

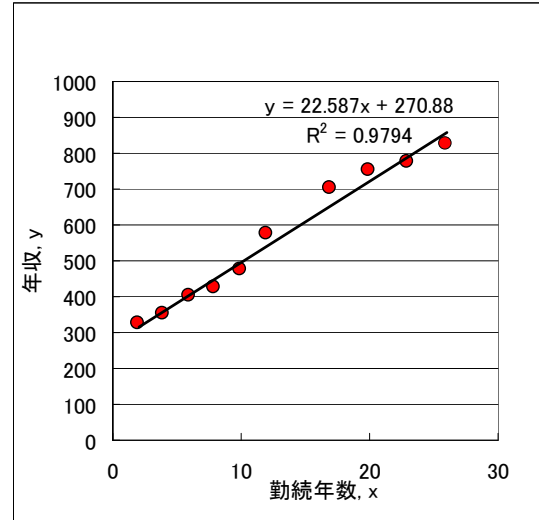
単回帰分析

1つの変数 x から、1つの変数 y を推定する。

例) 勤続年数と年収の関係を分析する。
直線で関係式を表現する。

$$y = ax + b \quad \begin{array}{l} x : \text{説明変数} \\ y : \text{目的変数} \end{array}$$

勤続年数, x	年収, y
2	325
4	350
6	400
8	425
10	475
12	575
17	700
20	750
23	775
26	825



1

単回帰分析 – 最小2乗法による –

推定の誤差の2乗和を考え、これを最小にするように a, b を決定する。

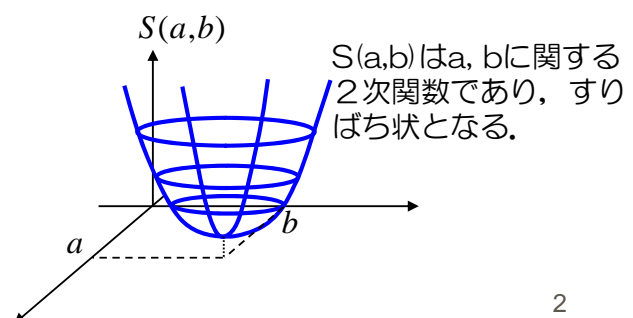
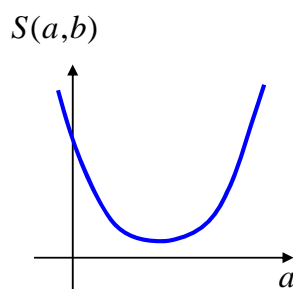
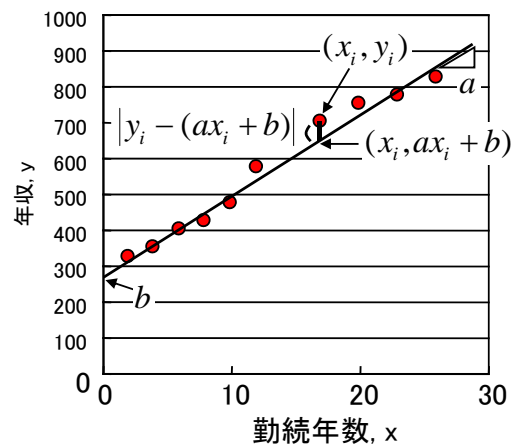
$$S(a, b) = \sum_{i=1}^n [y_i - (ax_i + b)]^2 \rightarrow \min.$$

観測値 推定値

a, b : 回帰パラメータ

条件:

$S(a, b)$ を a, b でそれぞれ微分したものが0でなければならない。



$S(a, b)$ は a, b に関する2次関数であり、すりばち状となる。

2

回帰関係の計算手順

推定の誤差の2乗和を考え、これを最小にするようにa,bを決定する。

$$S(a,b) = \sum_{i=1}^n [y_i - (ax_i + b)]^2 \quad \dots(1)$$

S(a,b)をa, bでそれぞれ微分したものが0でなければならない。

$$\frac{\partial S(a,b)}{\partial a} = \sum_i 2[y_i - (ax_i + b)](-x_i) = 0 \quad \dots(2)$$

$$\frac{\partial S(a,b)}{\partial b} = \sum_i 2[y_i - (ax_i + b)](-1) = 0 \quad \dots(3)$$

(2), (3) は以下のように書き直せる。

$$\sum_i [y_i - (ax_i + b)]x_i = 0 \quad \dots(4)$$

$$\sum_i [y_i - (ax_i + b)] = 0 \quad \dots(5)$$

(4), (5) を整理すると,

$$a \sum_i x_i^2 + b \sum_i x_i = \sum_i x_i y_i \quad \dots(6)$$

$$a \sum_i x_i + nb = \sum_i y_i \quad \dots(7)$$

これは、a,bに関する線形方程式になっている。これを正規方程式(normal equations)という。マトリクスで書けば、

$$\begin{bmatrix} \sum_i x_i^2 & \sum_i x_i \\ \sum_i x_i & n \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_i x_i y_i \\ \sum_i y_i \end{bmatrix} \quad \dots(8)$$

a,bは以下の式により計算できる。

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_i x_i^2 & \sum_i x_i \\ \sum_i x_i & n \end{bmatrix}^{-1} \begin{bmatrix} \sum_i x_i y_i \\ \sum_i y_i \end{bmatrix} \quad \dots(9)$$

3

重回帰分析

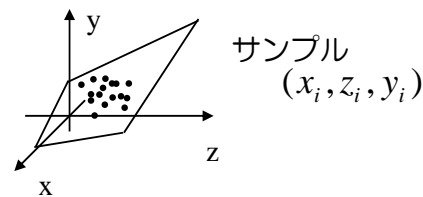
2変数x, zから、1つの変数yを推定する。

$$\bar{y}_{xz} = ax + bz + c \quad \dots(1)$$

1変数の場合と同様、推定の誤差の2乗を評価して、これを最小とするように、回帰係数a, b, cを決定する。

$$S(a,b,c) = ?$$

関数S(a,b,c)を極小とするa,b,cを求めるために、各変数による偏微分をとり、0とおいて解く。



整理すると

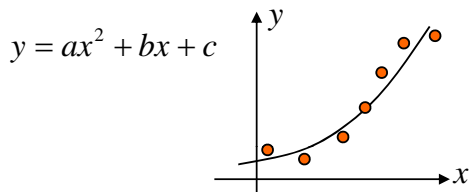
マトリクスで書けば

これを解いて回帰係数を得る。

4

重回帰分析 -高次項を用いた回帰-

課題：1変数の高次式から、1つの変数 y を推定する。



回帰係数を算出する計算式を導きなさい。

一決定係数と相関係数一

得られた標本について、回帰式（モデル式）によるあてはまりの程度を数値化する。

目的変数 y の分散は以下の式で表される。

$$S_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \dots(10)$$

一方、モデル式によって推定された y の値の、実測値からのばらつきを、以下の式によって評価する。

$$S_{y \cdot x}^2 = \frac{1}{n} \sum_{i=1}^n [y_i - (ax_i + b)]^2 \quad \dots(11)$$

(もし、モデル式による予測が完璧ならば、このばらつきは0になる。)

いま、

$$S_r^2 \equiv S_y^2 - S_{y \cdot x}^2 \quad \dots(12)$$

という測度を考え、(10)式の分散との比をとる。

$$r^2 = \frac{S_r^2}{S_y^2} = \frac{S_y^2 - S_{y \cdot x}^2}{S_y^2} = 1 - \frac{S_{y \cdot x}^2}{S_y^2} \quad \dots(13)$$

モデル式による予測が確からしいほど、 S_r^2 は S_y^2 に近づく。すなわち、 r^2 は1に近づく。

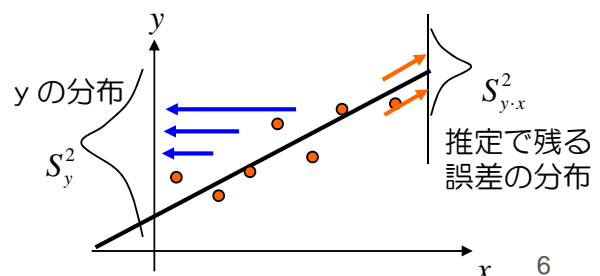
r^2 は決定係数(coefficient of determination)と呼ばれ、以下の範囲をとる。

$$0 \leq r^2 \leq 1$$

また

$$r = \pm \sqrt{r^2} \quad -1 \leq r \leq 1$$

を相関係数と呼ぶ。 r の符号は回帰係数 a の符号に合わせる。



回帰式の相関係数と2変数の相関係数との関係（補足資料）

回帰の誤差は、以下のように書き直せる。

$$\begin{aligned}
 S_{y \cdot x}^2 &= \frac{1}{n} \sum_i [y_i - (ax_i + b)]^2 \\
 &= \frac{1}{n} \sum_i [(y_i - \bar{y}) + \bar{y} - (ax_i + b)]^2 \quad \begin{array}{l} \text{最適化された} a, b \\ \text{に対して以下の} \\ \text{関係が成り立つ} \end{array} \\
 &= \frac{1}{n} \sum_i [(y_i - \bar{y}) + a\bar{x} + b - (ax_i + b)]^2 \quad \leftarrow \bar{y} = a\bar{x} + b \\
 &= \frac{1}{n} \sum_i [(y_i - \bar{y}) - a(x_i - \bar{x})]^2 \\
 &= S_y^2 - 2aS_{xy} + a^2S_x^2 \quad \leftarrow S_{xy} = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y}) \\
 &= S_y^2 - 2\frac{S_{xy}^2}{S_x^2} + \frac{S_{xy}^2}{S_x^4} S_x^2 \quad \leftarrow a = \frac{S_{xy}}{S_x^2} \\
 &= S_y^2 - \frac{S_{xy}^2}{S_x^2} = S_y^2 \left(1 - \frac{S_{xy}^2}{S_x^2 S_y^2} \right)
 \end{aligned}$$

したがって

$$\frac{S_{y \cdot x}^2}{S_y^2} = 1 - \frac{S_{xy}^2}{S_x^2 S_y^2} = 1 - r^2$$

すなわち、

$$r^2 = \frac{S_{xy}^2}{S_x^2 S_y^2}$$

これは、以前に定義した2変数間の相関係数、

$$r = \frac{1}{n} \sum_i \frac{(x_i - \bar{x})(y_i - \bar{y})}{S_x S_y} = \frac{S_{xy}}{S_x S_y}$$

と同等である。

7

重回帰分析 - 重決定係数, 重相関係数 -

2変数の場合も、単回帰と同様に、推定値の、実測値からの分散を考えることができる。

予測式、

$$\bar{y}_{xz} = ax + bz + c \quad \dots(1)$$

を用いて y を推定したときの、誤差の分散は

$$S_{y \cdot xz}^2 = \frac{1}{n} \sum_i [y_i - (ax_i + bz_i + c)]^2$$

いま、以下に示す2つの分散の差を考える。

$$S_r^2 = S_y^2 - S_{y \cdot xz}^2$$

ただし、

$$S_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \dots(10)$$

である。

以下のように、2つの分散の比を考える。

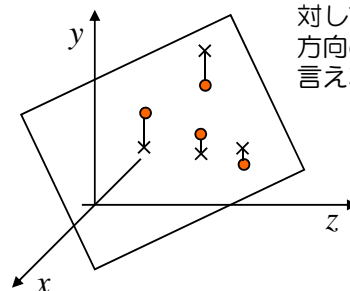
$$R^2 = \frac{S_r^2}{S_y^2} = \frac{S_y^2 - S_{y \cdot xz}^2}{S_y^2} = 1 - \frac{S_{y \cdot xz}^2}{S_y^2}$$

R^2 は、yの分散のうちxとzで説明される部分の割合を示している。

R^2 : 重決定係数

$$R = \sqrt{R^2} \quad \text{: 重相関係数} \quad 0 \leq R \leq 1$$

2変数以上の説明変数に対して、目的変数との方向の一致、不一致を言えないため。



8