

第 2 章 回帰 regression

浅川 伸一

2006 年 10 月 19 日

1 復習

中学生レベル (たぶん?)

問題 1. 二点 $(1, 4)$, $(3, 2)$ を通る直線の方程式を求めよ。

問題 2. 二点 $A(x_a, y_a)$, $B(x_b, y_b)$ を通る直線の方程式を求めよ。

2 回帰 regression

2 変数 (x, y) のデータが n 個あるとする。変数 x から y を予測することを考える。このとき直線で近似できると仮定し、最小自乗法 (最小二乗法ともいう) the least minimum method を用いる。各データ点からの距離の二乗を最小化する直線の方程式を求める。ガウスによって考案された。

2 次元上の直線は以下の式で表せる。

$$y = ax + b \quad (1)$$

n 個の点について

$$y_1 = ax_1 + b \quad (2)$$

$$y_2 = ax_2 + b \quad (3)$$

$$\vdots = \vdots$$

$$y_i = ax_i + b \quad (4)$$

$$\vdots = \vdots$$

$$y_n = ax_n + b \quad (5)$$

$n > 3$ ではすべての点をとる一次関数 (直線) を引くことができない。そこで、直線と各点との距離を最小化することを考える。 i 番目のデータの y の値は $\hat{y}_i = ax_i + b$ と予測できるとする。このとき $|y_i - \hat{y}_i|$ は実測値と予測値

を公式として載せているテキストもある。この式から回帰直線は x が平均のとき y の平均を通ることを示している。

ほとんど同じ手順を踏めば x の y 上への回帰係数を計算することができる²。この時の回帰係数は、 $r_{xy} \frac{s_x}{s_y}$ となる。 y の x への回帰係数を a_y 、 x の y への回帰係数を a_x とすればこれらの回帰係数と相関係数との関係は

$$r_{xy} = \sqrt{a_y a_x} \quad (19)$$

と幾何平均（相乗平均）になっている³。

3 分散

式 (9) へ式 (14),(15) を代入して誤差の分散を求めることを考える。

$$s_e^2 = \sum (y_i - \hat{y}_i)^2 \quad (20)$$

$$= \sum \{y_i - (ax_i + b)\}^2 \quad (21)$$

$$= \sum \left\{ (y_i - \bar{y}) - r_{xy} \frac{s_y}{s_x} (x_i - \bar{x}) \right\}^2 \quad (22)$$

$$= \sum (y_i - \bar{y})^2 - 2r \frac{s_y}{s_x} \sum (x_i - \bar{x})(y_i - \bar{y}) + r^2 \frac{s_y^2}{s_x^2} \sum (x_i - \bar{x})^2 \quad (23)$$

$$= s_y^2 - 2r \frac{s_y}{s_x} s_{xy} + r^2 \frac{s_y^2}{s_x^2} s_x^2 \quad (24)$$

$$= s_y^2 - 2 \frac{s_{xy}}{s_x s_y} \frac{s_y}{s_x} s_{xy} + \left(\frac{s_{xy}}{s_x s_y} \right)^2 \frac{s_y^2}{s_x^2} s_x^2 \quad (25)$$

$$= s_y^2 - \frac{s_{xy}^2}{s_x^2} \quad (26)$$

$$= s_y^2 \left(1 - \frac{1}{s_y^2} \frac{s_{xy}^2}{s_x^2} \right) \quad (27)$$

$$= s_y^2 (1 - r^2) \quad (28)$$

これを誤差変動、あるいは回帰直線周りの分散と呼ぶことがある。

また y の予測値 \hat{y} の分散は x の分散の a^2 倍であから

$$s_{\hat{y}}^2 = a^2 s_x^2 \quad (29)$$

$$= \left(r \frac{s_y}{s_y x} \right)^2 s_x^2 \quad (30)$$

$$= r^2 s_y^2 \quad (31)$$

² x と y とを入れ替えて考えてみれば良い

³全てを足し合わせてデータの個数で割る平均を相加平均、もしくは算術平均という

となる。

式 (28) と (31) の和が全分散 s_y^2 になることから、 y の分散が回帰で説明できる変動と誤差の変動とに分解できることが分かる。これを分散の分解定理というが次回以降でより一般的な形で取り上げる。