

統計十話

第5話 測定誤差・測定データの変動の評価なくして 統計的分析の意味はない——その2

林 知己夫

文部省統計数理研究所

第5話で、測定データ変動を無視するこわさを実例をもって示した。これをどう処理して正しい結論を得るかをここで示してみよう。

誤差の出方の説明

なお、これまで、あるいはこれから変動といった誤差といったしているが、分析のうえでは同じものとみなしていただきたい。誤差とは、われわれの情報を覆い隠す雑音という意味であるから、生体変動による変動であっても、平均値としての情報を覆い隠す誤差とみなすわけである。いずれにしても真の値（あるいは平均値）のまわりの変動を誤差という表現で述べてみる。したがって、測定者の読みの誤差もあるし、測定条件や生体の変動によるものも含まれているわけである。同じ時期に2回以上測定して標識づけたとき同一の値が示されない。すなわちデータに変動があるとき、これを“誤差があるため”と表現して取り扱う。誤差をこのように値の変動を含めた広い意味に用いることにする。

(1) 誤差の出方

前述のように全体の集団の分散は441程度であるが、これには測定誤差が含まれている。測定誤差は、変動の分散を前に述べた議論からよくコントロールされた条件下で一応50とみなしておき、かつ見通しをよくするため各人の変動の分散も一応同一とみなしておけば、測定誤差のないときの分散 $\equiv \sigma^2$ は、

$$\equiv \sigma^2 = 391$$

となる。前に述べた誤差分散と $\equiv \sigma^2$ との比は、0.13程度でかなり大きい。さて、こうした変動は一応ガウス分布とみなしておく。つまり平均0、分散 s^2 のガウス分布とみなしておく。 $s^2=50$ としておくわけである。この確率論的仮定も、データからみて、妥当な仮定と考えられる。さて、こうなると2回測定を行えば図のように、測定の誤差のないとき (x_0, y_0) 、ただし $x_0=y_0$ であるものが測定されたとき、それらの値はある確率をもって円形の範囲にばらつくことが考えられる。

このような (x, y) が観測されるはずである。このような人々が多くあれば、 $x_0=y_0$ であっても 45° の上にのらずに散らばってくることになる。T年およびT+1年において、測定誤差がなければ、等しい値をもつ、すなわち上の表現に従えば、各人について $x_0=y_0$ であるとしておこう。これは、前に示したT年、T+1年の平均、分散、分布の形の恒常性から一応認められよう。

(2) 測定誤差のない場合の測定値

さて、測定誤差のないときの血圧値の分布も簡単のため、ガウス分布としておく。これは、さきのデータからみるとやや無理な仮定であるようにみえる——実際は分布が少し流れ、ピアソンIII型の分布である——が、計算が容易で、見通しをよくするためひとまずこうしておく。ピアソンIII型にしても以下の計

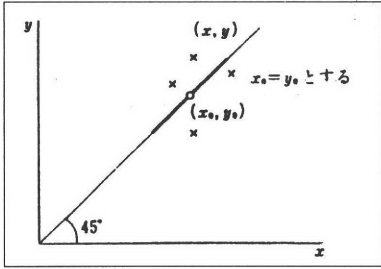


図1 誤差の出方

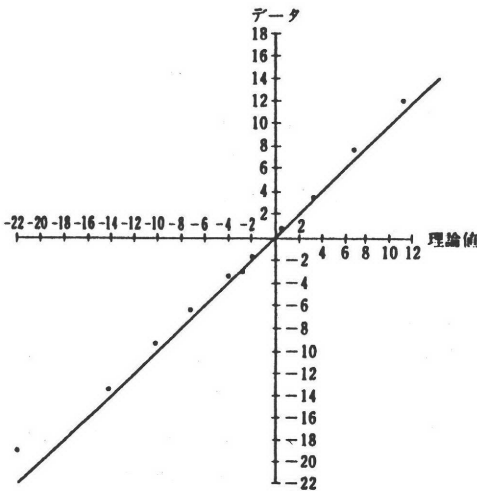


図2 理論値とデータ

算は可能であるが、見通しのよい形で書けな
いだけである。これを

$$\frac{1}{\sqrt{2\pi}\sigma} e^{-(x-M)^2/2\sigma^2} dx = (x|M, \sigma^2)$$

としておく。M は平均値、 σ^2 は分散である。
これまでの表現に従えば、本来 x, σ^2 でなく
 x_0, σ_0^2 を用いるべきであるが、煩雑を避ける
ため添字をおとして x としておく。T 年の測
定誤差のない血圧値の分布も、T+1 年のそ
れも同一であるとしておく。したがって、T 年
の M も σ^2 も T+1 の M も σ^2 も同一で
あり、かつ T 年である x を示すものは、T+1
年でもその x をもつとしておく。さて、測定
誤差の分布は平均 0、分散 s^2 のガウス分布
は、

$$\frac{1}{\sqrt{2\pi}s} e^{-\varepsilon^2/2s^2} d$$

であるとしておくので、測定誤差のないとき
血圧値 x を示すものの測定値 z の分布は、

$$\frac{1}{\sqrt{2\pi}s} e^{-(z-x)^2/2s^2} dz$$

となる。 $s^2=g(x)$ と、 x の関数であることも予
想されるが、いまは、一応 s^2 は常数としてお
く。 $s^2=g(x)$ でも計算はめんどうになるが可
能である。計算は、面倒なので省略するが(詳
しくは、林：行動計量学序説。朝倉書店、
1993, p79~81)、こうして一応の説明は可能
になったが、まだ不十分であった。これは、
 s^2 が一定であるとして計算したからである。
また、もとの血圧分布がガウス分布をしてい
ないこともその原因でもある。

より現実に近づける……

血圧値の大きいところに理論とデータに誤
差が高いので、さらに精密な手続きを用いて
検討を重ねてみた。つまり、血圧値の分布が
ガウス分布でないとして、データから求める。
血圧値の $s^2=g(x)$ についてデータから適当
とみなせる形を想定する、ということを行い、
誤差を含む現実の血圧分布に近づけた。

こうして計算した結果とデータを目盛って
みると図2のようにきわめてよい一致を得た。
以上の議論からみて、測定誤差のいたずら
によって思わぬ歪んだ様相があらわれてきた
ことが示された。つまり、測定誤差によって
われわれの得た“データの相関関係”の様相
がよく説明されたことがわかった。T 年も T
+1 年の測定も、もし測定誤差がなければま
ったく相等しい値を示すものとしても、T 年
からみた T+1 年の平均値は 45° の線上にの
らないで、偏った直線上にのることになるわ
けである。つまり、帰帰直線は 45° とならない
わけである。これが 45° にのったとしたら、む
しろ T 年と T+1 年とはまったく相等しい
値を示すものではない——測定誤差がないと
き——といえることになる。