

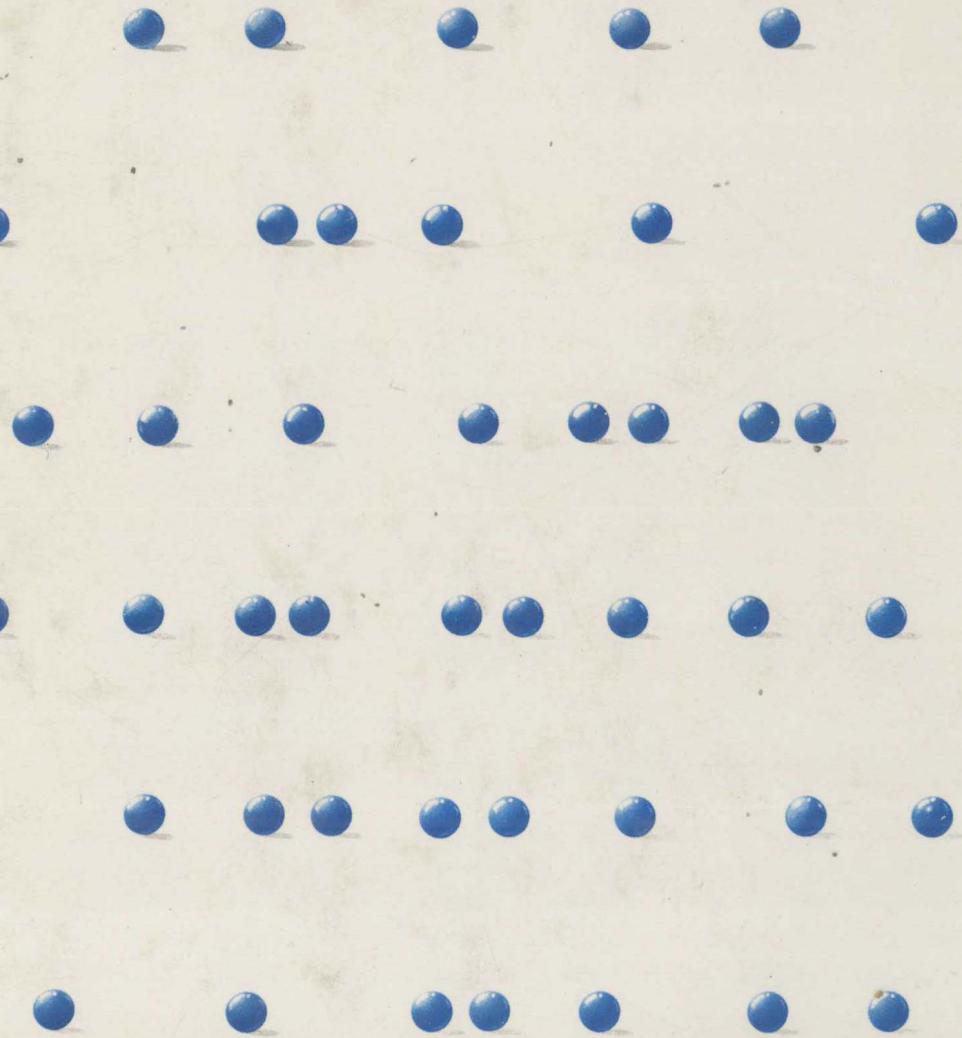
統計学入門

香川大学教授 木村 等

香川大学教授 大藪和雄

香川大学助教授 石川 浩

共著



香川大学教授 木村 等
香川大学教授 大藪和雄 共著
香川大学助教授 石川 浩



著者略歴

木村 等

1924年2月15日生
1947年9月 北海道帝国大学理学部数学科卒業
統計数理研究所等を経て
1959年4月 香川大学教授経済学部、現在に至る。

大藪 和雄

1937年10月1日生
1966年3月 一橋大学大学院経済学研究科修士課程修了
香川大学講師、助教授を経て
1980年10月 香川大学教授経済学部、現在に至る。

石川 浩

1941年11月29日生
1969年3月 京都大学大学院工学研究科博士課程修了
米国コロムビア大学客員研究員等を経て
1982年4月 香川大学助教授経済学部、現在に至る。
工学博士

統計学入門

NDC 417

1982年10月20日 第1刷発行 定価 1900円

木 村 等
著 者 大 藩 和 雄
石 川 浩
発 行 者 宇 野 豊 藏
印 刷 中 央 印 刷 株 式 会 社
製 本 株 式 会 社 若 林 製 本 所

発行所 実教出版社
東京都千代田区五番町五番地 〒102
電話 東京(263)0111(大代表)振替東京 4-183260

© H. KIMURA, K. OHYABU,
H. ISHIKAWA

1982

3041-2228-3205

まえがき

本書は大学初年級の学生を対象とする統計学の入門書として執筆したものである。草稿の段階で本学部学生の多くに目を通してもらい、彼らの初学者としての立場からの率直な意見を取り入れて稿を改めた。社会科学系の学部学生は数学をあまり得意とはしないが、本書の程度まで懇切に書けば、数理的な面も理解できるようである。しかしながら、勿論、標本分布の導出の積分などは、統計学自体の理解のためには必ずしも必要ではないので、数学の不得意な読者はその部分を飛ばしてもらって支障はない。要するに統計的な考え方を理解して欲しいのである。また、その意味で、各章末には基本的な演習問題を配し、巻末にその全解を入れた。読者は、この問題を、まず自ら解き、解答と参照することによってより深い理解を得ることと信ずる。

本書の構成は2部構成とした。第I部は記述統計に関するものであり、統計データの取り扱いの基本的手法について述べた。一方、第II部は統計的推論に関するもので、初めに確率を Frequency Theory としてとらえ、これから確率分布、期待値、分散の概念を構成し、さらに母集団と標本の概念、推定および検定の諸手法について述べた。

初学者に対する統計学の入門書としての著者らの目論見がどの程度達成されているかは、大方の厳正な御批判に待つ他はないが、多少なりとも学習の一助になりうるものとすれば望外の喜びである。

本書の出版に当たり、資料の準備段階から原稿の浄書、通読など、本学部学生諸君には多大の労をわざらわした。とくに本常功君、小出智子嬢には記してその労を多としたい。また実教出版の桜井瑞穂氏には、本書の企画段階から校了に至るまで、一貫して多大の御尽力を賜わった。同氏の御力添えがなければ本書が日の目をみることはなかつたであろうと思われる。ここに記して深甚の謝意を表する。

目 次

第Ⅰ部 記述統計

第1章 統計データ	2
1・1 はじめに (2)	
1・2 質的データと量的データ (3)	
1・3 静態統計と動態統計 (5)	
1・4 官庁統計と民間統計 (7)	
1・5 第1義統計と第2義統計 (7)	
1・6 1次統計と2次統計 (8)	
1・7 全数調査データと標本調査データ (8)	
1・8 タイム・シリーズ・データ（時系列データ）とクロス・セクション・データ（横断面データ） (9)	
1・9 行政区域統計とメッシュ統計 (11)	
1・10 実績統計と予測統計 (12)	
1・11 アクチュアル・データとユージュアル・データ (12)	
第2章 度数分布	14
2・1 度数分布表および度数分布図 (14)	
2・2 累積度数分布 (21)	
2・3 度数分布の形状 (22)	
第2章 演習問題	25
第3章 代表値	27
3・1 算術平均（相加平均） (27)	
3・2 幾何平均（相乗平均） (31)	
3・3 調和平均 (33)	
3・4 メディアン（中央値または中位数） (33)	
3・5 モード（最頻値） (35)	
第3章 演習問題	37

第4章 散布度（分散度）、歪度および尖度	38
4・1 散布度（分散度） (38)	
4・2 平均偏差 (39)	
4・3 分散および標準偏差 (40)	
4・4 範囲および四分位偏差 (43)	
4・5 相対的な散らばりの測度 (45)	
4・6 歪度と尖度 (47)	
第4章 演習問題	50
第5章 相関係数 (51)	
5・1 相関係数 (51)	
5・2 相関係数の計算 (55)	
5・3 相関係数の意味 (58)	
第5章 演習問題	61
第6章 統計的比例数 (62)	
6・1 構造比例数 (62)	
6・2 関係比例数 (62)	
6・3 指 数 (64)	
6・3・1 指数の概念 (64)	
6・3・2 総合指数の作り方 (65)	
6・3・3 指数算式の吟味 (68)	
6・3・4 指数作成のその他の問題点 (69)	
第6章 演習問題	70

第Ⅱ部 統計的推論

第7章 確 率 (72)	
7・1 序 説 (72)	
7・2 確率論のための基礎的な用語 (77)	
7・3 頻度論的確率論 (79)	
7・3・1 コレクティフと確率の考え方 (79)	
7・3・2 基本操作と確率の計算法則 (82)	
7・3・3 乱数表 (86)	
7・4 確率変数、確率分布、期待値および分散 (89)	

第7章 演習問題	97
第8章 基本的な確率分布	98
8・1 二項分布 (98)	
8・2 正規分布 (103)	
8・2・1 正規分布 (103)	
8・2・2 二項分布の正規分布による近似 (111)	
8・3 ポアソン分布 (112)	
第8章 演習問題	116
第9章 統計的推定	117
9・1 母集団と標本 (117)	
9・2 有限母集団の場合の区間推定 (119)	
9・2・1 平均値の推定 (119)	
9・2・2 比率の推定 (127)	
9・3 正規母集団における母数の推定 (131)	
9・3・1 平均の推定 (136)	
9・3・2 分散の推定 (137)	
9・4 二項分布における母数の推定 (139)	
9・5 その他の推定法 (140)	
9・6 ベイズの定理の推定問題への応用 (142)	
9・6・1 ベイズの定理 (142)	
9・6・2 経験的ベイズ推定 (143)	
9・6・3 個人確率を用いるベイズ推定 (147)	
9・6・4 決定関数の考え方によるベイズ推定 (151)	
9・7 おわりに (153)	
第9章 演習問題	154
第10章 統計的仮説検定	155
10・1 仮説検定の考え方 (155)	
10・2 2種類の誤り (157)	
10・3 平均の検定 (160)	
10・4 等分散の検定 (161)	
10・5 平均の差の検定 (164)	
10・5・1 2つの母集団の分散が等しい場合 (164)	

10・5・2 2つの母集団の分散が異なる場合 (165)	
10・5・3 標本数が大きい場合 (167)	
10・6 比率の検定 (168)	
10・7 比率の差の検定 (168)	
10・8 適合度検定 (170)	
10・9 独立性の検定 (172)	
10・10 小標本における 2×2 分割表と χ^2 検定 (176)	
10・11 ランダムネスの検定 (180)	
第10章 演習問題	183
〔数学的補遺〕	185
ガンマ関数, ベータ関数およびスターリングの公式	
演習問題解答	189
索引	229
〔付表 I~VII〕 (215~228)	
乱数表, 二項分布, ポアソン分布, 標準正規分布 $N(0, 1)$, t 分布, χ^2 分布, F 分布 (5% 点の表, 1% 点の表)	

第Ⅰ部 記述統計

第1章 統計データ

第2章 度数分布

第3章 代表値

第4章 散布度(分散度),

歪度および尖度

第5章 相関係数

第6章 統計的比例数

現実の数量的な側面をとらえるために、統計調査を行い、その結果を統計表の形にまとめる。統計表にまとめられたデータを用いて、そのデータの性質を直観的にとらえるためにグラフを描く。またそのデータを縮約した数値を求めて分析を進める。このような目的に用いられる方法が記述統計である。第Ⅰ部では、この記述統計について基本的事項を概説する。

第1章 統計データ

1・1 はじめに

統計データは通常、次のような過程を経て作成される。

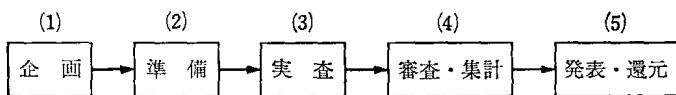


図 1・1 統計調査の作成過程

(1)の企画は、最も重要な過程である。調査目的をはっきりさせ、調査対象を定義し、具体的調査方法を選択し、調査の時と場所を決めなければならぬ。換言すれば、いわゆる統計調査の 5W を決めなければならない。ここに 5W とは、カウフマン(A. Kaufmann)によれば、調査対象(was), 調査方法(wie), 調査時点(wann), 調査の場所(wo), 調査の主体(wer)をいう。なお、人によってはさらに調査目的(worum)を加えて 6W とするときもある。

(2)の準備には、調査予算の獲得、調査員の選任、訓練、調査票の作成、調査用品の準備が含まれる。

(3)の実査は、調査対象と接触して情報を集める段階である。

(4)の審査・集計は、調査票の内容を検討し、内部矛盾や不備を見出し、調査対象に再度問い合わせて修正する審査の段階と、審査の終わった個票を統計表にまでまとめあげる集計の段階とを含んでいる。

(5)の発表・還元は、内容の公表と、印刷公刊、調査対象に対する得られた情報の周知還元をする段階である。

以上のような統計調査の実際を個別に知ることによって初めて、分析の対象となるデータの性格、利用の仕方がわかるのである。

統計データの性格、利用の仕方、データの体系、各種統計の比較などを、より一般的に考察する「統計データ論」というものが考えられるが、ここでは、統計データの各種の分類について述べ、データの性格についての一面にふれることに留める。

1・2 質的データと量的データ

統計データの分類の1つとして、質的データと量的データという分類がある。質的データとは、男女別人口、職業別就業者数などのように、分類の標識——集団を構成する個々の単位は、多種多様な性質をもっているが、これらの性質のうち、統計的観察の対象となるものを標識といふ——が、質的である場合であり、量的データとは、世帯人員別世帯数、住居の室数別世帯数、世帯の収入階級別世帯数、身長階級別学生数などにより分類の標識が量的である場合である。

量的データは、離散的なものと連続的なものとに分かれている。離散的とは1, 2, 3, …のように、とびとびの数として表される場合であり、連続的とは本来実数として表される場合である。世帯人員別世帯数、住居の室数別世帯数などは標識が離散的であるので離散的量的データと呼ばれる。一方、世帯の収入階級別世帯数、身長階級別学生数などは標識が連続的であると考えられるので連続的量的データと呼ばれる。数学的にいえば、世帯の収入は円の単位で測られており、本来は離散的であるが、1円は何百万円、何千万円に比べて小さいものであるから、統計的な取り扱いにおいては連続的なものとみなしている。

なお、これに関連して、計数統計と計量統計の区分がある。計数統計は質的データと離散的量的データを合わせたものをいい、計量統計は、連続的量的データに対応したものである。

質的データと量的データの区分については、上述の通りであるが、実際には、質的データを量的データで規定する場合もある。例えば、大企業、中小企業の区分は本来質的なものであるが、中小企業基本法によれば、鉱工業においては資本金1億円以下の会社、あるいは従業員300人以下の企業としている。この場合は質的な中小企業という概念を資本金あるいは従業員という量的なデータで規定している。

質的データの分類については、次の各種の統計用分類が用いられている。

- (1) 日本標準職業分類
- (2) 日本標準産業分類
- (3) 日本標準商品分類

4 第1章 統計データ

- (4) 輸出入統計品目表
- (5) 日本標準建築物用途分類
- (6) 疾病、傷害および死因の統計分類

このうち産業分類を例にとって、分類の問題を考えてみよう。産業分類では、分類の単位として事業所をとっている。事業所(establishment)とは、「物の生産またはサービスの提供が業として行われている個々の物理的な場所」であり、一般には工場、製練所、鉱山、商店、農家、病院、事務所などと呼ばれる、「一区画を占めて、経済活動を行っている場所」のことである。事業所の産業は、そこで行われている主要な業務によって決定することが原則であるが、事業所によってはいくつかの異なる経済活動を行っている場合もあり、このような場合には、原則として1ヵ年間の総収入額または総販売額の最も多い産業をその事業所の産業とする。

ところで、産業と商品とは密接に関連しているが、工業統計表をまとめの場合に、以下に示すように『産業編』と『品目編』の違いが生じる。仮想例として、いまM町に事業所がA, B, C, Dの4つあり、それぞれ、みそ、しょう油、清酒などを正在しているものとし、ある年の出荷額が表1・1のようであったとする。表1・1から、前述の事業所の産業分類の考え方従えば、A事業所は産業としてはみそ製造業として格付けされ、B事業所はしょう油製造業、C, D事業所は清酒製造業として格付けされる。したがって、この町の工業統計の『産業編』の数字は表1・2のようになる。その理由は、A事業所の主要な業務はみその製造であるから、A事業所の出荷額の総額はみそ製造業の出荷額に分類され、B, C, D事業所についても同様に取り扱われるからである¹⁾。また、表1・1から『品目編』は表1・3のように作られる。例え

表1・1 産業と品目の違い (単位 万円)

事業所\生産物	みそ	しょう油 食用アミノ酸	清酒	計
A	4 600	250	—	4 850
B	500	9 450	1 200	11 150
C	400	—	3 750	4 150
D	—	—	1 500	1 500
計	5 500	9 700	6 450	21 650

1) 統計法による秘密の保持のため、工業統計表ではある分類に1あるいは2事業所しか存在しない場合はxで表示することになっているので、上記の数字は全部xとなる。

表 1・2 工業統計表『産業編』(M町)

産業	出荷額(単位万円)
みそ製造業	4 850
しょう油・食用アミノ酸製造業	11 150
清酒製造業	5 650

表 1・3 工業統計表『品目編』(M町)

品目	出荷額(単位万円)
みそ	5 500
しょう油・食用アミノ酸	9 700
清酒	6 450

ばこの 2 つの表からみその出荷額についてみれば、みそ製造業の出荷額は現実のみその出荷額とは異なるものとなっていることがわかる。これは単なる仮想例ではあるが、実際の統計にもこのようなことがあるので統計表を見る際には十分注意しなければならない。このような問題は、現実にある事業所を 1 つの産業に分類することは本来不可能であるにもかかわらず、それをあえてすることから生ずるのであって、質的な分類でさえもこのような困難さを含んでいるのである。分類に関するもう 1 つの問題として、分類が現実の 1 つの側面をとらえてなされるために、他の側面を捨象してしまうことがある。少数の対象であれば、多方面から 1 度にとらえることもできるかも知れないが、統計調査が取り扱うような大量な対象については、人間の能力をもってしては、限られた側面からしかみることができず、その取り上げた側面についての分類を行うことによって情報を得ようとするのが統計学であって、この際他の面の情報が捨てられるのはやむをえないこととしているのである。

1・3 静態統計と動態統計

統計データの分類の 1 つとして、静態統計と動態統計の区別がある。対象が時間的に持続性をもつ場合、例えば、人口、資本ストック、国富などは、一時点とらえざるをえない。このような統計を静態統計という。他方、ある種の対象は瞬間的な出来事として存在し、持続性をもっていない。この種のものをとらえるためには、ある一定の期間を区切って観察しなければならない。例えば、出生、死亡、結婚、離婚、転入、転出、設備投資、生産高といったものは、このような性質をもっており動態統計と呼ばれている。

いま、静態統計と動態統計の違いを人口と出生・死亡を例にとって示すために生命線の図、図 1・2 を作る。○印は出生、●印は死亡を表すものとし、それらを結ぶ線分を生命線という。出生数をとらえるためには、出生が瞬間

6 第1章 統計データ

的な出来事であるから、ある期間を設けて、その期間内にある出生を数えなければならない。例えば、 t 年10月1日から $(t+1)$ 年9月30日までに含まれる○印の数を数えれば、4であるから、この1年間の出生数は4となる。他方、死亡は、この同じ期間に、2人について起こっているから、この期間の死亡数は2である。また、人口を把握するためには、人口

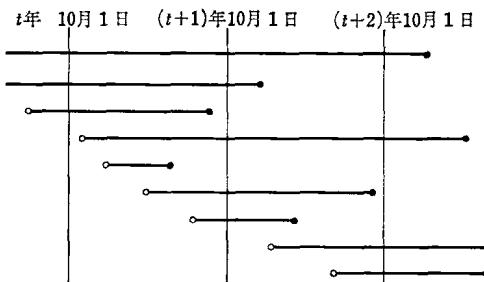


図 1・2 生命線と人口(静態統計と動態統計)

というものが持続性をもっていることから、例えば、 t 年10月1日を示す縦の線と生命線との交点を数えて、この場合3が t 年10月1日現在の人口であるとする。

さて、次に、静態統計と動態統計との関連を考えてみよう。

P_t : ある県の t 年10月1日の人口,

P_{t-1} : ある県の $(t-1)$ 年10月1日の人口,

B_t : ある県の $(t-1)$ 年10月1日から t 年9月30日までの出生数,

D_t : ある県の $(t-1)$ 年10月1日から t 年9月30日までの死亡数,

I_t : ある県の $(t-1)$ 年10月1日から t 年9月30日までの転入者数,

O_t : ある県の $(t-1)$ 年10月1日から t 年9月30日までの転出者数,
とすれば、これらの間には

$$P_t = P_{t-1} + B_t - D_t + I_t - O_t \quad (1 \cdot 1)$$

の関係が成り立つ。 $B_t - D_t$ は自然増加、 $I_t - O_t$ は社会増加と呼ばれている。ある県について、ある年の国勢調査人口を出発点として、将来の人口を推計しようとする場合を考えよう。初めに、 P_{t-1} にその年の国勢調査人口を代入する。 B_t は母親の年齢階級別特殊出生率(いわば、女性が子供を産む確率である)からある程度正確に推計することができるし、 D_t は生命表の死亡率を用いて相当正確に推計できる。難しいのは、転入・転出の推定である。過去の動向から、それらがどのように変動するかを予測する、この部分がうまく予測できさえすれば、 P_t の予測が可能となる。このようにして1年ごとに P_t を推定することを繰り返して将来の予測を行うことができる。逆に、

現時点から過去にさかのぼって考えると、 P_{t-1} は、 $P_{t-2}, B_{t-1}, D_{t-1}, I_{t-1}, O_{t-1}$ などから決まっている。このように、時点をさかのぼって考えると、結局、現在の人口は、過去の出生、死亡、転入、転出によって決定されていることがわかる。このような意味で、静態統計は動態統計の積み上げによって成り立っているということができよう。静態統計は経済学でいうストックに、動態統計は経済学でいうフローに対応している。

1・4 官庁統計と民間統計

調査主体によるデータの区分として、国や地方公共団体およびその他の公的機関が作成する統計を官庁統計、それ以外の調査主体が作成する統計を民間統計と呼んでいる。わが国の大部分の統計は官庁統計であり、民間統計には、労働組合の統計、民間業界の統計、新聞社・放送局などの統計、専門調査機関の統計などがある。調査主体によって調査目的、回収率なども異なり、結果の数字の正確さや意味が異なる場合がありうる。

1・5 第1義統計と第2義統計

これは調査の目的と調査方法に関する区分といってよいであろう。第1義統計は、本来統計をとる目的でとられた統計であり、本格的な統計といつてもよいであろう。これに対して、第2義統計は業務統計ともいわれるよう、他の目的、例えば、行政目的のために集められた資料をまとめて統計にしたものである。

建築基準法に基づく届出書から作成される「建築着工統計」、戸籍法に基づく各種届出書を調査票に転記することにより作成される「人口動態統計」、税関に提出される輸出入申告書などから作成される「貿易統計」、警察活動の各段階で作成される犯罪原票の結果を集計した「犯罪統計」、職業安定所の「職業紹介の統計」、日本国有鉄道の「輸送統計」などは、第2義統計の例である。

最近、統計調査の数が増えたこともある、調査される側の負担が増していること、および、プライバシーを重視する傾向が強くなつたために、調査される側がより慎重に対応するようになったこと、などのために、調査が困難になりつつある。このような情況の下で、業務統計の重要性が再認識され

始めている。例えば、行政管理庁では、サービス業に関する業務統計の整備を手がけようとしているが、今後、各方面での業務統計の発掘・整備が強く望まれる。

1・6 1次統計と2次統計

これは調査方法による区分である。1次統計がもともと得られた統計データであるのに対して、2次統計は加工統計といわれるよう、1次統計を加工して作ったデータである。2次統計は各種の経済指標を含んでいる。例えば、総理府統計局が発表している「消費者物価指数」 P_{ot} は、 p_{ot} を*i*品目の基準時の価格、 p_{ti} を*i*品目の比較時の価格、 $p_{ot}q_{ot}$ を基準時の*i*品目への支出金額、 n を採用品目の数として、次式

$$P_{ot} = \left(\sum_{i=1}^n \frac{p_{ti}}{p_{ot}} \cdot p_{ot}q_{ot} \right) \times 100 \quad (1 \cdot 2)$$

で計算されているが、 p_{ti}/p_{ot} の部分は「小売物価統計調査」という1次統計のデータを用いており、ウェイト $p_{ot}q_{ot}$ の部分は「家計調査」という1次統計のデータを用いている。このように、消費者物価指数は2次統計である。また経済企画庁が公表している「景気動向指数」は「鉱工業製品在庫率指数」、「卸売物価指数」、「原材料在庫指数」、「鉱工業生産指数」などの2次統計と、機械受注、輸入通関実績、全銀預貸率などの1次統計を用いて作成した2次統計である。また、国民所得勘定、産業連関表、資金循環表、国民貸借対照表、国際収支表といった社会会計上のデータは、非常に多くの1次統計を総動員して作成される2次統計の例である。

2次統計の正確さやデータの性質を知るためにには、そのもとになっている1次統計の正確さや性質を知る必要がある。また、2次統計に1次統計を組み込む方法の吟味が必要であるのはもちろんである。

1・7 全数調査データと標本調査データ

これも調査方法についての区分である。全数調査データは悉皆調査データとも呼ばれており、調査対象をことごとく調査して得られたデータであり、標本調査データは、調査対象の一部を調査して得られたデータである。

この2つのデータには一長一短がある。全数調査データは、調査対象全体

を調査して得られたものであるから、誤差を含まず正確なものであると思われている。しかしながら、調査の実施に当たっては、調査項目に関する概念規定の誤り、調査票などの設計上の誤り、調査対象名簿のもれや重複による誤差、申告者の意識的無意識的誤りなどを避けることができない。また、調査員を数多く動員しなければならず、主旨の徹底を欠くことなどに起因した調査員の誤解による誤差、調査員の記入誤りによる誤差なども生じる。その他、集計の際にも誤記・転記などの誤差、計算上の誤差などを生じる。以上述べたようないわゆる非標本誤差は、全数調査においても避けることができず、かえって誤差が大きくなる場合もあり得る。また、調査費用も膨大になり、集計時間も多くかかる。しかしながら全数調査は、標本調査のための基礎資料を与えるものであり、また数が多いことから精密な分析が可能であるという利点をもっている。

他方、標本調査は、調査対象の一部を調査するものであって、調査数も比較的少ないとことから、調査費用も安く、調査員の訓練も十分に行き届き、集計時間も少なくすみ、標本誤差は存在するが、その大きさをある程度管理できることもある、多くの統計が標本調査によって作られている。

全数調査のデータとしては、国勢調査、事業所統計調査、農(林)業センサス、工業統計調査、商業統計調査などがある。この中で、国勢調査は、特に調査が膨大であって、多大の費用がかかることから5年ごとに行われている。その間の時期における労働力の状況を把握するために、労働力調査、就業構造基本調査といった標本調査が実施されている。

1・8 タイム・シリーズ・データ(時系列データ)とクロス・セクション・データ(横断面データ)

調査の「時」に関連する区分である。タイム・シリーズ・データ (time-series data)は、時間の順序に並べられたデータを意味している。これに対して、1つの時点あるいは1つの期間について、階層別あるいは地域別などに分割されたデータをクロス・セクション・データ(cross-section data)という。

図1・3はこの関係を示すために、可処分所得を横軸にとり、消費支出金額を縦軸にとって作ったグラフである。実線は昭和45年と昭和55年のクロス・セクション・データであり、A点は昭和45年のクロス・セクション・