

電卓シリーズ (4)

相関係数の計算と統計図のプログラム

吉井 守正 (鈦 床 部)

今回は 多量なデータを処理する方法として筆者が考えたものについて述べ その具体例として相関係数と回帰直線を求める計算とプロッタを使ってのデータの図化とくにヒストグラム・相関図・三角図について記す。

1. 多量なデータを処理するための工夫

データ数が多くなると データの入出力に手間取るといふほかに さまざまな問題が派生して来る。そこで実際に経験した事柄をもとに 多量なデータへの対応策をいくつかご紹介し ご批判をおおぎたい。

1.1. 磁気テープを使う場合

1.1.1. ファイルの構成

データを磁気テープなどのいわゆる外部メモリーに貯えることによって 計算機本体のメモリー容量の何倍ものデータを処理する事ができる。筆者が使用している

YHP-20 計算機にも昨年からカセットテープレコーダが付属し 多量なデータをさばくのに大いに役立っている。そこで磁気テープの割り付けや使い方のコツなどつぎに示そう。

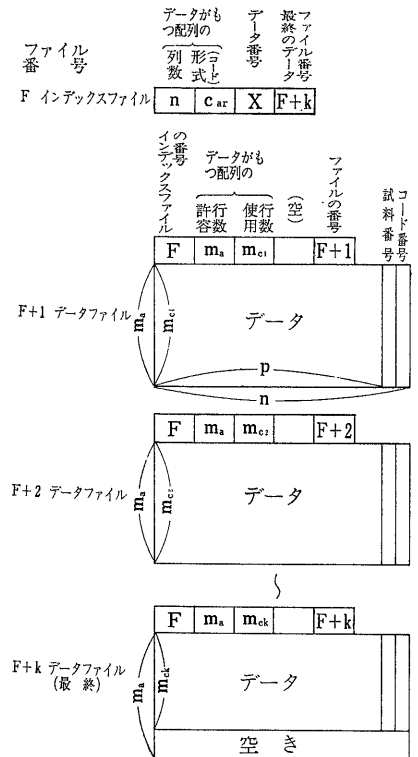
磁気テープは一般にはいくつかのファイルに区切って使われる。さて データの量がファイルの数個分にまたがるようなときには どのように設計しておけば能率がよくしかも安全であろうか。筆者が現在用いている方式はつぎのようなものである。

ファイルの構成と内容を第1図に示す。データを収めるファイルをデータファイルと呼ぶ事にする。データは一般に配列をもっているの で 考えやすいようにファイルを数表の形で表現しよう。ひとつのファイルに入れられるデータの配列は第1図に示すように最大 m_a 行 n 列の規模をもつものとする。そのファイルの大きさが その計算機のメモリー容量によってきめられることは申すまでもない。

データ中の1試料に属する要素の数 たとえば化学分析値の場合は成分数が この方式では配列の列数 P に対応する。各試料には試料番号と分類用のコード番号 (1.3 で述べる) が付けられるのでファイルの列数 n は $P+2$ 列となる。これらはデータによって定まった値なので メモリー容量すなわちデータレジスタ数を列数 n で割った値 (の端数を切り捨てたもの) が そのデータに対するファイルの許容行数 m_a となる。これはそのファイルに収容可能な試料数と同じ意味である。

データの量に応じてデータファイルの数をあらかじめ算出して ファイルを作っておかねばならない。このとき 試料数が配列の行数に対応する事を念頭におけば必要なファイル数はすぐ計算できる。なお 筆者の方式では データファイルの番号はつねに一連でなければならない。データをこれらのファイルに順次入力するためのプログラムも作られているが それについては省略する。データが入力されたあとのファイルの状況はデータを満載したファイルはその使用行数 m_c が許容行数 m_a と同じになるし ファイルに空きがある (一般には最終のファイル) 場合は 当然 $m_a > m_c$ となる。

これら一連のデータファイルの前に見出用のファイルを付けておく。これをインデックスファイルと呼ぶ事



第1図 磁気テープのファイルの構成

にする。このファイルには データ全体についての事項を記入しておく。すなわち データの列数・配列の形式を示すコード番号・データ全体の番号(データ名)・最終のデータファイル番号などである。これとは別に各データファイルにも おもにそのファイルに関する事項を記入する部分を作っておく。すなわちインデックスファイルの番号・データの許容行数と使用行数・そのデータファイルの番号などである。

計算を実行するときは まず必ずインデックスファイルの内容を入力し つぎにデータファイルの内容を必要に応じて順に入力する。そして現在計算機に入力されているデータがどのファイルのものか その配列規模がどの位かなどについては 上に述べた項目をチェックする行程をプログラムに付けておけば ファイルの取り違いや配列の計算処理の誤りを防ぐ事ができる。

1.1.2. テープ 走行 時間 の 節約

磁気テープの欠点のひとつは その走行時間がほかの計算処理時間に比べて圧倒的に長いという点である。筆者の経験では YHP-20 計算機に1000個の数値をテープから入力する場合 1回について約25秒かかる(ただし同社の新型機ではこれより数倍早いという話だが)。計算ループのかけ方によっては ひとつの計算の中で同じファイルの内容が繰り返えし入力されるような事が起きるが このような場合は 計算時間の大半がテープ走行時間に食われ しかもそれが実用にならないほどの長時間に及ぶ事さえまれでない。

そこで 計算の中間結果をメモリーに足し込むなど 適宜必要な処置をして できるだけテープの走行回数を減らすように努める必要がある。その具体例は 相関係数の計算(2.2)とヒストグラム(3.1.1)のところで述べる。

1.2. 数値が欠けているデータの取り扱い

実際の計算では データの中の数値がところどころ欠けたものを やむを得ず処理せねばならない事が起きる。たとえば化学分析データで いくつかの試料の中である成分の分析値がないという場合がこれに当たる。このようなときは分析値のあるものだけについて ともかく処理しようとするのが普通であろう。

この場合 数値がないという事と数値が0という事は当然区別されねばならない。そこでこの数値がない事を意味する“数値”をどう決めるかであるが これは一般には その計算に使われる数値の定義域の外にある数値を使えばよい。筆者はプログラムによってまちまちだが-90000以下とか -10^{99} 以下の数値が入力されたら

“数値なし”とする判断を付けている。“数値なし”の印字は“……”とするようにもしている。

化学分析値などマイナス値がないものは マイナスの値をこの意味に使える。これを応用して $-n$ を入力すると“以下 n 個空欄”の意味に使う事もできる。そのようにしたプログラムも作ったが 結構便利である。

1.3. コード番号による試料の分類と選択

磁気テープに収められているデータの一部についてだけ計算をしたい場合 必要なものを選び出すために 数ファイルにまたがる量のデータを編集しなおすのは 大変やっかいである。そこで 必要な試料だけを選び出す簡単な方法を考えてみよう。

筆者のやり方では 試料番号0の試料は計算から除くという判断が各プログラムに付けてある。各ファイルにあるデータの数値を変更するのは訂正用のプログラムを使うと容易なので 番号を0に改めればよい。

データが分類可能なときは 入力の際に使用者が任意に分類用のコード番号を設計して 各試料ごとに付けておき 計算の際にこのコード番号に合ったものだけを処理するという方法ができる。これについて説明しよう。

YHP-20 計算機の場合は 1個の数値が10けた以内に定められているので 筆者は各試料のコード番号を10けた以内で表現する事にした。この10けたの数値を3個に区切り 3種類のサブコードに分ける。それらはつぎのとおりである。

試料のコード番号→ $\frac{9}{\text{第1コード (6けた)}} \frac{8}{\text{第2コード (2けた)}} \frac{7}{\text{第3コード (2けた)}} \frac{6}{\text{第2コード (2けた)}} \frac{5}{\text{第3コード (2けた)}} \frac{4}{\text{第3コード (2けた)}} \frac{3}{\text{第2コード (2けた)}} \frac{2}{\text{第3コード (2けた)}} \frac{1}{\text{第3コード (2けた)}} \frac{0}{\text{第3コード (2けた)}}$

使用者は この3種のコードに任意に意味付けをしてよい。たとえば岩石試料なら 第1コードを産出地域(999999カ所に区分可能) 第2コードを時代(層準)(99に区分可能) 第3コードを 岩質(同じく) という具合に定義して分類したらよいだろう。

計算の際には 各サブコードの何番を選ぶかを指定する。もしそのサブコードにとくに着目しなくてよい場合は 0を指定する。たとえば第1コードから第3コードへ順に 21, 0, 14 を指定すると コード番号 $21 \times 10^2 + 0 \times 10 + 14$ をもつ試料が選ばれる(\times はこのけたの数字に無関係の意味)。データ全体を取り扱いたいときは 各コードとも0を指定すればよい。

つぎに 試料を選択する行程のために 筆者が考えた方法を示そう。試料のサブコードを d_1, d_2, d_3 使用者が指定したサブコードを S_1, S_2, S_3 とする。このとき 判断の行程はつぎのとおりである。

1. もし $d_1 \neq 0$ で $d_1 \neq S_1$ ならば 計算から除く
2. もし $d_2 \neq 0$ で $d_2 \neq S_2$ ならば 計算から除く
3. もし $d_3 \neq 0$ で $d_3 \neq S_3$ ならば 計算から除く

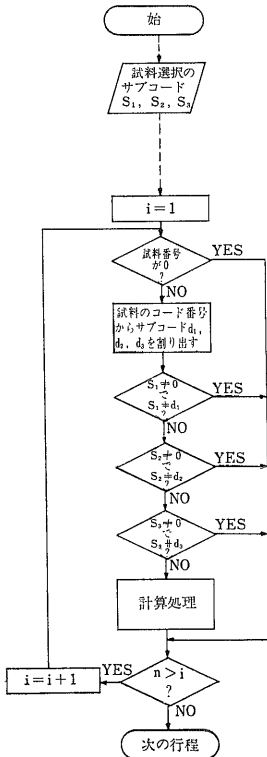
以上3段の判定で計算から除外されなかったものが計算のために選択された試料である。行程としては単純明快である。

第1コードから第3コードまでを上のようなケタ数に区切ったのは 使いやすさを考えての事であり各サブコードとも1ケタで良いのなら YHP-20の場合1レジスタあたり最大10種のサブコードを作ることができる。

使用者はこれを任意に組み合わせてよいので サブコードのうちどれとどれに着目して組み合わせるかという数だけを考えてもサブコードの種類が増すとその数は飛躍的にふえる。

しかし上に述べた方法によれば判断行程の数はサブコードの数にしかならない。行程をひとつにしてループ計算にすればプログラムも節約できる。したがって計算時間さえ問題にしなければ このようなやり方で十分多種の条件を重複させる事ができる。試料選択の行程を第2図に示す。

上の行程を眺めていると計算機がなるべく仕事をすまいと努めているように見えるので このやり方を“なま



第2図 試料を選択する行程

けもの法”とでも名付けたらどうだろう。

2 相関係数と回帰直線

相関係数を求めるプログラムは プログラミングの技術としては特筆すべきものはないのだが さまざまな目的で利用されたという実績を考えると ご紹介をする事にした。さらに その計算結果を評価する際に注意すべき点についても付け加えておいた。

筆者は統計学に対してとくに深い知識があるわけではないので 相関係数についての詳しい議論はさておき とくにプログラムと関係する要点だけを述べてみよう。

2.1. 相関係数とは

ひと口でいえば 相関係数はふたつの変量 X, Y の関係の強さを示すもので X, Y を同時に観測して得られた n 個の数値の組

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

を XY 平面上にプロットして得られた相関図の直線性の度合を表す値である。数式の上では X, Y の共分散を X と Y のそれぞれの標準偏差の積で割った値である。すなわち相関係数 r は

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \dots\dots\dots(1)$$

である。ここに \bar{x} は (x_1, x_2, \dots, x_n) の また \bar{y} は (y_1, y_2, \dots, y_n) のそれぞれ平均値である。

r は $0 \leq |r| \leq 1$ の範囲にあり $|r|$ が1に近いほど X, Y は強い相関関係にあり また $r > 0$ の場合は正(順)相関 $r < 0$ の場合は負(逆)相関という。

変量 X, Y に属する数値の組に直線的関係が認められる場合は 最小二乗法により相関図に回帰直線を引く事ができる。回帰直線は 任意の直線 AB を考え AB に点 $P_1(x_1, y_2), P_2(x_2, y_2), \dots, P_n(x_n, y_n)$ を通り Y 軸に平行な直線を引き AB との交点をそれぞれ $P_1'(x_1, y_1'), P_2'(x_2, y_2'), \dots, P_n'(x_n, y_n')$ とするとき y_i から y_i' までの長さの二乗の和 すなわち

$$L = \sum_{i=1}^n (y_i - y_i')^2 \dots\dots\dots(2)$$

を最小にするように引かされた直線の事である。

回帰直線の式を

$$y = ax + b \dots\dots\dots(3)$$

とする。このとき勾配 a は

$$a = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \dots\dots\dots(4)$$

切片 b はこの直線が点 (\bar{x}, \bar{y}) を通ることから

$$b = \bar{y} - a\bar{x} \dots\dots\dots(5)$$

としてそれぞれ求められる。これらの式をみちびく計算の詳細については 高橋ほか (1975) などにわかりやすく説明されているので ここでは省略する。

2.2. 相関係数を求めるプログラム

まず r を求める計算は (1)式を展開し

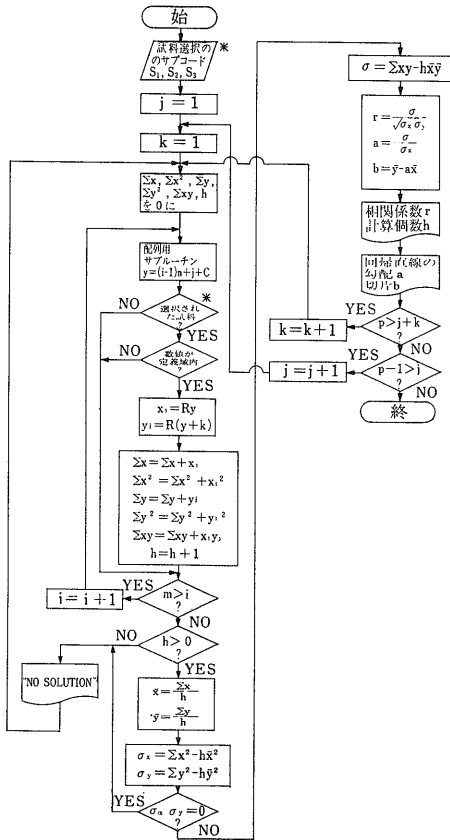
$$r = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2\right)\left(\sum_{i=1}^n y_i^2 - n\bar{y}^2\right)}} \dots\dots\dots(6)$$

という形にして行なう。 r の値を求めるには(6)式を構成する

$$n, \sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2, \sum_{i=1}^n y_i, \sum_{i=1}^n y_i^2, \sum_{i=1}^n x_i y_i \dots\dots\dots(7)$$

の各項の値をあらかじめ求める必要がある。

a についても同様に (4)式を展開し



第3図 相関係数の計算行程

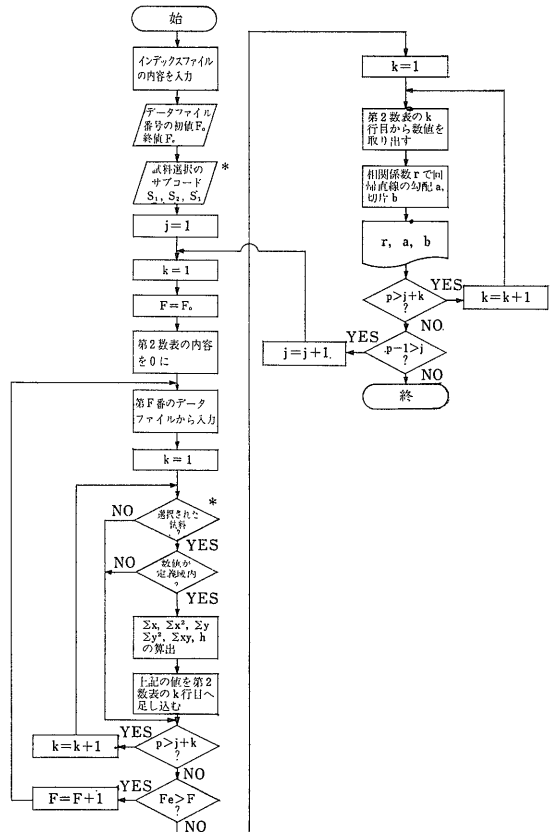
図には記さなかったが 回帰直線の a, b の値は $r \geq 0.9$ であって使用者が必要とするときに限って出力する。 m は配列の行数 p は変量の(列)数 n は配列の列数 ($p+2=n$ 試料番号とコード番号が加わるため)。 j は変量 c_j k はこれと組み合わせられる変量 c_{j+k} のそれぞれ列番号用の計数器。配列用サブルーチンについては吉井(1977)を見よ。 * 第2図参照

$$a = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \dots\dots\dots(8)$$

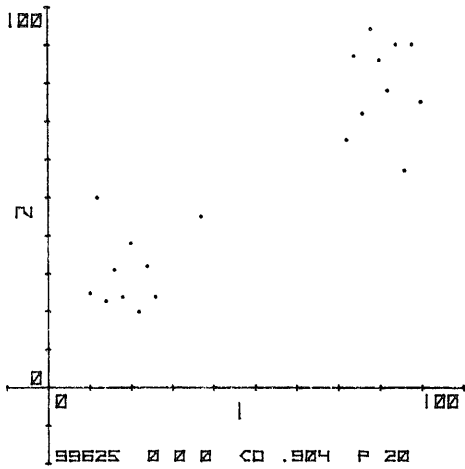
の式から求める。(8)式を構成する項は(7)式に含まれているので a は容易に算出できる。

筆者はプログラムでは n 行 p 列の配列をもつデータに対して その列要素を変量 C_1, C_2, \dots, C_p として そのすべての組について総当たりで相関係数を計算するようにしてある。すなわちその組み合わせと順序はまず C_1 対 C_2, C_3, \dots, C_p つぎに C_2 対 C_3, \dots, C_p そして最後は C_{p-1} 対 C_p である。その行程を第3図に示す。

データが数ファイルに及ぶときは テープ走行時間の節約のために C_j 対 $C_{j+k}, C_{j+k+1}, \dots, C_p$ についての計算を各ファイルごとに行ない。(7)式の各数値をメモリーの第2数表に記入しあとまでめて値を取り出すようにする。第2数表は k 行 6列の配列をもち 各列には(7)式の各数値が対応する。各行要素に対しては C_j 対 $C_{j+k}, C_{j+k+1}, \dots, C_p$ の計算結果が第 k 行に記



第4図 相関係数の計算で磁気テープを使う場合の計算行程
第2数表は k 行 6列の配列をもち その操作は k ループで行なう。合計3ヶ所に k ループが配置されている。 f はファイル番号用の計数器。計算の主要部については第3図を見よ。 * 第2図参照



第5図 相関係数を求め場合に問題となる例
 点の塊が2ヶ所に別れて分布するとき 全体としては $r=0.90$ となるが 個々の塊では左下のものが $r=0.21$ 右上のものが $r=0.02$ である。

入される。これを各ファイルについて繰り返して 足し込むと 最後のファイルの計算が終わったところで 第2数表には C_j 対 $C_{j+k}, C_{j+k+1}, \dots, C_p$ の相関係数を求めるのに必要な数値がそろう事になる。これらの組に対して一挙に相関係数を計算し結果を出力する。この行程を第4図に示す。

この方法によってもテープは $p-1$ 往復する事になり 能率の上ではまだ問題だが それでも積和計算のループの中に テープ走行を組み込むと実に $p(p-1)/2$ 回も往復する事になってしまう。

2.3. いくつかの問題点

実際の計算で考えておかねばならない点について つぎにいくつか述べる。数学的には $|r|$ が1に近い値ならば そのふたつの変数 X, Y には直線的関係があるといえる。だが 実際のデータでどのような場合にもこの結論を下してよいものか 少し検討しよう。

いまここに $r=0.90$ のデータがあり その相関図が第5図に示されるようなものであったとする。この図をみると点は離れた2ヶ所に塊を作って分布している。このようなパターンを示すときは必ず $|r|$ の値が1に近くなる。なぜならば計算上は点が2個しかないデータ ($|r|=1$) に近い結果となるからである。もし2ヶ所の塊が何かの条件(地学的データなら地域・層準・岩質など)の差を反映したためにできたものだとしたら 統計上は別の標本として扱うべきものかもわからない。個々の塊について計算すると 第5図の左下の塊りでは $r=0.21$ 右上の塊では $r=0.02$ となり 相関関係は認められない。

ここに示した例は データの分類を変えただけで相関係数が大きく変化する場合のある事を示している。このような例はしばしばみられるので注意を要する。

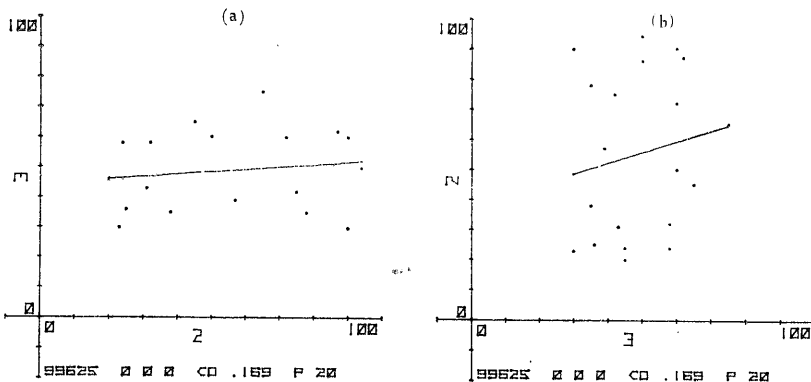
回帰直線を引く際にも問題がある。第6図にその例を示す。第6図(a)では一回帰直線が順当に引かれているように見える。しかし第6図(b)ではこの直線がまったく見当違いの方向に延びている。このふたつの図は実は同じデータをX軸とY軸を入れ替えて計算しただけなのである。ではなぜ図の向きによってこのような大きな違いが出るのだろうか。

この場合 r が0.17と大変低い値にある点が重要なヒントになる。 r を求める(1)式と a を求める(4)式は形がよく似ているので (4)式の分母分子に

$$\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \dots\dots\dots (9)$$

を掛けて整理すると

$$a = r \cdot \frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \dots\dots\dots (10)$$



第6図 回帰直線の問題を示す例

- (a) この図では相関図に一見正しく回帰直線が引かれているように見える。
- (b) この図は(a)図と同じデータをX軸とY軸を入れ替えてプロットしただけである。直線はここでは明らかに違った方向に引かれており 回帰直線とは認められない

となる。これで回帰直線の勾配 a が相関係数 r の関数になっている事がわかる。したがって $|r|$ が小さくなると回帰直線は点 (\bar{x}, \bar{y}) を中心に勾配がゆるやかになる方向に回転し $r=0$ では X 軸と平行になってしまうのである。

すなわち回帰直線は $|r|$ が 1 に近い場合以外はうかつに計算から求めて引いてはならないのである。筆者のプログラムでも $|r| \geq 0.9$ であって使用者が必要とする場合のほかは回帰直線の定数 a, b を算出しないようにしている(第3図と第4図ではこの判定を省略した)。

r の値から相関関係を論じ また回帰直線を引く場合は いずれもふたつの変量の間に直線的な関係があるという事を前提としている。だから理論式が一次式で表わされるときに観測値からその実験式をみちびくというような目的には 相関係数や回帰直線の考えが適合する。

しかしこれを化学成分間の相関関係の推定に応用するような場合は 上にあげたような種々の問題が起きる可能性がある。したがって数値だけから安易に結論を急ぐのは大変危険である。

これらの検討は必ず相関図を作って行なうべきである。

3. 統計 図

データが多いときは数表のままでは検討するよりも 図に表わして眺める方が そのデータのもつ性質をよく理解できる事は申すまでもない。しかし作図は手間のかかる仕事だし 人によって得手不得手もあり 一般にわずらわしい。そこで計算機にプロッタを接続して 必要な図を使用者の思うままに描かすと 能率的だし正確でもある。

では代表的な3種類の統計図すなわちヒストグラム・相関図・三角図を描くプログラムについて 目下筆者が実用化しているものをご紹介します。

これらのプログラムでのデータの取り扱い方は共通である。まずこれについて簡単に述べておく。

計算機へのデータの入りは各プログラムの中で磁気テープから行なわれるか または別に用意されたプログラムによってキーボードから行なわれる。計算処理はメモリー内の数表(吉井 1977)の列要素に対して行なわれる。その配列の列番号をここでは成分番号と呼ぶ事にする(化学分析値の各成分に対応するので)。

図には成分番号が意味する成分名を書く事もできる。これにはそれらに対応させるサブルーチンを作って 主プログラムに付け加えておく。同様に図の標題を書く事もできる。これらは使用者が適宜作るものとする。その詳細は省略する。

3.1. ヒストグラム

ひとつの変量 X に属する n 個の数値 x_1, x_2, \dots, x_n をいくつかの階級(区間)に分けて 各階級に属する数値の度数(個数)を求め 棒グラフにしたものをヒストグラム(度数柱状図)という。 XY 座標の X 軸に階級を Y 軸に度数をそれぞれ取るものとする。度数百分率(=度数 $\times 100/n$)による表示もしばしば用いられる。階級の幅(X 軸のきざみ)は 一般には階級ごとに異なっているがよいのだが ここでは各階級が同じ幅をもつ場合に限る事にしよう。

3.1.1. プログラム

まず作図条件の入力が必要である。それにはつぎのようなものがある。

1. X 軸 側
 - a 図示範囲の最小値 X_{\min} と最大値 X_{\max}
 - b 階級の幅 d (軸の目盛を兼ねる)
2. Y 軸 側
 - a 度数の最大値 f_{\max}
 - b 軸の目盛

これらと図との関係を第7図に示す。

ヒストグラムのプログラムの行程を第8図に示す。

ヒストグラムを描くには まず各階級ごとの度数を算出する必要がある。一般には各階級の境界値を第7図のように a_0, a_1, \dots, a_k とするとき x_1, x_2, \dots, x_n の中から $a_{i+1} > x_j \geq a_i$ となる数値の個数(これが度数)を求める。

プログラムでは 各階級をはさむふたつの境界値のうち値の大きい方を a_i 小さい方を a_s として もし

$$a_i > x_j \geq a_s$$

ならば度数用の計数器 f_i を1だけ足し上げるという操作を x_1, x_2, \dots, x_n に対して繰り返す。階級の移動は初値をそれぞれ

$$a_s = X_{\min} \quad a_i = X_{\min} + d$$

として 以下順次

$$a_s = a_i$$

$$a_i = a_s + d$$

の操作で行なう。これらによって度数 f_1, f_2, \dots, f_k の算出が行なわれる。

つぎにヒストグラムを描く行程について述べる。いま座標内の任意の点 (x_1, y_1) から別の点 (x_2, y_2) まで直線を引く命令を

$$PLT(x_1, y_1) \rightarrow (x_2, y_2)$$

と表わす事にすると $a_s = a_i \quad a_i = a_{i+1}$ の階級では

$PLT(a_i, f_i) \rightarrow (a_i, f_i)$

という命令になる。プロッタのペンを下げたままにしてこの操作を繰り返せばよく、各階級の境界値のため線も自動的に引かれる。ヒストグラムの右端にあるたて線はこれらのループから脱け出したところで ($a_i = a_k$ のとき)

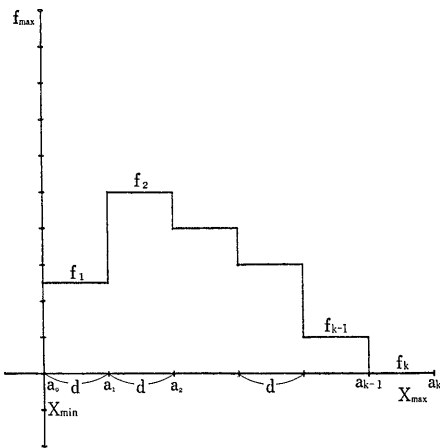
$PLT(a_i, f_k) \rightarrow (a_i, 0)$

の命令で引かせる。

さてヒストグラムは原理的には度数 f_i を算出するごとにその階級の分を描くようにすればよい。しかしデータが磁気テープの数ファイルに及ぶ場合はこの方法では階級の数だけテープが往復しその走行時間が長過ぎて実用にならないおそれがある。そこで各ファイルのデータごとに度数 f_1, f_2, \dots, f_k を求めてメモリーに足し込みあとで全体について一挙にヒストグラムを描くという方法をとるとよい。

この方法を積極的に利用するとヒストグラムの重ね描きができる。すなわち作図条件が同じの複数のデータについて各度数を積算し前のヒストグラムの上に積み重ねてゆくのである。これとコード番号による試料の選択とを組み合わせるとヒストグラムの合併などが容易となりデータの解析に役立つ。この重ね描きの例を第9図に示す。

描かれた図の下にはデータの番号・選択コード・作図条件・累積度数など必要事項を書き添える。それらを同時にプリントアウトしておく図の検討や再生に便利である。その例を第10図に示す。



第7図 ヒストグラム

a_0, a_1, \dots, a_k 階級の境界値 d 階級の幅 f_1, f_2, \dots, f_k 度数

第8図

ヒストグラムを描くプログラムの行程

*1. 磁気テープに入っているデータの入力

*2. コード番号による試料の選択

*3. 使用者が別途用意するサブルーチン。成分名用のサブルーチンを付けたいときは図に成分番号が書かれる。

3.1.2. 階級の幅と図形の変化

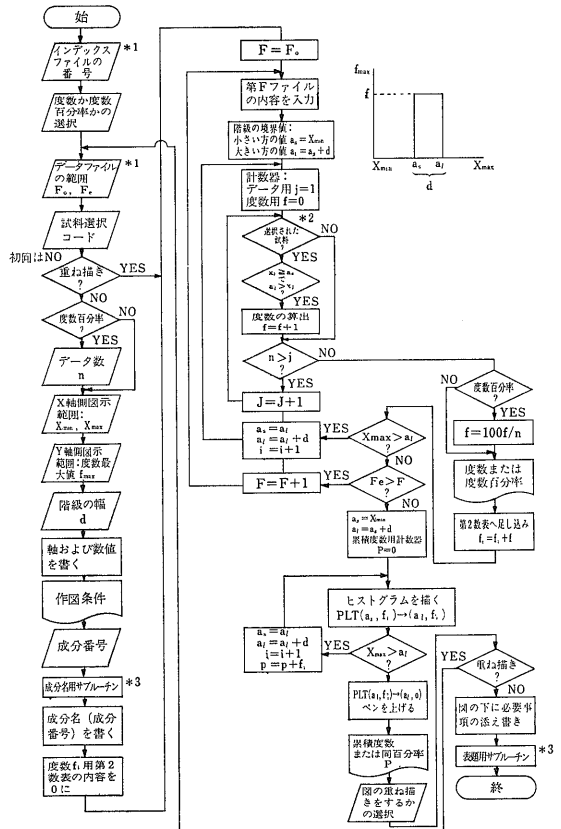
ヒストグラムの形は階級の幅 d の大きさによってかなり違ってくる。同じデータに対して d の値を 10, 15, 20, 30 と変化させて得られた図形を第11図に示す。 d の値は大き過ぎていけないのは申すまでもないが小さ過ぎててもでこぼこが目立って全体の状況がわからなくなる。適当な値がどの辺にあるかはデータによって異なるので試行錯誤せねばならない。これを手作業で行なうのは大変な手間であるが計算機に順次異なった d の値を入れては図を描かせもつとも説得力のある図形を選択するのは容易である。その上度数の算出は正確だからこの辺は計算機の威力である。

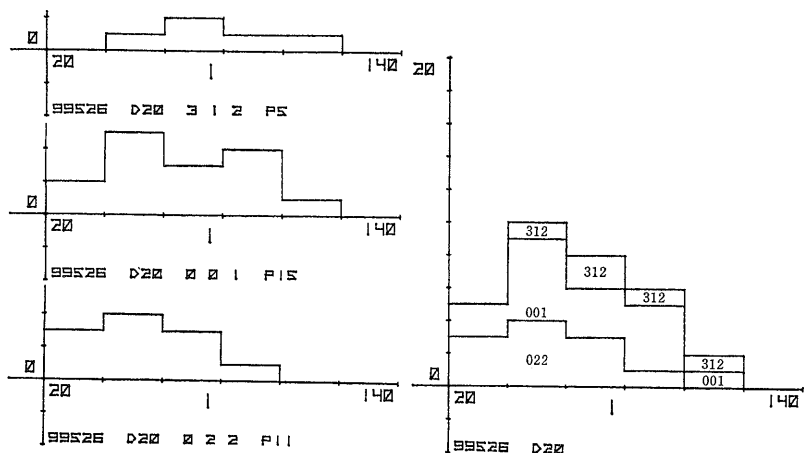
3.2. 相 関 図

2 種の変量 $X(x_1, x_2, \dots, x_n)$ と $Y(y_1, y_2, \dots, y_n)$ の対応する数値 (x_i, y_i) の組の相関関係を調べるために XY 座標に表わしたものを相関図という。

3.2.1. プログラム

相関図の行程を第12図に示す。





第9図
ヒストグラムの重ね描き
ひとつのデータをコード番号によって3つに分けた例を図の左側に示す。図の下の数字は左からデータ番号・階級の幅・選択コード(3字)・累積度数を示すヒストグラムを 選択コード022 001 312の順に下から積み重ねてできた図を右側に示す。

作図条件としてはX軸側Y軸側とも つぎのものを入力する。

1. 図示範囲

- a 最小値 X_{min}, Y_{min}
- b 最大値 X_{max}, Y_{max}

2. 軸の目盛

データの処理でほかのプログラムと共通な点は省略する。 関連図の例を第13図に示す。

関連図の図示範囲は 使用者が任意に指定できるために データの一部が図示範囲外に飛び出す場合もあり得る。 これを使用者が承知している場合はよいのだがもし図外となる数値が出たのを知らないでいると問題がおきる。 そこで図外になりプロットからはずされた値は別にプリントアウトしておく。 こうすれば 使用者

への警告にもなり 図示範囲の検討をする上での参考にもなる。 このプリントアウトには作図条件そのほかの必要事項も印刷するその例を第13図(右)に示す。

数値の組 (x_i, y_i) をプロットする行程は
 $PLT(x_i, y_i);$ ペンを上げる。

という命令で これを繰り返せばよい。 上の命令のあとに停止 (STOP) 命令を付けて 1点打つごとに動作が停止し プログラムを進めるキーを押すと次の1点を打ちまた動作が止まるという方法もできるようにしておく。 これを利用すると 打たれた点に使用者が番号そのほか必要なメモを そのつど書き添える事ができて便利である。 図の重ね描きもできる。 関連図の場合は重ね描きと言っても点の追加に過ぎないので プログラムもヒストグラムに比べてずっと簡単である。

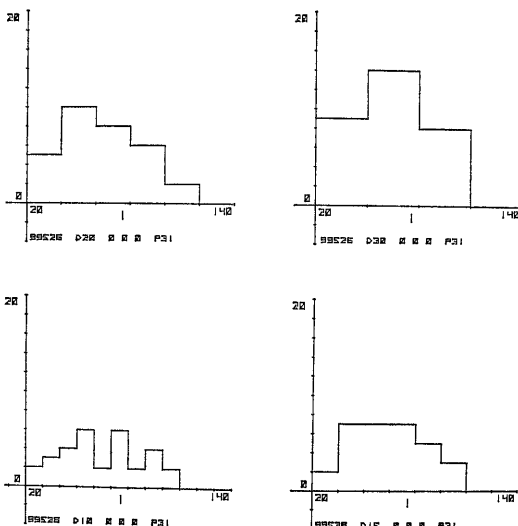
また ペンを上げる命令を略して ペンをおろしたま

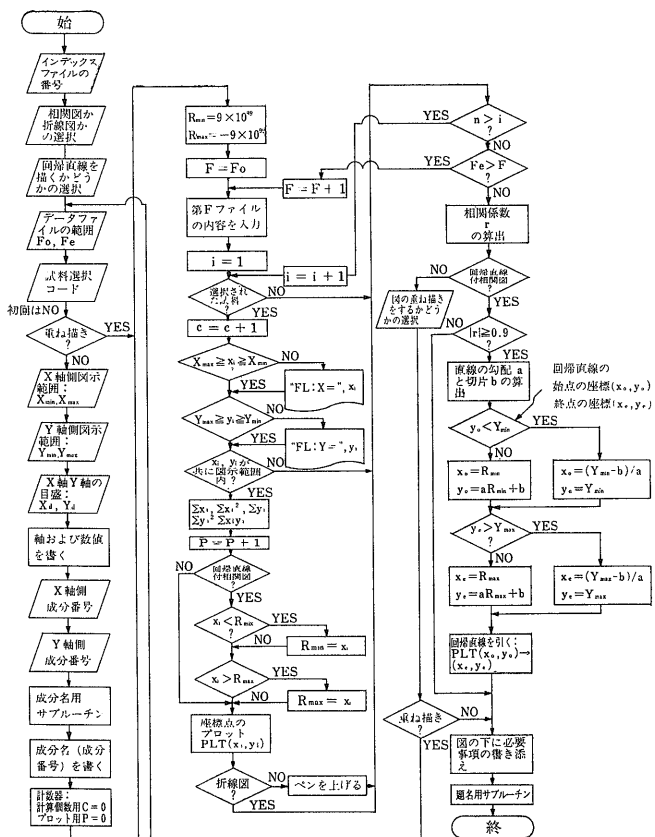
```

JOB=
FILE:FRM=,TQ=
SELECT=
X:FRM=
TQ:X=,Y=
DIV:X=,Y=
CMP=
2,000
3,000
4,000
6,000
2,000
6,000
2,000
4,000
2,000
0,000
0,000
0,000
T=
31,000
    
```

第10図
ヒストグラムの記録用プリントアウト
上からデータ番号(作業番号)・磁気テープのファイルの区間・選択コード・X軸側の図示範囲の最小値・両軸の最大値・軸のきざみ(X軸側は階級の幅を意味する)。 入力データの成分番号・各階級の度数およびその合計(第11図左上のヒストグラムの記録)。

第11図
階級の幅の違いによるヒストグラムの変化
同じデータでも階級の幅によってパターンが異なる。
階級の幅は 左下10 右下15 左上20 右上30



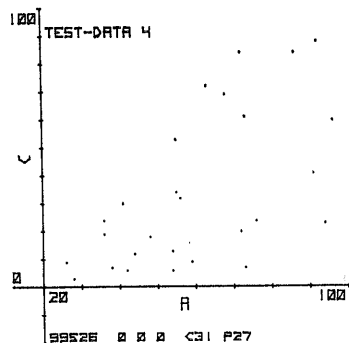


第12図 相関図(回帰直線付き)のプログラムの行程

ま上の操作をすると折線図が描ける。なおこの場合一般には片方のデータの組がある規則的な順序に並べられている必要がある。

3.2.2. 回帰直線

相関図に回帰直線を引く事もできる。ただし筆者の



```

DOT
JOB=
FILE:FRM=,TO= 99526
SELECT=
X:FRM=,TO= 28
Y:FRM=,TO= 100
DIV: X=Y= 10
CMP: X=Y= 13
FL:Y= 280.000
FL:Y= 227.000
FL:X= 109.000
FL:X= 109.000
C,P=
31
27
    
```

プログラムでは 相関係数の絶対値が0.9以上で 使用者が希望するとき という条件を付けている。 そのわけは前の章(2.3)で説明した。 回帰直線を引いた例を第14図に示す。

回帰直線の勾配 a と切片 b を求めるには相関係数の計算が必要なので その行程を付け加えておく。 この計算は図示範囲内にある点の数値に対してだけ行なう。

回帰直線の線分は図中の点の分布する範囲にだけ引くように工夫するとよい。 それには図示範囲内にある点の X 座標の最小値と最大値をそれぞれメモリー (R_{min} および R_{max}) に記憶させる。 いま最小値を求める方法を考えてみよう。 初値は $R_{min}=9 \times 10^{99}$ (十分に大きい値) とし x_i が $x_i < R_{min}$ ならば R_{min} へ x_i の値を入れる。 この操作を x_1, x_2, \dots, x_n について行なった結果 R_{min} にはその最小値が入れている。 最大値も同じような手順で R_{max} に入れる事ができる。

回帰直線の式を

$$y = ax + b \quad (a \neq 0) \dots\dots\dots (1)$$

とする。 いまこの線分の始点 (x_0, y_0) について考えてみると その座標の値は

$$x_0 = R_{min} \quad y_0 = aR_{min} + b \dots\dots\dots (1')$$

である。 このうち x_0 値はつねに図示範囲内にあるが y_0 は必ずしもそうとは限らない。

いま仮りに y_0 の値が図の下方へ飛び出す ($y_0 < Y_{min}$) とすれば始点を上の直線上で移動して図の下限の境界線上にもって来る必要がある。 すなわちこの場合の始点は

$$x_0 = \frac{Y_{min} - b}{a} \quad y_0 = Y_{min} \dots\dots\dots (1'')$$

とする。 $y_0 > Y_{max}$ の場合は上の式の Y_{min} を Y_{max} に置き換える。 同様にして終点 (x_e, y_e) についても必要な処理をする。 回帰直線を引く命令は

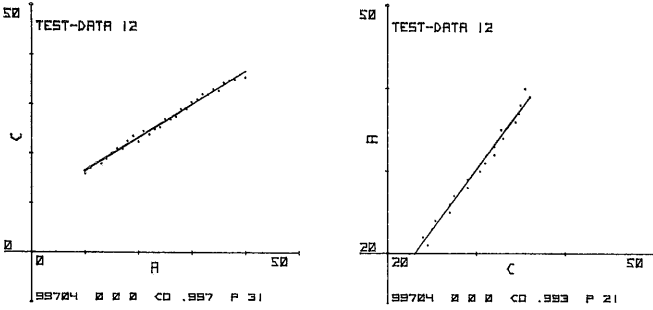
$$PLT(x_0, y_0) \rightarrow (x_e, y_e); \text{ペンを上げる}$$

となる。 始点が図の下方になる例を第14図(右)に示す。

3.3. 三角図

3 種の変量 $A(a_1, a_2, \dots, a_n) B(b_1, b_2, \dots, b_n) C(c_1, c_2, \dots, c_n)$ の対応する数値の組 (a_i, b_i, c_i) の相互の比率を正三角形内の点に対応させて表わしたものを

第13図 相関図とそのプリントアウト
相関図の成分名(AとC)は成分名用サブルーチンによる。 図の下の数値でC23は計算数 P27はプロット数で両者に差があるのは図示範囲から4個の数値が飛び出した事を示す。
プリントアウトの成分番号(CMP)1と3は図のAとCの成分と対応する。 その下に図示範囲から出た数値4個が示され それがどちらの軸で起きたかもわかる。



第14図 回帰直線付き相関図

右の図は左の図の成分を入れ換え 作図条件も変更したもので これにより回帰直線の始点が図外となった。この場合はX軸上に始点を求め直して回帰直線を引いている

が三角図である。地学関係の研究ではしばしば使われてなじみの深い図でもある。

三角図は 正三角形内の任意の点Pから各辺におろした垂線の長さ a, b, c の和がその正三角形の高さ h に等しいという性質を利用したものである。

蛇足ながら これを証明しておこう。証明法はさまざまあるので各自試みられたい。筆者が考えたものをつぎに記す。

まず1辺の長さが2の正三角形 ABC を考える。この内部に任意の点 P をとり P から辺 BC, CA, AB におろした垂線の長さをそれぞれ a, b, c とする。また P と頂点 A, B, C をそれぞれ直線で結ぶ。これらを第15図に示す。このとき 各辺の長さは

$$BC=CA=AB=2$$

だから $\triangle BCP, \triangle CAP, \triangle ABP$ の面積はそれぞれ a, b, c となる。また $\triangle ABC$ の高さは h だから $\triangle ABC$ 面積は h となる。図から明らかに

$$\triangle ABC = \triangle BCP + \triangle CAP + \triangle ABP$$

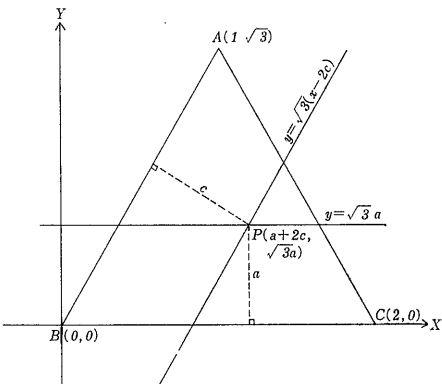
したがって

$$h = a + b + c$$

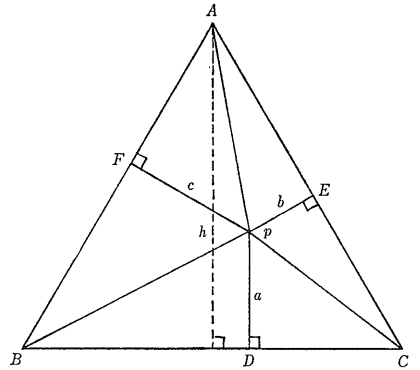
これで証明が終った。

3.3.1. プログラム

三角図では作図条件の指定はなく 3種の成分番号を入力するだけなので 使い方は簡単で プログラムも比較的短かい。ほかのプログラムと共通する点の説明は



第16図 三角図のプログラムでの正三角形の配置とプロットすべき点Pの求め方



第15図 三角図の原理
正三角形では $h = a + b + c$

省略する。

ではXY座標上に正三角形をどのように置き プロットする点をどのように決めたら良いか述べよう。筆者がいくつか検討したところでは この正三角形は頂点のひとつを座標の原点に置き 1辺をX軸上に置いてその長さを2とすると 式が簡単になるようだ。これを第16図に示す。すなわち各頂点の座標は $A(1, \sqrt{3})$ $B(0, 0)$ $C(2, 0)$ とする。

つぎに点Pの座標の求め方について述べる。

数値の組 a_i, b_i, c_i (ただし $a_i \geq 0, b_i \geq 0, c_i \geq 0$) の和

$$S = a_i + b_i + c_i$$

を求め $S > 0$ ならば ($S = 0$ のときはあとで述べる) S に対する a_i, b_i, c_i の割合を

$$a = \frac{a_i}{S}, \quad b = \frac{b_i}{S}, \quad c = \frac{c_i}{S}$$

とする。これらを第16図の正三角形の辺 BC, CA, AB から P まで距離に対応させる。

まず a は BC に平行な直線を表わし それが点 $A(1, \sqrt{3})$ を通るから その方程式は

$$y = \sqrt{3}a$$

と書ける。同様に c は AB に平行な直線で $C(2, 0)$ を通るので その方程式は

$$y = \sqrt{3}(x - 2c)$$

となる。この2直線の交点が P であり その座標は $P(a + 2c, \sqrt{3}a)$

として求められる。この点のプロット命令は

$$PLT(a + 2c, \sqrt{3}a); \text{ペンを上げる}$$

である。

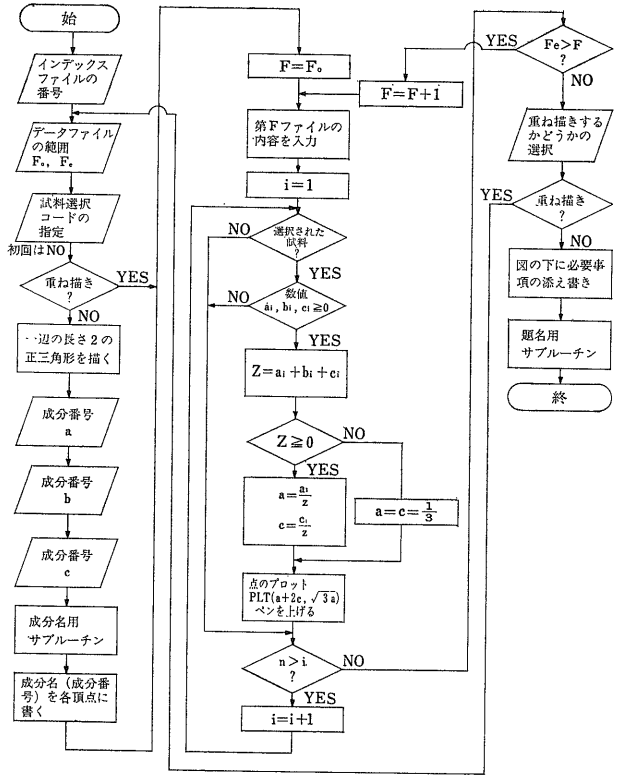
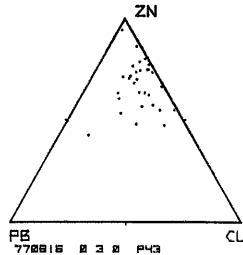
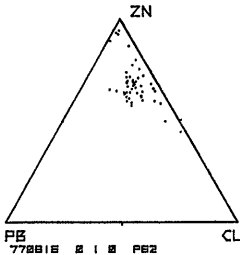
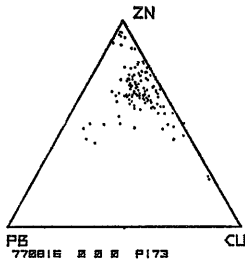
なお $S = 0$ のときは 分母が0となるため上の分数計算ができない。しかし $a = b = c$ の関係には変わりないのでこの場合は便宜的に

$$a = c = \frac{1}{3}$$

として取扱う。 $S = 0$ の場合の対策をしておかないと 不条理演算となって計算が停止してしまうのは申すまでもない。 三角図のプログラムの行程を第17図に示す。 また三角図の例を第18図に示す。

参考文献

- 高橋宏一・藤本和昌・平野勝臣 (1975) : 統計学要論 206 p. 共立出版
 吉井守正 (1977) : 卓上型電子計算機によるいくつかの計算例 その1 配列をもつデータの処理。 地質ニュース no. 275 p. 16-19



第17図 三角図のプログラムの行程
成分番号 a, b, c の順で正三角形の上の頂点から左回りに 各成分の位置が定められる

第18図 三角図の例
ある岩石の化学分析値をもとに描いた例。 左上が全データ (選択コード000) 左下(010)と右下 (030)はそれを岩質ごとに表わしたものである。 コード番号を活用するとデータの解析が容易になる

日曜の地学 4

東京の地質をめぐって

本書は 東京という郊外の清純な自然を楽しむ喜びと 都会の地下に眠る自然発達史が綴るロマンとを結びつけて 改めて東京の自然環境を見直すのに役立つ目的で 編集されている (同書のまえがきより)。

「東京の地質」というと まずはコンクリートで固められた大都会東京に 地質の見られるところがあるのかといふが 少したってから奥多摩も伊豆大島も東京都だったと考えなおしそこへ案内されるかと思うかもしれない。ところが本書の第1章では 私達が毎日通勤している電車や地下鉄に乗せられることになる。コンクリートを選り抜いて通らなかった編集者の苦心と意気込みに感心させられる。

本書は21人の執筆者 (中・高校 大学教師や地質調査所・文化庁などの職員) により書かれた 23のそれぞれ独立した章から成る。日曜の地学の名のとおり 一部を除き 模式地へ日

帰り程度で案内する形式をとっている。各章ごとに詳しいコース案内図がついているので 地形図は必ずしも持参しなくても 迷うことはない。1~11章は都心とその周辺の第四紀層関係 12~16章は関東山地の基盤岩類 17~20章は伊豆・小笠原諸島めぐり 21~24章は 化石・地形・地史等の総括的な解説からなる。これとは別に 地質調査の心得・公害問題・動植物の話など12項目の記事が挿入されている。説明は平易だが 記述は最新の知識を盛り込み厳密である。

文字通り日曜のある日 本書をたずさえて知的レジャーを楽しむもよし また気軽に通読していても東京の自然環境について教えられるところが多いであろう。(垣見俊弘)

書名: 日曜の地学 4 東京の地質をめぐって
 編著: 大森昌衛
 B 6判 184ページ 980円
 発行所: 築地書館
 ①04 東京都中央区築地2-10-12
 ☎ 03-542-3731 振替東京1-19057