

雑感 データの追加と分散の変化

■ あるデータの平均値や分散などが分かっているとす。そこに1つのデータを付け加えたとき、平均や分散はどのように変化するだろうか。

もとのデータを $\{x_1, x_2, \dots, x_n\}$ とし、平均値を \bar{x} 、分散を σ^2 とす。そこに、データ x_{n+1} を付け加えたときの平均を \bar{x}' 、分散を σ'^2 とす。

■ 平均値については容易であり、 $(n+1)\bar{x}' = n\bar{x} + x_{n+1}$ であるから、

$$\bar{x}' = \frac{n\bar{x} + x_{n+1}}{n+1} \text{ である。}$$

■ 分散は容易ではない。ちなみに、 $\sigma^2 = \frac{1}{n} \sum_{k=1}^n x_k^2 - \left(\frac{1}{n} \sum_{k=1}^n x_k \right)^2$ で

$$\text{あり、} \sigma'^2 = \frac{1}{n+1} \sum_{k=1}^{n+1} x_k^2 - \left(\frac{1}{n+1} \sum_{k=1}^{n+1} x_k \right)^2 \text{ である。}$$

これらの関係を明らかにするには、こつこつ計算してもよいのだが、以前の「雑感 128 合併2集団の分散」の次の結果を利用するとよい。

集団 X : データ数 m , 平均値 \bar{x} , 標準偏差 σ_X
 集団 Y : データ数 n , 平均値 \bar{y} , 標準偏差 σ_Y
 とし、2集団 X, Y を合併した集団 Z に対して、
 集団 Z : データ数 $m+n$, 平均値 M , 標準偏差 σ
 とするとき、

$$\sigma^2 = \frac{1}{m+n} \{ \sigma_X^2 m + \sigma_Y^2 n + (\bar{x} - M)^2 m + (\bar{y} - M)^2 n \}$$

これによれば、集団 X を元の集団、 $Y = \{x_{n+1}\}$ とし、 $\sigma_Y = 0$ に注意して、

$$\begin{aligned} \sigma'^2 &= \frac{1}{n+1} \{ \sigma^2 n + (\bar{x} - \bar{x}')^2 n + (x_{n+1} - \bar{x}')^2 \} \\ &= \frac{n}{n+1} \left\{ \sigma^2 + \left(\frac{\bar{x} - x_{n+1}}{n+1} \right)^2 + \left(\frac{\bar{x} - x_{n+1}}{n+1} \right)^2 n \right\} \text{ となる。} \end{aligned}$$

■ さて、データ x_{n+1} によって、分散が増加する場合と減少する場合があるが、それはどのような場合であろうか。

$x_{n+1} = \bar{x} + k\sigma$ ($k > 0$) とおいて、 k についての条件を求める。

$$\sigma'^2 - \sigma^2 \text{ を計算すると、} \sigma'^2 = \frac{\sigma^2 n}{(n+1)^2} (n+1+k^2) \text{ であるから、}$$

$$\sigma'^2 - \sigma^2 = \frac{n(n+1+k^2) - (n+1)^2}{(n+1)^2} \sigma^2 \text{ より}$$

$$\sigma' > \sigma \Leftrightarrow k > \sqrt{\frac{n+1}{n}}, \quad \sigma' < \sigma \Leftrightarrow k < \sqrt{\frac{n+1}{n}}$$

を得る。

■ なるほど！ 素晴らしい結果が得られたことになる。

n が一定大きければ、

$\sqrt{\frac{n+1}{n}} \doteq 1$ であるから、右のようなイメージになる。

区間 $[\bar{x} - \sigma, \bar{x} + \sigma]$ に追加すれば分散は小さくなり、それ以外(両サイド)に追加すれば分散は大きくなる。

この結果は、標準偏差 σ の意味を再認識させる。

■ なお、この結果は、以前の「雑感 93 データの分析のオリジナル問題」の(2)も関連しているものと考えられる。

