

平均の再定義と拡張

——中央値，レンジと最小二乗法——

かめさん @cogitoergosumkm

x_1, x_2, \dots, x_n をデータとする.

定義 0.1 (平均と分散).

平均 μ , 分散 σ は

$$\mu = \frac{1}{n} \sum_{k=1}^n x_k \quad (0.1)$$

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \mu)^2 \quad (0.2)$$

と定義される.

1 平均の再定義

定義 1.1 (平均と分散).

$$g(\mu) = \sum_{k=1}^n (x_k - \mu)^2 \quad (1.1)$$

とおく. $g(\mu)$ を最小にする μ を x_1, x_2, \dots, x_n の平均, $\min_{\mu} \{g(\mu)/n\}$ を分散と定義する.

先の定義と同じことを確認する.

$$g(\mu) = \sum_{k=1}^n (x_k - \mu)^2 \quad (1.2)$$

$$= n\mu^2 - 2\mu \sum_{k=1}^n x_k + \sum_{k=1}^n x_k^2 \quad (1.3)$$

$$= n \left(\mu - \frac{1}{n} \sum_{k=1}^n x_k \right)^2 + \sum_{k=1}^n x_k^2 - \frac{1}{n} \left(\sum_{k=1}^n x_k \right)^2 \quad (1.4)$$

したがって, $\mu = \frac{1}{n} \sum_{k=1}^n x_k$ のときに $g(\mu)$ は最小値をとる. その値を計算すると以下ようになる.

$$\frac{1}{n} g \left(\frac{1}{n} \sum_{k=1}^n x_k \right) = \frac{1}{n} \sum_{k=1}^n x_k^2 - \left(\frac{1}{n} \sum_{k=1}^n x_k \right)^2 \quad (1.5)$$

$$= \frac{1}{n} \sum_{k=1}^n \left(x_k - \frac{1}{n} \sum_{k=1}^n x_k \right)^2 \quad (1.6)$$

2 平均の拡張 ——中央値，レンジ——

この節においては, データを並べ替えて, $x_1 \leq x_2 \leq \dots \leq x_n$ としておく.

定義 2.1 (p ノルム).

$\mathbf{v} = (v_1, v_2, \dots, v_n)$ に対して, $1 \leq p < \infty$ として, p ノルムを

$$\|\mathbf{v}\|_p = \left(\sum_{k=1}^n |v_k|^p \right)^{\frac{1}{p}} \quad (2.1)$$

と定義し, $p \rightarrow \infty$ を考えて,

$$\|\mathbf{v}\|_\infty = \max\{v_1, v_2, \dots, v_n\} \quad (2.2)$$

$\mathbf{x} = (x_1, x_2, \dots, x_n)$, $\boldsymbol{\eta} = (1, 1, \dots, 1)$ とする. $g_p(\mu) = \|\mathbf{x} - \mu\boldsymbol{\eta}\|_p$ とすると, 前節の $g(\mu)$ は $g_2(\mu)$ と等しい.

定義 2.2.

$1 \leq p \leq \infty$ に対して, $g_p(\mu)$ を最小にする μ を μ_p と定義する. さらに, $\sigma_p = n^{-\frac{1}{p}} g_p(\mu_p)$ と定義する.

$p = 2$ のとき, μ_2 は平均となり, σ_2 は標準偏差, つまり σ_2^2 は分散となる.

ほかによく用いるものとして, $p = 1$ を考えると,

$$g_1(\mu) = \sum_{k=1}^n |x_k - \mu| \quad (2.3)$$

である.

$$\frac{d}{d\mu} |x_k - \mu| = \begin{cases} -1 & (\mu < x_k) \\ 1 & (\mu > x_k) \end{cases} \quad (2.4)$$

であるから, n が奇数なら, $\mu_1 = x_{\frac{n+1}{2}}$ となり, n が偶数なら, μ_1 は区間 $[x_{\frac{n}{2}}, x_{\frac{n}{2}+1}]$ 内の任意の数である. 偶数の時は $\mu_1 = \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$ をとれば, μ_1 は中央値である.

さらに, $p = \infty$ のときは,

$$g_\infty(\mu) = \max_{1 \leq i \leq n} |x_i - \mu| \quad (2.5)$$

だから, $\mu_\infty = \frac{x_1 + x_n}{2}$ であり, $p \rightarrow \infty$ で $n^{-\frac{1}{p}} \rightarrow 0$ だから, $\sigma_\infty = g_\infty(\mu_\infty) = \frac{x_n - x_1}{2}$ である. これはレンジ $x_n - x_1$ の半分になっている.

3 平均の拡張 —— 最小二乗法 ——

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ をデータとする. X_1, X_2, \dots, X_n の平均を \bar{X} , 分散を σ_X^2 と書くこととする. また, Y_1, Y_2, \dots, Y_n もあるとき, 共分散を

$$\text{Cov}(X, Y) := \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})(Y_k - \bar{Y}) \quad (3.1)$$

$$= \overline{(X - \bar{X})(Y - \bar{Y})} \quad (3.2)$$

$$= \bar{X} \bar{Y} - \bar{X} \bar{Y} \quad (3.3)$$

と定義する. $g(\mu) = \sum_{k=1}^n (x_k - \mu)^2$ とおけば, g を最小化する μ として x の平均が得られ, $g(\mu) = \sum_{k=1}^n (y_k - \mu)^2$ とおけば, g を最小化する μ として y の平均が得られるが, これでは (x, y) とペアでデータを得たことが活かされない. そこで μ を関数だと考える. g は汎関数となる.

$$g[f] = \sum_{k=1}^n (y_k - f(x_k))^2 \quad (3.4)$$

とはいえ、任意の関数をとって良ければ、データに (x, y) と (x, y') のような、 x は等しく、 y は異なるデータが含まれない限り g の最小値は 0 となり、そのような f は無数に得られてしまう。 $f_1(x), f_2(x), \dots, f_m(x)$ をあらかじめ決めておき、

$$f(x; a_1, a_2, \dots, a_m) = \sum_{k=1}^m a_k f_k(x) \quad (3.5)$$

と定義して、

$$g(a_1, a_2, \dots, a_m) = \sum_{k=1}^n (y_k - f(x_k; a_1, a_2, \dots, a_m))^2 \quad (3.6)$$

を最小化する a_1, a_2, \dots, a_m の組を求めて近似関数を得る。このような方法を最小二乗法という。ただし、 $m < n$ である。もし、 $m \geq n$ ならば、やはり g の最小値は 0 になってしまい、強引に形を合わせたに過ぎないことになってしまう。

特に、直線で近似する場合を考える。

$$f(x; a, b) = ax + b \quad (3.7)$$

$$g(a, b) = \sum_{k=1}^n (y_k - ax_k - b)^2 \quad (3.8)$$

まず、 $X = y - ax$ として、 g を変形する。式 (1.4) より、

$$g(a, b) = \sum_{k=1}^n (X_k - b)^2 \quad (3.9)$$

$$= n \left(b - \frac{1}{n} \sum_{k=1}^n X_k \right)^2 + \sum_{k=1}^n X_k^2 - \frac{1}{n} \left(\sum_{k=1}^n X_k \right)^2 \quad (3.10)$$

$$\frac{1}{n} g(a, b) = \{b - (\bar{y} - a\bar{x})\}^2 + \overline{X^2} - \bar{X}^2 \quad (3.11)$$

さらに、

$$\overline{X^2} - \bar{X}^2 = \overline{(y - ax)^2} - (\bar{y} - a\bar{x})^2 \quad (3.12)$$

$$= a^2(\overline{x^2} - \bar{x}^2) - 2a(\overline{xy} - \bar{x}\bar{y}) + (\bar{y}^2 - \bar{y}^2) \quad (3.13)$$

$$= \sigma_x^2 a^2 - 2\text{Cov}(x, y)a + \sigma_y^2 \quad (3.14)$$

$$= \sigma_x^2 \left(a - \frac{\text{Cov}(x, y)}{\sigma_x^2} \right)^2 - \frac{(\text{Cov}(x, y))^2}{\sigma_x^2} + \sigma_y^2 \quad (3.15)$$

まとめれば、

$$\frac{1}{n} g(a, b) = (b - \bar{y} + a\bar{x})^2 + \sigma_x^2 \left(a - \frac{\text{Cov}(x, y)}{\sigma_x^2} \right)^2 - \frac{(\text{Cov}(x, y))^2}{\sigma_x^2} + \sigma_y^2 \quad (3.16)$$

したがって、

$$a = \frac{\text{Cov}(x, y)}{\sigma_x^2} \quad (3.17)$$

$$b = \bar{y} - \frac{\text{Cov}(x, y)}{\sigma_x^2} \bar{x} \quad (3.18)$$

のとき、 g は最小値をとり、 $y = ax + b$ は直線での近似となる。