値を最大とするからである。

図8.に性のランダムな語集合と、最適な語集合とでの覆率の比較を示す。破線で示した補助線は、ランダムな語集合の曲線上の各点(例えば $N = N_0$ の点)を、水平に $N = k(N_0)$ の位置まで移動したものである。例

えば、ランダムに6500個の語集合をとると、それは点Aで示された覆率をもつ。この語集合に含まれる異なり語は点Bの横座標で示される数(すなわち3000個)だけある。しかし、この異なり語の集合も最適な語集合ではないために、点Cより約10%低い覆率しかもっていない。

図8.から次のことがわかる。ランダムな語集合は、最適な語集合にくらべて、覆率は小さいが、その差において、

- (1) 語集合の大きさ N が大きくなると、 N と異なり語数 $k(N)$ とのちがいに起因する差が大部分を占める。
- (2) N が小さい所では、逆に異なり語の組合せがランダムであることによる差が大部分を占める。

5. おわりに

本稿の提案する推定法は、予備調査として1万～2万程度の語彙調査を必要とする。この予備調査自身、そう簡単ではない。しかし、数10万～数100万の語彙調査を行なうことは、10数人が数年を要する作業となり、かつ調査の過程で誤りの混入を防ぐことも至難の業となる。したがって、調査結果のもつ価値と、所要コストとの明確な見通しがない限り、実施が困難である。本方法は、①対象語類の性格のちがいが、どのように作業量に影響してくるか、②得られた結果が、母集団に対してもつ覆率はどうか、について回答を与える試みである。

今後、さらに次の検討を行なってゆきたいと考えている。

- (1) 本方法は、予備調査に基づく累積出現率を外挿する方法である。その妥当性

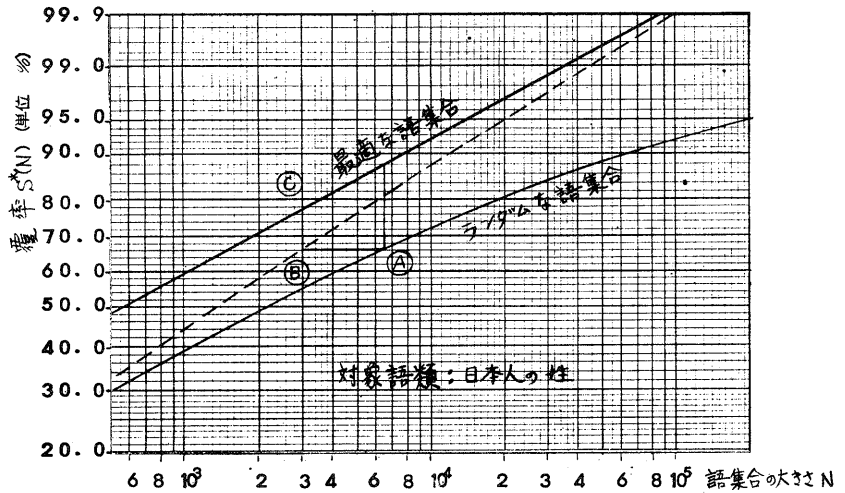


図8. 語集合の覆率の比較

日立製作所システム開発研究所

は、間接的には確認されているが、機会をとらえ、直接的確認も行ないたい。
 (2)本推定の各段階、すなわち、予備調査のサンプル抽出、統計過程、モデル式の係数決定、などの段階でもち込まれる誤差が、推定結果ほどの程度の誤差を与えるかについての吟味。

参考文献

- 1) 国立国語研究所 : 婦人雑誌の用語 : 国研報告 4 (1953)
- 2) " : 総合雑誌の用語前編 : 国研報告 12 (1957)
- 3) " : " 後編 : 国研報告 13 (1958)
- 4) " : 現代雑誌九十種の用語用字(1)総記および語彙表 : 国研報告 21 (昭 37-9)
- 5) " : 現代雑誌九十種の用語用字(2)漢字表 : 国研報告 22 (昭 38-3)
- 6) Paul H. Klingbiel : Multimillion Word Data Bases - A Preliminary Report, Volume 1 : Defence Document Center, Distributed by NTIS, AD-777-200 (April 1974)
- 7) 国立国語研究所 : 電子計算機による新聞の語彙調査 : 国研報告 37 (昭 45-2)
- 8) " : 現代新聞の漢字調査(中間報告) : 国研資料 8 (昭 46-2)
- 9) 佐久間 英 : 日本人の姓 : 六芸書房
- 10) 田中 康仁 : 日本人の姓と名に使われた漢字 : 言語生活, 267 (昭 48-12)
- 11) " : 日本人の姓と名の統計 : 言語生活, 254 (昭 47-11)
- 12) 江連 隆 : 漢文教科書の漢字調査 : 計量国語学, No. 52 (1970)
- 13) 田中 靖康 : 全国地名の漢字使用頻度調査 : Computer Report (1973-10)
- 14) 国立国語研究所 : 電子計算機による新聞の語彙調査(IV) : 国研報告 48 (1973)
- 15) Herdan, G. : Language as Choice and Chance : P. Nordhoff N.V., Groningen (1956)
- 16) 松岡 : 語彙出現率曲線の近似法について : 昭 50 年度情処才 16 回大会 (昭 50-11)
- 17) 松岡 : 日本語の語彙統計に関する一考察 : 昭 51 年度電子通信学会全国大会 (昭 51-3)