

白河市言語調査に於ける サンプリング・調査計画

林 知己夫
石 田 正 次

此の研究は科学試験研究費によるものの一部である。国立国語研究所・民族学研究所との共同研究である。国研の中には、特別調査班（委員会）が設置せられ、国語研究所の所員、民族学研究所の研究員、統計数理研究所の所員によつて構成せられた。以下はこの委員会の中で主として統計数理研究所の研究員が行つた成果の一部である。此の研究には外に丸山文行、西平重喜、村岡充子、田熊雅子が参加し、別の部分を担当した。

§ 1. 調査対象の決定と母集団の構成

白河市に居住するものの言語の実態を把握するのが調査の目的であるから調査対象は一先づ昭和24年10月に於ける白河市在住者と云ふことになる。

なお社会生活の調査実施の面を考慮し我々の場合は対象をさらに限定し白河市に居住する15才から69才迄の男女と定めた。この様な人々の織りなす言語生活の Pattern を我々は調査することになるのである。然しこの調査では Pattern のすべてを見極めようとするものではなく一定の基準によつて作られた調査票（個人の経歴、社会環境

* 白河市とは旧白河市を指すものとする。

及び言語関係の調査項目——主として共通語化の程度をみるためのもの——)に基づきその項目にもられているものについて、知識を得ようとするのである。かう考えてくると我々の場合の母集団は次の様に定義するのが妥当であろうと思はれる。

「調査票にもられている調査事項に対する反応を標識とし且つ各個人の抽出確率が等しいと定められた、昭和24年10月白河市に居住する数え年15才から69才迄の男女」

我々の場合は調査は性格上当然面接によらねばならず且つ又微妙な内容をもつ関係上調査員は一定の素養を持ちさらにある程度の訓練を必要とするものであつて、為に適した調査員をさう多く得ることは出来ない。

その他調査にまつわる外的制約(予算、日時、等)を考へ併せるとさう2方に重人とする母集団のすべてを抽出調査することはまづ不可能と言わざるを得ない。ここに当然サンプル調査が行われることになるであろう。我々の場合は云うまでもなく調査の精度(サンプルから母集団を推定する時の)の科学的保証のあるランダムサンプリングの方法によるのである。

§ 2 サンプル数の決定 その第一次的意味調査の実施は純理では行うことはできないのは去るまでもないが、一定の制約の下で最適な結果を得る様理論的に考へてゆかねばならない。

此の調査では調査員は22名程度(東京より出向くもの約15人、現地で得られるもの約7人)と考えられ予算の関係で約4日位に最後の調査を完了しなければならない。調査票の内容よりみて被調査者1名の調査時間は平均約1時間* 従つて調査員1名は1日に7人の調査を行えるものとみてよい。

* 調査時間は長きに亘るとき調査員にも被調査者にも苦痛を興へ調査内容は不確かなものとなる虞れがある。したがつて調査時間が長くとも1時間程度に終る様調査票はつくられてある。

こう考えると調査は4日間に約700人の調査が行えることになる。
 しかし、これは全く順調に調査が進む場合の事であり、実際の場合こうは
 あり得ない。訪問する被調査者は常に在肥するとはかぎらず、しかも
 あくまで此を追究調査しなければならぬ。職場へ行き或は又夜間訪
 向をする場合もあるであろう。しかも調査員は無理を強いてはならぬ
 い。(無理を強いるとき調査内容は不確かなものとなる)ことを思ふと
 サンプル数は高々500程度が限度であろうと思われる。それでは
 500のランダムサンプルで調査の精度は如何であろうか。

調査事項の反応が百分比率であらわされるものとする。サンプルの
 百分比率から母集団の百分比率を推定する時の精度は近似的に、

$$\text{信頼度 } 95\% \text{ で } \pm 1.96 \sqrt{\frac{N-n}{N-1} \frac{p'(1-p')}{n}}$$

によつてあたえられる。但し N は母集団の大きさ我々の場合 17890、
 n はサンプルの大きさ、我々の場合 500、 p' はサンプルに於ける調査
 事項のある反応を示す比率である。

見透をよくするために信頼巾を示す。

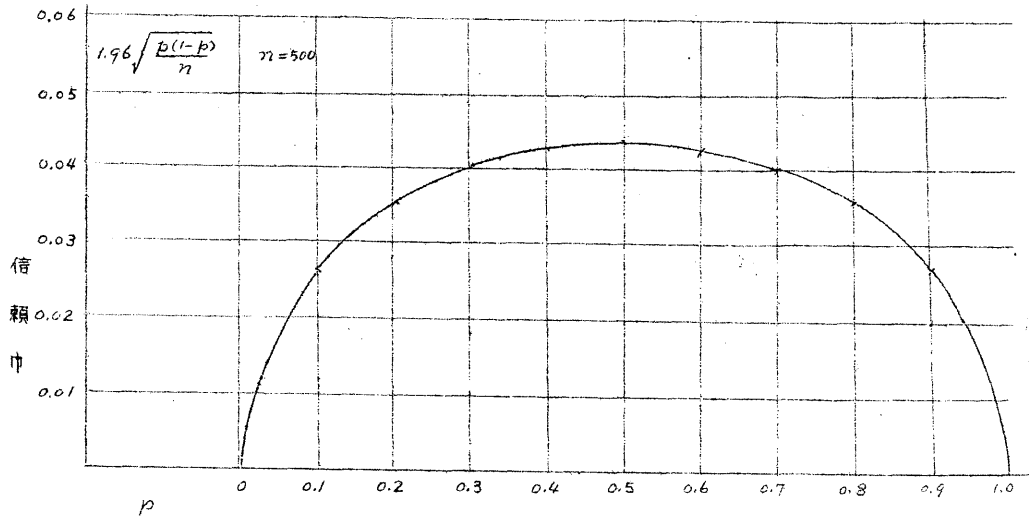
$$1.96 \sqrt{\frac{N-n}{N-1} \frac{p(1-p)}{n}}$$

と p との関係を図示してみよう。こゝに p は母集団の比率を表す
 ものとする。

我々の場合 N が n に比して十分大であるから

$$1.96 \sqrt{\frac{p(1-p)}{n}}$$

を用いた。



此から直にわかる様に500のサンプルがあればまず満足すべき結果が得られるであろう。

然し我々は更に精度をあげるために層化を行う事を考えよう。

§ 3. 層化とサンプルの第一次抽出

我々の場合調査事項の反応形に強い影響を與えたと考えられるものとしては性、年齢、学歴、職業、居住区、家庭環境等々である。

したがって母集団を此等の標識を用いて層化するのがよい。然し層化の客観的資料として容易に用い得られるものは、性、年齢、職業、現住所である。従つて此等を用いて層化を行はうと思う。

このためにはまず母集団を表すリストがなくてはならない。此のリストとして最も適当なものは物資配給台帳である。然し此には今の場合対象外のものをも含むであるし、一世帯が同一頁に記載されてあり一人一人別になつて居らない。従つて此を層化するためには対象たる17890人のカードを新につくり此を層化する手続をとらねばならない。此の事は事実上*不可能である。それは如何にすればよいか。

重ね板きの方法を用いてみる。

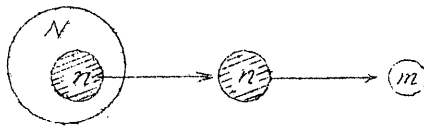
註次頁

つまり第1次的に母集団からある人数をランダムに抽出、しかる後総する500人のサンプルをランダムに抽出する方法(此の時層化をほどこして)を採用することにする。この様な方法をとるとき欲するサンプルを得られることを証明しながら以下サンプルの抽出を述べてみよう。

(i) 母集団の大きさを N とする。その標識は X_1, X_2, \dots, X_N であらわされ、又母集団分散は σ^2 によつて示されるものとする。 N 個からなる母集団からまず n 個のサンプルを抽出する。この時サンプル平均 $\bar{x} = \frac{1}{n} \sum x_i$ の精度を示す分散 s^2 はどうか、なおこの \bar{x} は母集団平均の偏りのない推定量

になつてゐるのは云うまでもない。

此は副次抽出であるから



$$s^2 = \frac{n-m}{n-1} \frac{\sigma_w^2}{m} + \sigma_b^2 \quad \text{となる。}$$

但し σ_w^2 は内分散 σ_b^2 は外分散であり次の関係をもつ

$$\sigma^2 = \sigma_w^2 + \sigma_b^2$$

$$\sigma_b^2 = \frac{N-n}{N-1} \frac{\sigma^2}{n}$$

したがつて

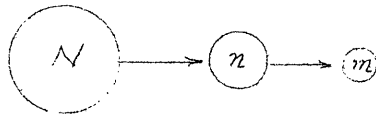
$$s^2 = \frac{n-m}{n-1} \frac{1}{m} (\sigma^2 - \sigma_b^2) + \sigma_b^2$$

註. *日時労力の上から不可能であり又莫大な犠牲をばらい此を実行することは調査の全体的構想よりみる時意味がない。

$$\begin{aligned}
&= \frac{n-m}{n-1} \frac{1}{m} \sigma^2 + \sigma_b^2 \left(1 - \frac{n-m}{n-1} \frac{1}{m} \right) \\
&= \sigma^2 \left(\frac{n-m}{n-1} \frac{1}{m} + \frac{N-n}{N-1} \frac{1}{n} \cdot \frac{n(m-1)}{n-1} \frac{1}{m} \right) \\
&= \frac{N-m}{N-1} \frac{\sigma^2}{m}
\end{aligned}$$

となる。此は N 個の要素を持つ母集団から m 個のサンプルを直接ランダムに抽出した時の分散である。即ち二度ぬきしても精度のかわらない率を示してある。此処で大切なのは S^2 が第1次抽出の個数 n に全く依存しないと言う事である。

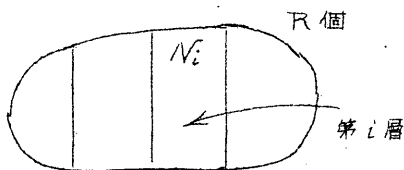
(ii)



前の場合と同様であるが、但し n が確率変数であるときどうなるか。(但し $n \geq m$)

先づ n を個定めて S^2 をもとめてみればよいのであるが S^2 は全く n に関係しないので此の時の分散も亦母集団から直接 m 個のサンプルを抽出した時と同一となる。

(i)(ii)より次の様な考えが導かれる。まづ母集団が R 個の層に分けられたとする。



N_i は i 層の母集の大きさ。今全体から n 個のサンプルをランダムに抽出すると i 層に属するサンプルを n_i だけ得る。

此の n_i は勿論確率変数であるが、 n_i の期待値

$$E(n_i) = n \times \frac{N_i}{N}$$

である。 n_i から欲するサンプル数 m_i 個を抽出 ($n_i \geq m_i$)
 此からサンプル平均をつくる時、その分散は n_i を無視し N_i から
 直接 m_i のサンプルを抽出した時の分散に等しい。

こういうわけであるから全体のサンプル数 m を定め、第 i 層
 の割当数 m_i をきめておき各層で上述の様なサンプリング方法を用
 いて (母集団全体を層化せず、全体よりランダムにとつた第1次サ
 ンプル n を層化しその中から定められた m_i のサンプルをランカ
 ムに抽出する) 普通の場合の様に推定値をつくれれば完全な層化抽出
 法の効果をあげうることになる。

云いかえれば層化抽出を行って効果をあげるのに母集団を層化せ
 ず比較^的少数の第1次サンプルを層化しその中より欲するサンプルを
 抽出しウエイトを用い偏りない推定値をつくれればよいことになる。

層化の手續は著しく簡単になることであろう。

我々の場合、各層へのサンプルの割当は人口比例割当法によるこ
 とにする。蓋し此の割当法は^新機軸な分析を自由に簡単に行いうる
 という利点があるからである。したがって、

$$m_i = m p_i$$

$$\text{但し } p_i = \frac{N_i}{N} \quad \text{とする。}$$

m_i が決めきまるためには p_i がわかつておらねばならない。

白河市の場合残念ながら層化に用い様とする要因すべてについて
 p_i は判っておらない。

p_i は男女 × 年齢についてのみわかっている。即ち、

年	15	20	25	30	35	40	45	50	55	60	65	
令	19	24	29	34	39	44	49	54	59	64	69	
男	9.52	6.84	5.31	4.28	4.06	4.28	3.89	3.18	2.29	1.97	1.28	
女	8.90	7.75	6.74	5.24	5.25	4.76	4.04	3.24	2.70	2.13	1.75	百分率

したがって m_i を明確に定めるのは男女×年令についてであるから第1次層化は男女×年令についてのみ行うことにした。

次に問題となるのは第1次サンプルを n を幾何に定めたいのかの問題である。通常の層化サンプリングと同じ効果をもつサンプリングを行いたいためには常に

$$n_i \geq m_i = n p_i$$

でなくてはならない。此が成立する様 n を定めなければならない。

n_i は確率変数であり

$$E(n_i) = n p_i$$

である。

n_i の分散 $\sigma_{n_i}^2$ は

$$\sigma_{n_i}^2 = \frac{N-n}{N-1} \frac{p_i(1-p_i)}{n}$$

であたえられるから 99.7% (99.9%) の信頼を以て、
 $n_i \geq n p_i$ が成立するためには、

$$n \left(p_i - 3(4) \sqrt{\frac{N-n}{N-1} \frac{p_i(1-p_i)}{n}} \right) \geq m p_i$$

なる様 n が定められねばならない。

我々の場合 $m = 500$ である。一般に n と p との関係を
図示してみると次の様になる。(別紙折込)

此による n としてまづ 2000 程度抽出すればよい事が解るであ
らう。

さらに又、次の様なことも考えておこう。

$m_i = m p_i$ とするとき p_i の小さな層の m_i は小となりすぎ
て i 層そのものの精度があまり落すぎ他との比較が意味をもたなく
なることもあり得るので比例割当法でなく p の小さな所に $m p$ の
2倍のサンプルを割当てることをするとしてよう。

2倍であるから全体の結果を出す場合さう複雑でないであろう。

此が出来うるために n はどうであろうか。

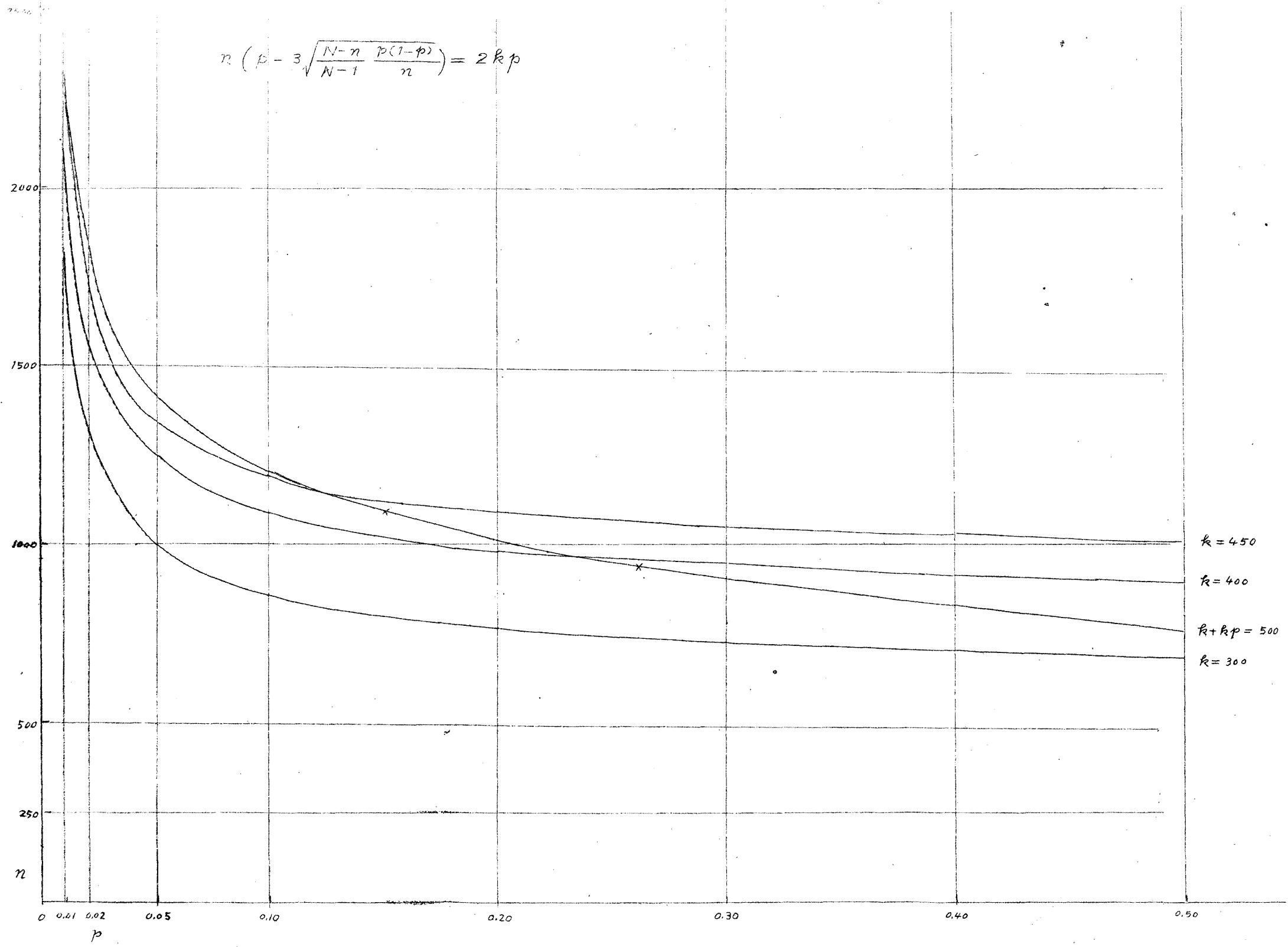
此の時 k 個のサンプルは前のおり比例割当てするものとする。

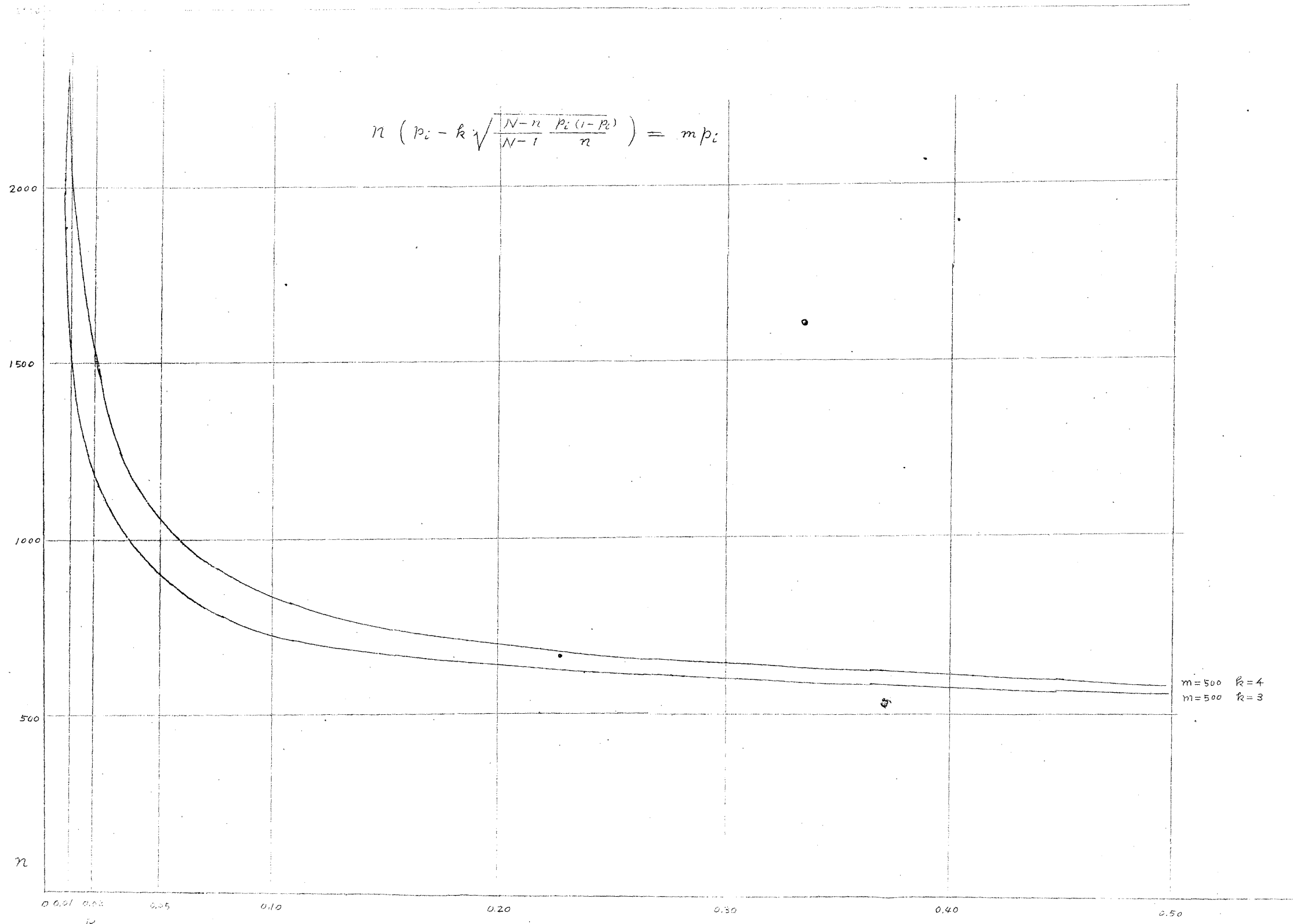
次に、 p の小なところは $2k p$ のサンプルを割当てるものとする
さうすると前と同様に

$$n \left(p - 3 \sqrt{\frac{N-n}{N-1} \frac{p(1-p)}{n}} \right) \geq 2k p$$

でなくてはならない。 k と p との関係を図示してみると次の様
になる。又 $k + k p = 500$ の関係があるとき即ち全サ
ンプルを 500 と定めた時の関係も同時に示しておいた。

$$n \left(p - 3 \sqrt{\frac{N-n}{N-1} \frac{p(1-p)}{n}} \right) = 2kp$$





以上総合し略2000程度の第1次サンプルを抽出すればまづ十分であろうかと思われる。

2000のサンプルを抽出し此をカードにとり層化することは調査の実際的観点よりみて容易に行えるところである。

対象は17890人である。此より2000人を抽出するのであるから抽出の便利及びシステマティックサンプリンクの精度のよさ(よみ書き能力調査に於ける標本調査計画参照)を考へてランダム数を出発番号とし抽出間隔8を以て数えてゆく方法をとつた。

蓋し此の方法は抽出速度がきわめて大きいからである。^{*}

因みに、白河市の平均家族人数は4.9人であり抽出間隔8は此と同類していない。此の結果サンプル数2184を得た。第1次層別を示すと次の様になる。

なお被抽出者を記入するカードには次の様なものを用いた。

				女
町	番地	方		男
				才
				職 業
			家の職業	

註. ^{*} 対象外を除外して数える事は甚だ煩雑であり間違も起し易い

年 令	男		女	
	Sample	Population	Sample	Population
15—19	228	208	176	194
20—24	132	149	205	169
25—29	137	116	140	147
30—34	86	94	107	115
35—39	77	89	105	115
40—44	85	93	99	104
45—49	99	85	79	88
50—54	79	69	81	71
55—59	65	63	60	59
60—64	37	43	45	47
65—69	29	28	33	38
計	1054	1037	1130	1147

こゝろみに母集団の構成と χ^2 検定を行うと

$$\chi^2 = 29.8 \quad \text{自由度 } 21$$

此より大なる χ^2 を得る確率は 10% 程度である。此の 2184 のサンプルは母集団からのランダムサンプルと考えられる程度のものであり強い偏倚を与えていないものと思われる。

§ 4. 第 2 次サンプルの抽出

§ 3 で示した様な各層に属している第 1 次サンプルからいよいよ第 2 次サンプルを抽出することになる。さてこの様な男女 × 年令の層へのサンプルの割当は比例割当によるのであるから各層へのサンプル数は

$$m_i = m p_i \quad m = 500$$

によってあてられるのである。 m_i は次の様になる。

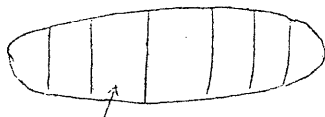
年 令	15 19	20 24	25 29	30 34	35 39	40 44	45 49	50 54	55 59	60 64	65 69	計
男	48	34	27	21	20	21	19	16	15	10	6	237
女	45	39	34	26	26	24	19	16	14	11	9	263
計	93	73	61	47	46	45	38	32	29	21	15	500

此から見てわかる様は明らかに $np_i \geq m_i$ は満足せられている。

各層から m_i 人のサンプルを再びランダム数を出発番号としてステイマティックに抽出する。

此をなす時再び工夫を試みよう。

各層の中を職業、現住所の性格（町単位とし町をその性格により、商業地帯（中心部）、商業地帯（非中心部）、住宅地帯、農業地帯の4段階とした）によつて第2次的に層化し、此を一列に配列し此の中からシ



性 x 年齢一定

ステイマティックに所要数のサンプルを抽出することにした。此の時は完全な層化抽出法ではない。何と云はば

第2次層への割当数は母集団の比率を以てなされるのではなく第1次サンプルの比率を以てなされるからである。かくするとき精度を示す分散はどうなるであろうか、順を追って示してみよう。

§ 5. 此のサンプリク法による精度の計算

(i) 層化抽出による分散（比例割当）

1' 層の母集団の大きさ N_i

2' 層の分散 σ_i

i 層の平均 \bar{X}_i
 i 層のサンプルの大きさ m_i
 母集団の分散 σ^2
 " 平均 \bar{X}
 層の数 R

$$\sum_{i=1}^R N_i = N$$

$$\sum_{i=1}^R m_i = m$$

$$\frac{N_i}{N} = p_i$$

$$\text{又 } m_i = m p_i$$

とする。

此の時

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

の分散

S^2 は近似的に

$$S^2 = \frac{N-m}{N-1} \frac{\sigma^2}{m} - \frac{N-m}{N-1} \frac{1}{m} \sum p_i (\bar{X}_i - \bar{X})^2$$

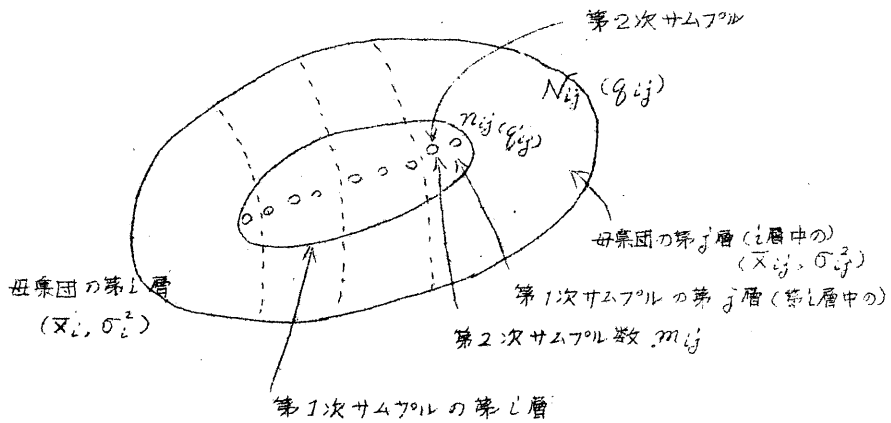
によってあたえられる。

第2項が層化の効果である。

我々の場合 男女×年令 の層化については此の様な精度を得る。但し \bar{x} は推定比率（調査事項に対する反応形体を示す比率）である。重ね抜きサンプルリングを行つても普通の層化抽出法と同一の分散を示す事が證明せられてあるからである。

(ii) 第1次層化による層化の効果

第2次サンプルを層化する。第1次サンプルを層化した時此の層の数を R_i とする。今簡単のため第1層を固定して考える。



i 層中第 j 層の母集団の大きさを N_{ij} とする。

$$\sum N_{ij} = N_i \quad \frac{N_{ij}}{N_i} = g_{ij}$$

とする。

g_{ij} のサンプルの値 (第1次サンプル) を g'_{ij} とする。

$$g'_{ij} = \frac{n_{ij}}{n_i}$$

n_i, n_{ij} は第1次のサンプルの第 i 層に属しているものの数, 第 i 層中の第 j 層に属しているものの数である。

勿論 $n_i = \sum n_{ij}$ である。 g'_{ij} の分散は近似的に安全目に見積って

$$\frac{N_i - n p_i}{N_i - 1} \frac{g_{ij}(1 - g_{ij})}{n p_i} \approx \frac{g_{ij}(1 - g_{ij})}{n p_i}$$

である。

さて此から m 個のサンプルを抽出するのであるが層を一行に並べ此からミステイマティックサンプリクによって抽出するので i 層へのサンプルのサンプル割当は略 $m_i g'_{ij} = m_{ij}$ に等しくなる。

此の様なサンプルから推定値

$$\bar{x}_i = \frac{1}{m} \sum_{k=1}^m x_{ik}$$

をつくと此の分散はどうなるであろうか。

$$\begin{array}{ll} j \text{ 層の母集団平均を } & \bar{X}_{ij} \quad \text{分散を } \sigma_{ij}^2 \\ \text{全体の} & \bar{X}_i \quad \text{" } \sigma_i^2 \end{array}$$

とする。

$$\bar{x}_i = \frac{1}{m} \sum_{k=1}^m x_{ik} = \sum q'_{ij} \bar{x}'_{ij}$$

\bar{x}'_{ij} は j 層のサンプル平均である。

勿論

$$E(\bar{x}_i) = \sum q_{ij} \bar{X}_{ij} = \bar{X}_i \quad \text{である。}$$

(不偏推定値)

\bar{x}_i の分散を考えてみよう。 ことから煩雑をさけるため一応 i の指標をおとして考えてゆこう。

$$\begin{aligned} & E(\bar{x} - \sum q_j \bar{x}_j)^2 \\ &= E(\sum q'_j \bar{x}'_j - \sum q_j \bar{x}_j)^2 \\ &= E(\sum q'_j \bar{x}'_j)^2 - (\sum q_j \bar{x}_j)^2 \\ &= E(\sum q'^2_j \bar{x}'^2_j) - (\sum q_j^2 \bar{x}_j^2) + E(\sum_k \sum_{j \neq k} q'_k q'_j \bar{x}'_k \bar{x}'_j) \\ &\quad - \sum_k \sum_{j \neq k} q_k q_j \bar{x}_k \bar{x}_j \\ &= \sum E(q'_j \bar{x}'^2_j - q_j^2 \bar{x}_j^2) \\ &\quad + E(\sum_k \sum_{j \neq k} (q'_j - q_k)(q'_j - q_j) \bar{x}'_k \bar{x}'_j) \\ &\quad + \sum_k \sum_{j \neq k} q_k q_j \bar{x}_k \bar{x}_j - \sum_k \sum_{j \neq k} q_k q_j \bar{x}_k \bar{x}_j \end{aligned}$$

$$= \sum_{j=1}^R E (q_j' \bar{x}_j' - q_j \bar{x}_j)^2$$

$$+ E \left(\sum_K \sum_{j \neq K} (q_K' - q_K) (q_j' - q_j) \bar{x}_K' \bar{x}_j' \right)$$

第1項の \sum の中

$$E (q_j' \bar{x}_j' - q_j \bar{x}_j)^2$$

$$= E (q_j' \bar{x}_j' - q_j' \bar{x}_j + q_j' \bar{x}_j - q_j \bar{x}_j)^2$$

$$= E q_j'^2 (\bar{x}_j' - \bar{x}_j)^2 + \bar{x}_j'^2 E (q_j' - q_j)^2$$

こゝで再び i の指標を入れ § 3. の結果を用ふれば (此の場合と N_{ij} から抽出推定値をつくる時と分散が同一という事を用うる)
~~直接 m_{ij} を~~

$$E q_{ij}'^2 \frac{N_{ij} - m_{ij}}{N_{ij} - 1} \frac{\sigma_{ij}^2}{m_{ij}} + \bar{X}_j^2 \frac{N_i - n p_i}{N_i - 1} \frac{q_{ij} (1 - q_{ij})}{n p_i}$$

$$= E q_{ij}' \frac{N_{ij} - m_{ij}}{N_{ij} - 1} \frac{\sigma_{ij}^2}{m_i} + \bar{X}_j^2 \frac{N_i - n p_i}{N_i - 1} \frac{q_{ij} (1 - q_{ij})}{n p_i}$$

$$\frac{N_{ij} - m_{ij}}{N_{ij} - 1} \doteq 1 \text{ と考へ } \frac{N_i - m_i}{N_i - 1} \text{ と近似すれば}$$

$$\doteq \frac{N_i - m_i}{N_i - 1} \frac{1}{m_i} q_{ij}'^2 \sigma_{ij}^2 + \bar{X}_j^2 \frac{N_i - n p_i}{N_i - 1} \frac{q_{ij} (1 - q_{ij})}{n p_i}$$

となる。次に第2項について考へる。近似的に安全目にて成立する

$$E (q_{ij}' - q_{ij}) (q_{ik}' - q_{ik}) = - \frac{N_i - n p_i}{N_i - 1} \frac{q_{ij} q_{ik}}{n p_i}$$

なる関係を用うれば、第2項は

$$\begin{aligned}
 & - \frac{1}{np_i} \frac{N_i - np_i}{N_i - 1} \sum_k \sum_{j \neq k} g_{ik} g_{ij} \bar{X}_{ik} \bar{X}_{ij} \\
 = & - \frac{1}{np_i} \frac{N_i - np_i}{N_i - 1} \sum_k g_{ik} \bar{X}_{ik} \sum_{j \neq k} g_{ij} \bar{X}_{ij} \\
 & \sum_{j \neq k} g_{ij} \bar{X}_{ij} = \bar{X}_i
 \end{aligned}$$

したがって

$$\begin{aligned}
 & \sum_{j \neq k}^R g_{ij} \bar{X}_{ij} = \bar{X}_i - g_{ik} \bar{X}_{ik} \quad \text{互用される} \\
 = & \frac{N_i - np_i}{N_i - 1} \frac{1}{np_i} \left(-\bar{X}_i \sum_k g_{ik} \bar{X}_{ik} + \sum_k g_{ik}^2 \bar{X}_{ik}^2 \right) \\
 = & \frac{N_i - np_i}{N_i - 1} \frac{1}{np_i} \left(-\bar{X}_i^2 + \sum_k g_{ik}^2 \bar{X}_{ik}^2 \right) \\
 = & \frac{N_i - np_i}{N_i - 1} \frac{1}{np_i} \left(-\bar{X}_i^2 + \sum g_{ik} \bar{X}_{ik}^2 - \sum g_{ik} \bar{X}_{ik}^2 + \sum g_{ik}^2 \bar{X}_{ik}^2 \right) \\
 = & \frac{N_i - np_i}{N_i - 1} \frac{1}{np_i} \left(\sigma_{bi}^2 - \sum_k g_{ik} (1 - g_{ik}) \bar{X}_{ik}^2 \right)
 \end{aligned}$$

但し σ_{bi}^2 は i 層中の各層間の外分散 (between の Variance) である。

以上綜合すれば、 \bar{X}_i の分散 S_i^2 は近似的に

$$\begin{aligned}
 S_i^2 & \doteq \sum_j \frac{N_i - m_i}{N_i - 1} \frac{1}{m_i} g_{ij} \sigma_{ij}^2 + \sum_j \bar{X}_{ij}^2 \frac{g_{ij}(1 - g_{ij})}{np_i} \frac{N_i - np_i}{N_i - 1} \\
 & + \frac{N_i - np_i}{N_i - 1} \frac{1}{np_i} \sigma_{bi}^2 - \frac{N_i - np_i}{N_i - 1} \sum_k \bar{X}_{ik}^2 \frac{g_{ik}(1 - g_{ik})}{np_i} \\
 & \doteq \left\{ \frac{\sigma_i^2}{m_i} - \sigma_{bi}^2 \left(\frac{1}{m_i} - \frac{1}{np_i} \right) \right\}
 \end{aligned}$$

となる。

常に $\frac{1}{m_i} > \frac{1}{np_i}$ であるから此の層化によつて第2項の分式

分散が小となり得るようになる。

註、

$$S_i^2 = \frac{\sigma_i^2}{m_i} - \sigma_{bi}^2 \left(\frac{1}{m_i} - \frac{1}{np_i} \right)$$

を考へてみる。

もし $\sigma_{bi}^2 = 0$ ならば結果は0となる。

又 $m_i = np_i$ に於ても同様である。母集団の比率 q_{ij} がわかっており $q'_{ij} = q_{ij}$ となしうる場合は

$$S_i^2 = \frac{\sigma_i^2}{m_i} - \frac{\sigma_{bi}^2}{m_i}$$

となりもつとも効果がある場合となる。 $-\frac{1}{np_i}$ の項は q'_{ij} が確率変数であるために起る損失である。

さてそれで全体の推定値

$$\bar{x} = \sum p_i \bar{x}_i$$

はどうなるか、

此は勿論 \bar{x} の偏りのない推定値である。

この時の分散 S^2 は、近似的に

$$\begin{aligned} S^2 &= \sum_{i=1}^R p_i^2 S_i^2 \\ &= \sum_{i=1}^R \frac{N_i - m_i}{N_i - 1} \frac{1}{m_i} p_i^2 \sigma_i^2 - \frac{N - n}{N - 1} \sum \left(1 - \frac{m_i}{np_i} \right) \frac{p_i^2 \sigma_{bi}^2}{m_i} \\ &= \frac{\sigma^2}{m} - \frac{1}{m} \sum p_i (\bar{x} - \bar{x}_i)^2 - \sum \left(1 - \frac{m_i}{np_i} \right) \frac{\sigma_{bi}^2}{m} p_i \\ &= \frac{\sigma^2}{m} - \frac{1}{m} \left(\sum p_i (\bar{x} - \bar{x}_i)^2 + \sum \left(1 - \frac{m_i}{np_i} \right) \sigma_{bi}^2 p_i \right) \end{aligned}$$

によつてあはえられる。層化の効果は第2項によつてあはえられている。即ち完全な層化抽出法でなくとも一般に効果のあることが保證された。

我々の場合 $\frac{m_i}{n p_i} \doteq \frac{1}{4}$ であるから

$$S^2 \doteq \frac{\sigma^2}{m} - \frac{1}{m} \left(\sum p_i (\bar{X} - \bar{X}_i)^2 + \frac{3}{4} \sum \sigma_{bi}^2 p_i \right)$$

とあらはせる。

以上の論議の保證によつて我々は各層の中を層化し欲する天のサムアルをランダムに抽出することにした。

此によつて相当分散を小にせし得、500のサンプルであつても層化しない時のより大きなサンプルに相当する精度をもたせることが出来るであらう。云いかえしてみるならば §2 でのべたよりもさらによい精度を得ることになるであらう。

§ 5. その他

確実にサムアル500を得るために止むを得ざる調査不能者を見込んで各層から補助サムアルを抽出、万一の場合にそなへた。

此の数は116である。

然し此はサンプルの追及調査を忽にしてよいと云ふ事ではない。

上述の様に細かく組むたサムアルであるから、調査はますます厳密にし、あくまでも追及調査し、移転、死亡、重病、長期の旅行等、眞に止むを得ず調査できぬもの丈サムアルの差し替へを行い、予定のサムアル数を確保したいが爲である。

§ 6. 此の様な計画は実際に可能であらうか、一番困難と思われするのは面接に行つた時不在であるものを追及することである。このためあらかじめ在宅率と言ふものをしらべておく事は大切なことである。

準備調査のとき此を行つてみた。(実際には同時に講和問題につき、世論調査を行つて政治的関心をもあわせて調査しておいた。)

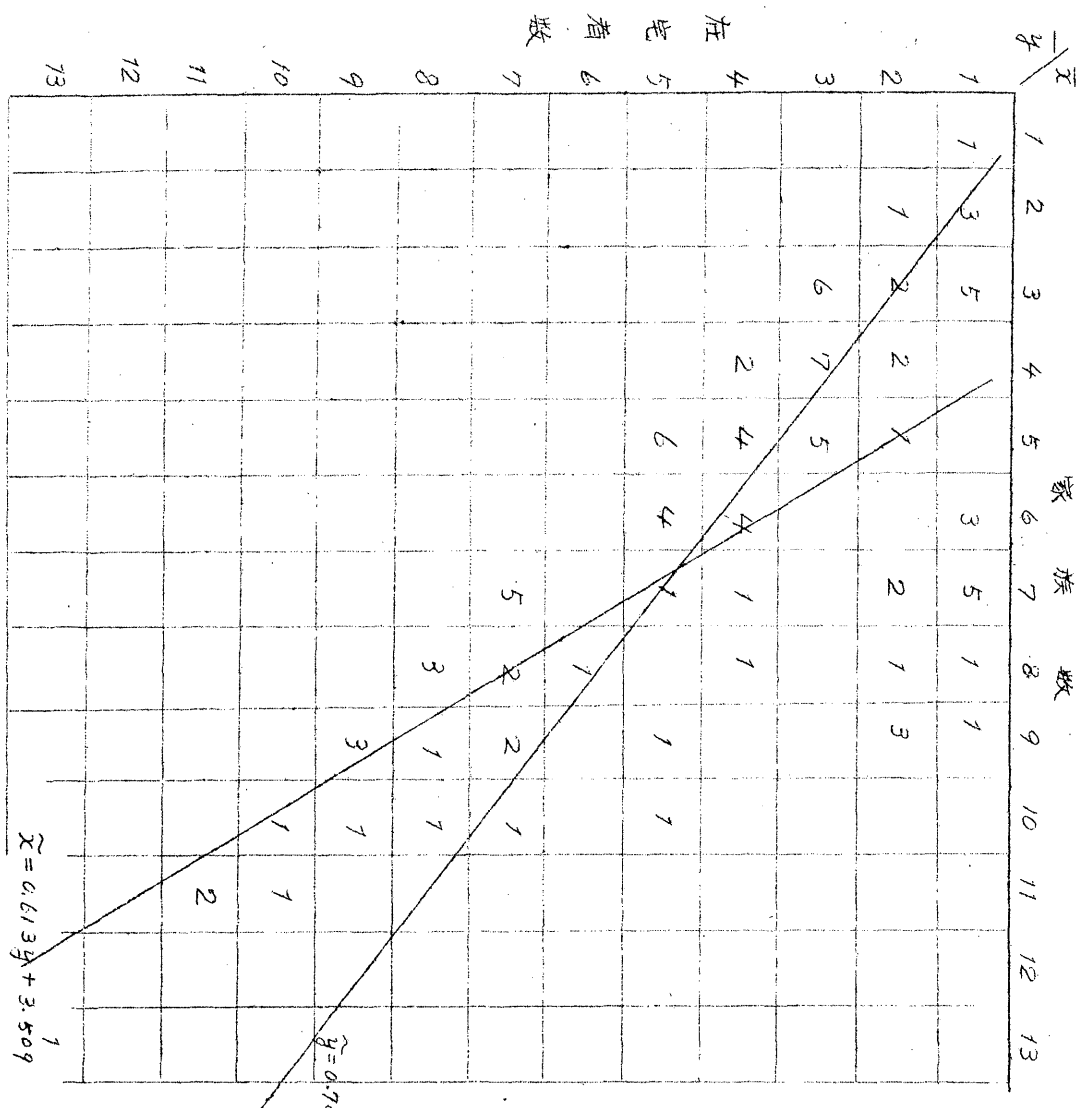
サンプリクの方法は各世帯を調査単位とした。即ち等しい確率を以て世帯を抽出した。

各町を層としてこゝに含まれている世帯数に比例してサンプルを割あてた。この方法によつて今の在宅者と家族人数を面接によつて調査した。こうしてから在宅者率を Regression Estimate を用いて推定することにした。此の方法によるとき夕時—4時迄(此の時間が実際調査した所の時間)の間の平均在宅者率(時間平均*)が求められるであろう。

さて、

家族人数と調査時に於ける在宅者数との相関関係を求めてみると次の様になる。(次頁)

註。*各人の訪門される時刻は等しい確率で決定されるとの確定をしていることにはなる。



住宅者数

族数

$$P = 0.684$$

$$\bar{x} = 6.091$$

$$\bar{y} = 4.212$$

$$\sigma_x = 2.499$$

$$\sigma_y = 2.790$$

$$\hat{y} = 0.7632x - 0.435$$

$$\hat{x} = 0.613y + 3.509$$

y の x に対する回帰直線は

$$\tilde{y} = 0.763x - 0.435 \quad \text{となる。}$$

$\bar{x} = 4.9$ と云ふ事（白河市の一世帯当り平均家族人数）は既知であるから平均在宅者（一世帯当り）が回帰線からよみとれる比は 3.3 である。こゝで \bar{y} の分散を計算してみると、

$$V(\bar{y}) = \frac{\sigma_y^2(1-p^2)}{n} = \frac{7.78(1-0.684^2)}{99} = 0.042$$

となる。したがつてその在宅率は

$$A = \frac{3.3 \pm 0.4}{4.9} = 0.67 \pm 0.08$$

（信頼中 95%）

と推定される。なお此の様な推定が正規分布にしたがふ事が石田によつて説明（講究録、火災推定参照）されてあるのでそれにしたがふ信頼中がつけられてある。

これから本調査の場合 30% の不在者を思はしめるが農閑期に入るからこれよりはいくぶん在宅率は上るであろう。

以上の様な目安によつて期的にみて此の調査計画はまづ無理はないものとみてよい。

§ 7. 附 録

調査結果からみたサムプリングの精度

調査の主な目的は共通語化の程度と言ふ事であつた。その共通語化の程度は音韻の反応によつて測定された。しかゝその程度は一応共通語的反応を示したところの項目の数によつてあらわされることになつた。

したがつて此の標識をもとにしてサムプリングの精度を計算しておこう。勿論推定は母集団平均の推定についてである、即ちサンプル平均*
*の信頼中

まづ、層別の効果について考へてゆこう。

此を出すためには各層の分散を計算してゆかねばならない。しかし重ねぬきをした時考へる層のえのについてはサンプル数が過少であつたため分散を推定することは意味が少いと考へたので男女×年令の層別についてのみ話をすゝめてゆくことにする。重ねぬきの方の時用いた所の層に就ても層の平均の間に標識の上で差のあることは結果の分析から判明してゐるので以上の様な検討は相当控へ目なものであることが言われるであらう。

さて層別した時の分散を計算すると

$$\sigma_{st}^2 = \frac{N-n}{N-1} \frac{1}{n} \sum_{i=1}^R p_i \sigma_i^2$$

$$N = 17290$$

$$n = 476$$

p_i は既述

σ_i は i 層の分散

$$\sigma_{st}^2 = 0.03236 \quad \sigma_{st} = 0.1799$$

層別しないとした時の分散 σ_2^2 は

$$\sigma_2^2 = 0.03854 \quad \sigma_2 = 0.1963$$

となる。層別の効果は、標準偏差の意味で

$$\frac{\sigma_n - \sigma_{st}}{\sigma_n} = 0.08$$

となる。即ち約 8% 効果があった事になる。

層別の効果が比較的少なかったのは意外であった。

さて母集団平均の推定値 \bar{x} は

$$\bar{x} = 19.6$$

であるから、サンプル平均の変異係数をもとめると

$$\frac{\sigma_{st}}{\bar{x}} = 0.01$$

となり、サンプル平均の精度は相当高いと言われるであろう。

何とならば、信頼中をしめす相対精度

$$k \frac{\sigma_{st}}{\bar{x}}$$

(k は此の場合、信頼度 95% で 1.96 と言へてよい) は、2% 以下となるからである。

サンプル平均の結果はまず満足すべきものと思う。