

## 多変量解析 目次

1.	回帰分析	1
1. 1	単回帰分析	1
1. 2	重回帰分析	3
1. 3	標準偏回帰係数	6
1. 4	相関係数と決定係数	6
1. 5	回帰式の信頼性	10
1. 6	標準誤差 (Standard Error)	14
1. 7	偏回帰係数の検定	15
1. 8	多重共線性について	18
1. 9	良い重回帰式を作成する	19
1. 11	残差 $\varepsilon$ について	21
2.	重回帰分析例題	25
3.	判別分析	39
3. 1	線形判別式を使用する方法	39
3. 2	ボックスM検定	45
3. 3	マハラノビスの距離による判別	46
3. 4	多変量における2群の母平均の差に関する検定	49
3. 5	判別分析の的中率	50
3. 6	誤判別の確率	50
3. 7	説明変量の寄与	51
3. 8	よい判別式を作成する	52
4.	判別分析例題	55
5.	主成分分析	69
5. 1	主成分を求める	69
5. 2	例題について	74
5. 3	寄与率	76
5. 4	主成分負荷量	77
5. 5	採用する主成分の数について	78
6.	主成分分析例題	79
7.	正準相関分析	89
7. 1	正準相関係数を求める	89
7. 2	正準相関係数の検定	92
8.	正準相関分析例題	93
9.	数量化 I 類	100
9. 1	予測線形式を求める	100
9. 2	カテゴリ数量の基準化	105
9. 3	重相関係数と偏相関係数	106

10.	数量化Ⅰ類例題	108
11.	数量化Ⅱ類	114
11. 1	判別式を求める	114
11. 2	行列を使用して判別式を求める	116
11. 3	カテゴリ数量の基準化	119
11. 4	外的基準に与えるアイテムの影響力について	119
12.	数量化Ⅱ類例題	122
13.	数量化Ⅲ類	126
13. 1	サンプルスコア・カテゴリスコアを求める	126
13. 2	行列を使用して、サンプルスコア・カテゴリスコアを求める	130
13. 3	アイテム・カテゴリ方式	132
14.	数量化Ⅲ類例題	136
15.	EXCELでの行列演算	141
15. 1	関数ウィザードを使用する	141
15. 2	逆行列を求める	143
15. 3	もとの行列が対角行列である時の逆行列を求める	145
15. 4.	行列式を求める	146
15. 5.	行列の積を求める	146
15. 6	行列を使用して、連立方程式を解く	147
15. 7	固有値を求める	147

## 1. 回帰分析

何名かの体重と身長が分かっているとき、体重の値は分かっているが、身長が不明の人がいるとする。このようなとき、すでに得ているデータから身長と体重の関係を調べ、その相関を求め、身長不明の人の身長を予測する。このような分析方法を回帰分析という。

求めるものは身長であり、これを目的変数と呼ぶ。身長の値を予測するのは、体重からであるので、この体重のことを説明変数と呼ぶ。説明変数が1つの時を単回帰分析といい、説明変数が2つ以上の時を多重回帰分析という。

回帰分析では、説明変数は量的データであり、また目的変数も量的データである。

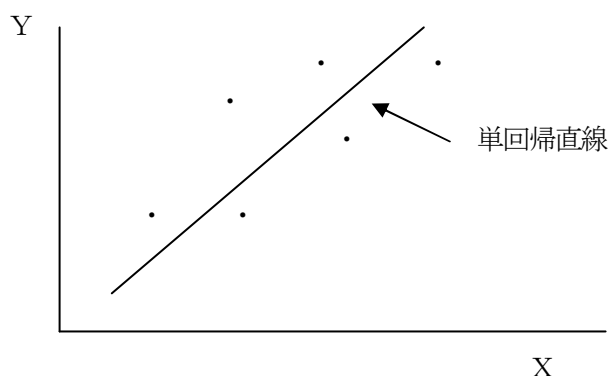
なお、回帰式で予測をするときには、説明変数の範囲内で予測することが望ましい。説明変数の範囲を大きく越えたところで予測すると誤差が大きくなり実用に適さなくなる。

### 1. 1 単回帰分析

正規母集団から抽出して得られた標本データ  $x \cdot y$  が下表のようにあり、 $x \cdot y$  間にある関係があるものとする。

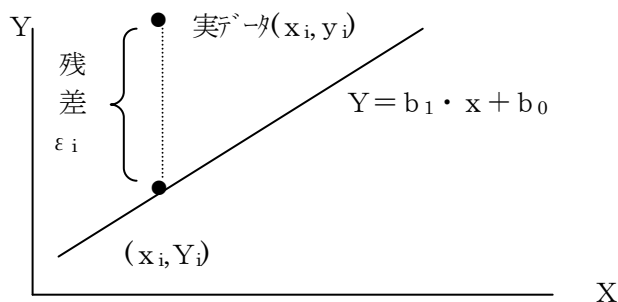
標本	説明変数 $x$	目的変数 $y$
1	$x_1$	$y_1$
2	$x_2$	$y_2$
...	...	...
$n$	$x_n$	$y_n$

以上の標本データをXYグラフで描くいて、下のようになったとする



標本データ  $x \cdot y$  の間には右上がりの関係がありそうなので、 $x$  と  $y$  の関係を表す適当な直線を考える。目的変数  $y$  と説明変数  $x$  との間に相関があるとき、

$Y = b_1 \cdot x + b_0$  なる直線を1本考え、実データとこの直線上の値との差を  $\epsilon$  とする。



$Y = b_1 \cdot x + b_0$  なる直線は全ての標本データについて、その残差が最小になるようにひく必

要がある。この直線から各標本データとのズレ具合を計るために、各残差の平方和をとり、この平方和を最小にするようにする。このような方法を最小2乗法という。

標本データは、直線  $Y = b_1 \cdot x + b_0$  から残差 ( $\varepsilon$ ) 分ずれているので、標本データは  $y = b_1 \cdot x + b_0 + \varepsilon$  と表す。

このことから線形回帰モデルを

$$y_i = \beta_1 \cdot x_i + \beta_0 + \varepsilon_i \quad (i=1,2,\dots,n) \text{ とすると}$$

残差  $\varepsilon$  について、

- $$\left\{ \begin{array}{l} \text{① } \varepsilon_i \text{ と } \varepsilon_j \text{ はお互いに独立であり、正規分布 } N(0, \sigma^2) \text{ に従う。} \\ \text{② } \varepsilon_i \text{ の平均値 (期待値) は } 0 \text{ である。} \\ \text{③ } \varepsilon_i \text{ の分散は一定である。} \end{array} \right.$$

このような仮定下で単回帰式を  $Y = b_1 \cdot x + b_0$  とする。

いま、残差  $\varepsilon$  に注目すると  $\varepsilon_i = y_i - Y_i$   $\varepsilon_i = y_i - b_1 \cdot x_i - b_0$  である。

この残差を全ての標本データについて合計し、その合計値を最小にするような  $b_0 \cdot b_1$  を求め、この単回帰式を得る。

$$\sum \varepsilon_i^2 = \sum (y_i - b_1 \cdot x_i - b_0)^2 \text{ であるから}$$

$$f = \sum (y_i - b_1 \cdot x_i - b_0)^2 \text{ とすると}$$

この式を  $b_0, b_1$  で偏微分して、0とおくことにより、正規方程式を得て、式  $f$  を最小にする  $b_0 \cdot b_1$  を得ることができる。

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial b_1} = -2 \sum x_i \cdot (y_i - b_1 \cdot x_i - b_0) = 0 \\ \frac{\partial f}{\partial b_0} = -2 \sum (y_i - b_1 \cdot x_i - b_0) = 0 \end{array} \right.$$

これから

$$\left\{ \begin{array}{l} b_1 = \frac{n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i}{n \cdot \sum x_i^2 - (\sum x_i)^2} \\ b_0 = \frac{\sum x_i^2 \cdot \sum y_i - \sum x_i \cdot y_i \cdot \sum x_i}{n \cdot \sum x_i^2 - (\sum x_i)^2} \end{array} \right.$$

また

$$b_1 = \frac{n \cdot \sum x_i \cdot y_i - \sum x_i \cdot \sum y_i}{n \cdot \sum x_i^2 - (\sum x_i)^2} = \frac{\sum x_i \cdot y_i - \frac{\sum x_i \cdot \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

ただし、 $x_i \cdot y_i$  の偏差平方和・積和をそれぞれ  $S_{xx} \cdot S_{yy} \cdot S_{xy}$  とすると

$$S_{xx} = \sum (x_i - \bar{x})^2 \quad S_{yy} = \sum (y_i - \bar{y})^2$$

$$S_{xy} = \sum (x_i - \bar{x}) \cdot (y_i - \bar{y}) \text{ である。}$$

以上から、単回帰式は

$$Y - \bar{y} = b_1 (x - \bar{x}) \rightarrow Y = \frac{S_{xy}}{S_{xx}} (x - \bar{x}) + \bar{y} \text{ と表される。}$$

また  $x$  と  $y$  の相関係数を  $R_{xy}$  とすると

$$R_{xy} = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} \quad \text{であるから、}$$

相関係数を使用して単回帰式を表すと

$$Y = R \cdot \frac{\sqrt{S_{yy}}}{\sqrt{S_{xx}}} (x - \bar{x}) + \bar{y} \quad \text{と書くことができる。}$$

### 1. 2 重回帰分析

それでは次に説明変数が  $x_1 \cdot x_2$  の 2 変量になったときの回帰式を求める。

標本	説明変数 $x_1$	説明変数 $x_2$	目的変数 $y$
1	$X_{11}$	$X_{21}$	$Y_1$
2	$X_{12}$	$X_{22}$	$Y_2$
...	...	...	...
$n$	$X_{1n}$	$X_{2n}$	$Y_n$

説明変数が 2 変量あるので、単純に説明 2 変量 ( $x_1$  と  $x_2$ ) の平均値をとって、その値と目的変数 ( $y$ ) との相関を求めても、平均値をとる段階で失う情報量が大きいため正しい回帰式を得ることができない。

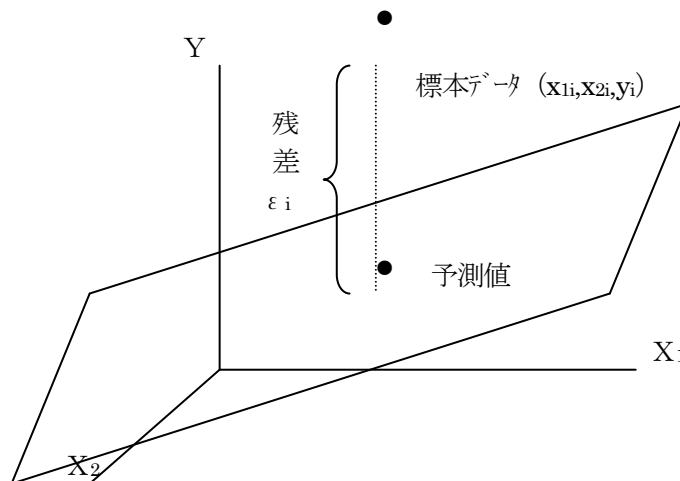
このように説明変数が 2 つ以上ある時の回帰分析を重回帰分析という。

#### 1.2.1 重回帰式を求める。

2 説明変数が次のようになっているときの重回帰直線を求める。

標本	説明変数 $x_1$	説明変数 $x_2$	目的変数 $y$
1	$X_{11}$	$X_{21}$	$Y_1$
2	$X_{12}$	$X_{22}$	$Y_2$
...	...	...	...
$n$	$X_{1n}$	$X_{2n}$	$Y_n$
平均	$\bar{x}_1$	$\bar{x}_2$	$\bar{y}$

この関係を図で表すと



説明変量 ( $x_1, x_2$ ) と目的変量 ( $y$ ) との間に相関関係があるとき

$$Y = b_1 \cdot x_1 + b_2 \cdot x_2 + b_0$$

なる平面を考え、実際の標本データからこのこの平面上への残差を  $\varepsilon$  とすると  
説明変量が2つある時の重回帰式は

$$y_i = b_1 \cdot x_{1i} + b_2 \cdot x_{2i} + b_0 + \varepsilon_i \quad \text{と表される。}$$

残差  $\varepsilon$  に注目すると  $\varepsilon_i = y_i - Y_i$

$$\varepsilon_i = y_i - (b_1 \cdot x_{1i} + b_2 \cdot x_{2i} + b_0) \quad \text{であるから}$$

この残差平方和を求め、残差平方和が最小にするような  $b_0 \cdot b_1 \cdot b_2$  を求めると、重回帰式を得ることができる。

一般に説明変量が  $p$  個ある時の線形重回帰モデルは

$$y_i = \beta_1 \cdot x_{1i} + \beta_2 \cdot x_{2i} + \dots + \beta_p \cdot x_{pi} + \beta_0 + \varepsilon_i \quad (i=1,2 \dots n)$$

と表される。この時単回帰分析と同様に

残差  $\varepsilon$  について、

- $$\left\{ \begin{array}{l} \textcircled{1} \varepsilon_i \text{ と } \varepsilon_j \text{ はお互いに独立であり、正規分布 } N(0, \sigma^2) \text{ に従う。} \\ \textcircled{2} \varepsilon_i \text{ の平均値 (期待値) は } 0 \text{ である。} \\ \textcircled{3} \varepsilon_i \text{ の分散は一定である。} \end{array} \right.$$

との仮定下で重回帰予測式を

$$Y_i = b_1 \cdot x_{1i} + b_2 \cdot x_{2i} + \dots + b_p \cdot x_{pi} + b_0 \quad \text{とする。}$$

$b_1 \cdot b_2 \dots b_p$  を偏回帰係数といい、 $\beta_1 \cdot \beta_2 \dots \beta_p$  を母偏回帰係数という。

[残差平方和  $\sum (\varepsilon_i)^2$  を最小にするような  $b_0 \cdot b_1 \cdot b_2$  を求める。]

$\sum (\varepsilon_i)^2 = \sum \{y_i - (b_1 \cdot x_{1i} + b_2 \cdot x_{2i} + b_0)\}^2$  を最小にする  $b_0 \cdot b_1 \cdot b_2$  を求める。

$f = \sum (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0)^2$  とし、この式を  $b_0 \cdot b_1 \cdot b_2$  で偏微分する。

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial b_1} = -2 \sum x_{1i} \cdot (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \\ \frac{\partial f}{\partial b_2} = -2 \sum x_{2i} \cdot (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \\ \frac{\partial f}{\partial b_0} = -2 \sum (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \end{array} \right.$$

これより

$$\left\{ \begin{array}{l} \sum x_{1i} \cdot (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \quad \dots \textcircled{1} \\ \sum x_{2i} \cdot (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \quad \dots \textcircled{2} \\ \sum (y_i - b_1 \cdot x_{1i} - b_2 \cdot x_{2i} - b_0) = 0 \quad \dots \textcircled{3} \end{array} \right.$$

上の式を正規方程式という

$$\textcircled{3} \text{ から } \sum y_i - b_1 \cdot \sum x_{1i} - b_2 \cdot \sum x_{2i} - \sum b_0 = 0$$

$\sum b_0 = n \cdot b_0$  であるから

$$b_0 = \frac{\sum y_i - b_1 \cdot \sum x_{1i} - b_2 \cdot \sum x_{2i}}{n} = \overline{y} - b_1 \cdot \overline{x_1} - b_2 \cdot \overline{x_2}$$

これを①②に代入して

$$b_1 \cdot (\sum x_{1i}^2 - n \cdot \overline{x_1}^2) + b_2 \cdot (\sum x_{1i} \cdot x_{2i} - n \cdot \overline{x_1} \cdot \overline{x_2}) = \sum x_{1i} \cdot y_i - n \cdot \overline{x_1} \cdot \overline{y}$$

$$b_1 \cdot (\sum x_{1i} x_{2i} - n \cdot \overline{x_1} \cdot \overline{x_2}) + b_2 \cdot (\sum x_{2i}^2 - n \cdot \overline{x_2}^2) = \sum x_{2i} \cdot y_i - n \cdot \overline{x_2} \cdot \overline{y}$$

これより  $b_0 \cdot b_1 \cdot b_2$  を求めると、重回帰式の係数を得ることができる。

### 1.2.2 偏差平方和・積和から重回帰式を求める

(1)説明変数が2個の時

説明変数  $x_1 \cdot x_2$  の偏差平方和それぞれ  $S_{11} \cdot S_{22}$ 、偏差積和を  $S_{12}$  とすると

$$\begin{cases} S_{11} = \sum (x_{1i} - \bar{x}_1)^2 = \sum x_{1i}^2 - n \cdot \bar{x}_1^2 \\ S_{22} = \sum (x_{2i} - \bar{x}_2)^2 = \sum x_{2i}^2 - n \cdot \bar{x}_2^2 \\ S_{12} = \sum (x_{1i} - \bar{x}_1) \cdot (x_{2i} - \bar{x}_2) = \sum x_{1i} \cdot x_{2i} - n \cdot \bar{x}_1 \cdot \bar{x}_2 \end{cases}$$

また目的変数  $y$  と説明変数  $x_1$  との偏差積和を  $S_{y1}$  とすると

$$S_{y1} = \sum (x_{1i} - \bar{x}_1) \cdot (y_i - \bar{y}) = \sum x_{1i} \cdot y_i - n \cdot \bar{x}_1 \cdot \bar{y}$$

目的変数  $y$  と説明変数  $x_2$  との偏差積和を  $S_{y2}$  とすると

$$S_{y2} = \sum (x_{2i} - \bar{x}_2) \cdot (y_i - \bar{y}) = \sum x_{2i} \cdot y_i - n \cdot \bar{x}_2 \cdot \bar{y}$$

以上から前の式は

$$\begin{cases} S_{11} \cdot b_1 + S_{12} \cdot b_2 = S_{y1} \\ S_{12} \cdot b_1 + S_{22} \cdot b_2 = S_{y2} \end{cases}$$

となるので、これから係数  $b_0 \cdot b_1 \cdot b_2$  を求める。

また  $b_0 = \bar{y} - (b_1 \cdot \bar{x}_1 + b_2 \cdot \bar{x}_2)$  である。

(2)説明変数が  $p$  個ある時

平方和・積和行列を

$$\mathbf{S} = \begin{pmatrix} S_{11} & S_{12} & \cdots & S_{1p} \\ S_{21} & S_{22} & \cdots & S_{2p} \\ \cdots & \cdots & \cdots & \cdots \\ S_{p1} & S_{p2} & \cdots & S_{pp} \end{pmatrix} \quad \text{とする。}$$

求める重回帰式を、 $Y_i = b_1 \cdot x_{1i} + b_2 \cdot x_{2i} + \cdots + b_p \cdot x_{pi} + b_0$  とする、この回帰式の係数  $b_0 \cdot b_1 \cdots b_p$  は、下の連立方程式の解として与えられる。

$$\begin{cases} S_{11} \cdot b_1 + S_{12} \cdot b_2 + \cdots + S_{1p} \cdot b_p = S_{y1} \\ S_{21} \cdot b_1 + S_{22} \cdot b_2 + \cdots + S_{2p} \cdot b_p = S_{y2} \\ \cdots \\ S_{p1} \cdot b_1 + S_{p2} \cdot b_2 + \cdots + S_{pp} \cdot b_p = S_{yp} \end{cases}$$

係数  $b_0$  は、 $b_0 = \bar{y} - (b_1 \cdot \bar{x}_1 + b_2 \cdot \bar{x}_2 + \cdots + b_p \cdot \bar{x}_p)$  である。

また行列では

$$\mathbf{y} = \begin{pmatrix} y_{11} \\ y_{12} \\ \cdots \\ y_{1n} \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} b_0 \\ b_1 \\ \cdots \\ b_n \end{pmatrix}$$

とすると

係数  $\mathbf{b}$  は、 $\mathbf{b} = (\mathbf{x}' \cdot \mathbf{x})^{-1} \cdot \mathbf{x}' \cdot \mathbf{y}$  で求めることができる。

### 1. 3 標準偏回帰係数

説明変数がどれくらい目的変数に影響を与えているか（寄与しているか）を見るには、求めた重回帰式の偏回帰係数を見ればよい。通常、偏回帰係数が大きいほど目的変数に与える影響が大きいので多く寄与しているといえる。しかし、説明変数間で単位が異なるときには、単位の影響を受けるので、単純に偏回帰係数の大小比較して決めることはできない。単位の影響を除くには、標本データを標準化する。データを標準化することにより、平均=0・分散=1となり単位の影響を受けなくなるので、標準化したデータから偏回帰係数を求めるようにする。このように標準化したデータから得られた偏回帰係数を、標準偏回帰係数という。

標準偏回帰係数の大きいほど、目的変数に与える影響が大きく、寄与の大きい変数であるといえる。

通常説明変数が2つの時の重回帰式は

$$Y - \bar{y} = b_1 (x_1 - \bar{x}_1) + b_2 (x_2 - \bar{x}_2) \text{ と書ける。}$$

いま、目的変数の標準偏差を $\sqrt{S_{yy}}$ 、説明変数 $x_1$ の標準偏差を $\sqrt{S_{11}}$ 、説明変数 $x_2$ の標準偏差を $\sqrt{S_{22}}$  とすると

データの標準化は

$$Y \rightarrow \frac{Y - \bar{y}}{\sqrt{S_{yy}}} \quad x_1 \rightarrow \frac{x_1 - \bar{x}_1}{\sqrt{S_{11}}} \quad x_2 \rightarrow \frac{x_2 - \bar{x}_2}{\sqrt{S_{22}}} \text{ を行うことである。}$$

$$\frac{Y - \bar{y}}{\sqrt{S_{yy}}} = b_1 \cdot \frac{x_1 - \bar{x}_1}{\sqrt{S_{yy}}} + b_2 \cdot \frac{x_2 - \bar{x}_2}{\sqrt{S_{yy}}} = b_1 \cdot \frac{\sqrt{S_{11}}}{\sqrt{S_{yy}}} \cdot \frac{x_1 - \bar{x}_1}{\sqrt{S_{11}}} + b_2 \cdot \frac{\sqrt{S_{22}}}{\sqrt{S_{yy}}} \cdot \frac{x_2 - \bar{x}_2}{\sqrt{S_{22}}}$$

データを標準化して得られる重回帰式の係数 $b_1'$ ・ $b_2'$ は

$$b_1' = b_1 \cdot \frac{\sqrt{S_{11}}}{\sqrt{S_{yy}}} \quad b_2' = b_2 \cdot \frac{\sqrt{S_{22}}}{\sqrt{S_{yy}}} \text{ と表すことができる。}$$

### 1. 4 相関係数と決定係数

#### 1.4.1 単回帰式における相関係数と決定係数

標本	説明変数 x	目的変数 y
1	$x_1$	$y_1$
2	$x_2$	$y_2$
...	...	...
n	$x_n$	$y_n$
平均	$\bar{x}$	$\bar{y}$



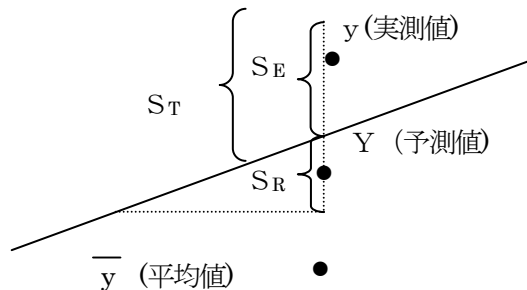
説明変数  $x$  の変化に従って目的変数  $y$  が変化する（相関関係にある）とき  $x$  と  $y$  の間の相関係数を  $R$  とすると

$$R = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}}$$

相関係数  $R$  は  $-1 \leq R \leq 1$  の値をとり

$$\begin{cases} R > 0 & \cdots \text{ 正の相関がある} \\ R = 1 & \cdots \text{ 無相関} \\ R < 0 & \cdots \text{ 負の相関がある} \end{cases}$$

いま説明変数  $x$  と実測値  $y$  との関係が  $r$  である時、これから求めた単回帰式を  $Y = b_1 \cdot x + b_0$  とすると



実測値  $y$  は、単回帰直線の付近にばらついて散在している。このばらつきの小さいほど単回帰式のあてはまりがよい（精度が高い）直線といえる。また説明変数  $x$  の目的変数に与える影響が大きいといえる。つまり決定力が大きいといえる。

分散状況を見ると、全分散 ( $S_T$ ) は、実測値  $y_i$  が平均値  $\bar{y}$  からどれ位分散しているかであるので、 $\sum (y_i - \bar{y})^2$ 。回帰で説明可能な部分の分散 ( $S_R$ )、つまり予測値が平均値からどれ位分散しているかは、 $\sum (Y_i - \bar{y})^2$ 。回帰で説明できない残差部分の分散 ( $S_E$ ) つまり実測値が予測値からどれ位分散しているかは、 $\sum (y_i - Y_i)^2$  である。

これらの変動の間には、

$$S_T = S_R + S_E \quad \text{つまり} \quad \sum (y_i - \bar{y})^2 = \sum (Y_i - \bar{y})^2 + \sum (y_i - Y_i)^2 \quad \text{なる関係がある。}$$

この両辺を  $\sum (y_i - \bar{y})^2$  で割ると

$$1 = \frac{\sum (Y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} + \frac{\sum (y_i - Y_i)^2}{\sum (y_i - \bar{y})^2} \quad 1 - \frac{\sum (y_i - Y_i)^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (Y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{S_R}{S_T} = R^2$$

この  $R^2$  のことを決定係数という。この決定係数は  $0 \leq R^2 \leq 1$  の値をとる。

また、この決定係数  $R^2$  は相関係数  $R$  の 2 乗に等しい。

### 1.4.2 重回帰式における相関係数と決定係数

#### (1) 重相関係数と決定係数

標本No	説明変量				実測値	予測値
	$x_1$	$x_2$	...	$x_p$	$y$	$Y$
1	$x_{11}$	$x_{21}$	...	$x_{p1}$	$y_1$	$Y_1$
2	$x_{12}$	$x_{22}$	...	$x_{p2}$	$y_2$	$Y_2$
...			...		...	...
n	$x_{1n}$	$x_{2n}$	...	$x_{pn}$	$y_n$	$Y_n$
平均	$\bar{x}_1$	$\bar{x}_2$	...	$\bar{x}_p$	$\bar{y}$	$\bar{Y}$

重相関係数  $R$  は、実測値データ  $y$  と重回帰式から求めた予測値データ  $Y$  との相関係数である。

$$R = \frac{\sum (y_i - \bar{y}) \cdot (Y_i - \bar{Y})}{\sqrt{\sum (y_i - \bar{y})^2} \cdot \sqrt{\sum (Y_i - \bar{Y})^2}} = \sqrt{\frac{\sum (Y_i - \bar{Y})^2}{\sum (y_i - \bar{y})^2}}$$

また単回帰のときと同様に、相関係数の2乗を決定係数と呼び、やはり  $0 \leq R^2 \leq 1$  の値をとる。 $R^2$  が1に近いほど重回帰式の精度が高いといえる。

$$R^2 = 1 - \frac{\sum (y_i - Y_i)^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (Y_i - \bar{Y})^2}{\sum (y_i - \bar{y})^2}$$

[重相関係数の検定]

標本から得られた重相関係数について、その母重相関係数 ( $\rho$ ) が無相関かどうかの検定を行う。標本から得られた重相関係数を  $R$  とする時、その母相関係数 ( $\rho$ ) について  $\rho = 0$  の仮説につき、検定統計量を  $F$  とすると

$$F = \frac{R^2/p}{(1-R^2)/(n-p-1)} \quad \text{ただし } p : \text{説明変量の個数 } n : \text{標本数 } R : \text{重相関係数}$$

は、自由度  $p$ ,  $n-p-1$  の  $F$  分布に従うことを利用して検定を行う。

検定をおこなう

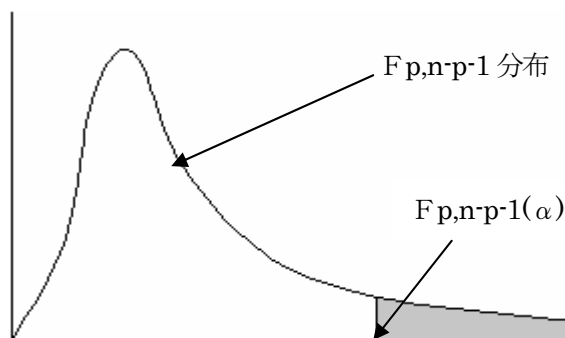
(1) 仮説をたてる

仮説  $H_0 : \rho = 0$  (母重相関係数は無相関である)

対立仮説  $H_1 : \rho \neq 0$  (母重相関係数は無相関ではない)

(2) 検定統計量  $F$  は自由度  $p$ ,  $n-p-1$  の  $F$  分布に従う。

(3) 有為水準  $\alpha$  で検定を実行する。



$F \geq F_{p, n-p-1}(\alpha)$ であれば、仮説を棄却する。つまり、母重相関係数は有効であり、実測値と予測値の間には相関があるといえる。

重相関係数は、実測値  $y$  と予測値  $Y$  との相関係数である。これに対して単純に2変量間の相関係数を単相関係数という。多変量データにおいて、2変量間の相関係数が本当に正しい相関を示すとは限らない。多変量においては2変量間の相関係数を求めても、その2変量以外の変量がこの2変量に影響を与えるからである。よって、多変量間における2変量の正しい相関係数を求めるには、相関係数を求める2変量以外の変量の影響を取り除いて（一定にして）相関係数を求める必要がある。このようにして求めた相関係数を偏相関係数という。

### (2) 偏相関係数

多変量データにおいて、任意の2変量間の単純な相関係数を単相関係数というが、これは相関をとる2変量以外の変量が、その2変量に影響を与えている相関係数である。これに対し、相関を求める2変量以外の他の変量の影響を取り除いた2変量間の相関係数を偏相関係数という。

いま  $P$  変量の任意の2変量間の単相関係数を  $r_{ij}$  とする。

	$X_1$	$X_2$	...	$X_p$
$X_1$	$r_{11}$	$r_{21}$	...	$r_{p1}$
$X_2$	$r_{21}$	$r_{22}$	...	$r_{p2}$
...			...	
$X_p$	$r_{p1}$	$r_{2p}$	...	$r_{pp}$

単相関行列を  $R$  とすると

$$R = \begin{pmatrix} r_{11} & r_{21} & \dots & r_{p1} \\ r_{21} & r_{22} & \dots & r_{p2} \\ \dots & \dots & \dots & \dots \\ r_{p1} & r_{2p} & \dots & r_{pp} \end{pmatrix}$$

逆行列  $R^{-1}$  を

$$R^{-1} = \begin{pmatrix} r^{11} & r^{21} & \dots & r^{p1} \\ r^{21} & r^{22} & \dots & r^{p2} \\ \dots & \dots & \dots & \dots \\ r^{p1} & r^{2p} & \dots & r^{pp} \end{pmatrix}$$

成分  $ij$  以外の変量を一定にした成分  $i \cdot j$  間の偏相関係数を  $r_{ij \cdot p \dots q}$  とする。

$$r_{ij \cdot p \dots q} = \frac{-r^{ij}}{\sqrt{r^{ii} \cdot r^{jj}}}$$

### (3) 自由度調整済み決定係数

決定係数や重相関係数は、説明変数の数を増やすと単純に増加する傾向がある。

そこで、単純に説明変数の数を増やしても、決定係数が単純に増加しないように調整した自由度調整済み決定係数という。通常標本数が  $n$  個、説明変数が  $n - 1$  個のものは分析することができない。必ず説明変数が  $n - 2$  個以下にする必要がある。

自由度調整済み決定係数を  $R'^2$  とすると

$$R'^2 = 1 - \frac{\frac{S_E}{n-p-1}}{\frac{S_T}{n-1}}$$

n : 標本数 P : 説明変数の個数 (P = n - 1 の時には分母が 0 になってしまう)  
 また書き換えると

$$R'^2 = 1 - \frac{n-1}{n-p-1} \cdot (1-R^2)$$

自由度調整済み重相関係数を R' とすると  $R' = \sqrt{R'^2}$  である。

### 1. 5 回帰式の信頼性

回帰式を使用して説明変数から目的変数の値を予測する時、その予測値がどのくらい信頼性があるのかを検定する方法に、分散分析を用いる方法と相関係数を用いる方法がある。

#### 1.5.1 分散分析を用いる場合

##### (1) 単回帰のとき

説明変数 x と実測値 y と単回帰式から求めた予測値 Y が下表のようである時

標本	説明変数 x	実測値 y	予測値 Y
1	x <sub>1</sub>	y <sub>1</sub>	Y <sub>1</sub>
2	x <sub>2</sub>	y <sub>2</sub>	Y <sub>2</sub>
...	...	...	...
n	x <sub>n</sub>	y <sub>n</sub>	Y <sub>n</sub>

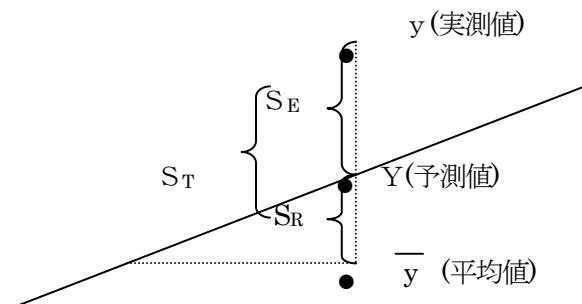
予測値 Y<sub>i</sub> は、 $Y = b_1 \cdot x + b_0$  の回帰式から求めた値

以上のデータをもとに、分散分析表を作成し回帰式の信頼性を検定する。

全体の変動 (S<sub>T</sub>) を、回帰による変動 (S<sub>R</sub>) と残差による変動 (S<sub>E</sub>) とに分け、回帰による変動が残差による変動よりも小さいようであれば、回帰直線で求めた予測値は残差による影響の方が大きいので予測には役立たないを考える。

実測値の変動 (S<sub>T</sub>) = 回帰による変動 (S<sub>R</sub>) + 残差による変動 (S<sub>E</sub>)

残差が小さいほど「実測値の変動」≒「回帰による変動」となり、よい予測値を得られる。



##### (1) 変動を求める

① 実測値の全変動 (S<sub>T</sub>) … 実測値の各値 y<sub>i</sub> が、実測値の平均  $\bar{y}$  からどれ位ばらついている

かである。

$$S_T = \sum (y_i - \bar{y})^2$$

②回帰による変動 ( $S_R$ ) …回帰直線によって求めた予測値 $Y_i$ が、実測値の平均 $y$ からどれ位ばらついているかである。  $S_R = \sum (Y_i - \bar{y})^2$

③残差による変動 ( $S_E$ )  $S_E = \sum (y_i - Y_i)^2$

(2)自由度を求める

①回帰による変動の自由度 ( $f_R$ )  $f_R = 2 - 1 = 1$

②残差による変動の自由度 ( $f_E$ )  $f_E = n - 2$

③全変動の自由度 ( $f_T$ )  $f_T = f_R + f_E = n - 1$

(3)不偏分散を求める

①回帰による変動の不偏分散 ( $V_R$ )  $V_R = \frac{S_R}{f_R} = \frac{S_R}{1}$

②残差による変動の不偏分散 ( $V_E$ )  $V_E = \frac{S_E}{f_E} = \frac{S_E}{n-2}$

③全変動の不偏分散 ( $V_T$ )  $V_T = \frac{S_T}{f_T} = \frac{S_T}{n-1}$

(4)分散比Fを求める

$F = \frac{V_R}{V_E} = \frac{S_R}{S_E/(n-2)}$  は自由度1,  $n-2$ のF分布に従う。

右片側検定を行い、 $V_R$ が $V_E$ より大きいかどうか検定する。 $V_R > V_E$ であれば、回帰による変動が残差による変動よりも全変動に与える影響が大きいため、回帰直線は予測に役立つといえる。

(5)検定を行う

(1)仮説をたてる

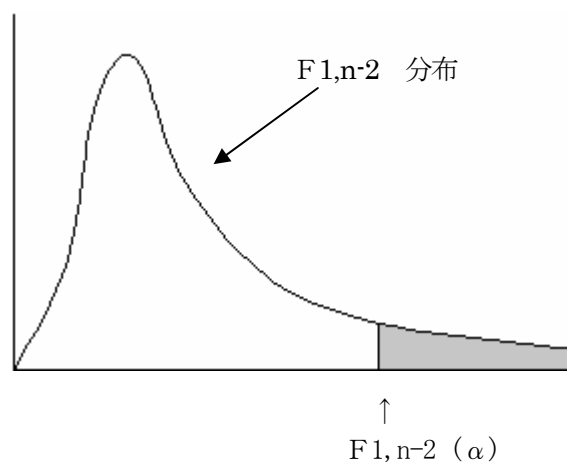
仮説  $H_0$ : 回帰直線は予測に役立たない ( $V_R \leq V_E$ )

対立仮説  $H_1$ : 回帰直線は予測に役立つ ( $V_R > V_E$ )

(2)検定統計量Fを求める

$F = \frac{V_R}{V_E}$  は自由度1,  $n-2$ のF分布に従う

(3)有為水準 $\alpha$ で右片側検定を行う



$F \geq F_{1, n-2}(\alpha)$ であれば、仮説 $H_0$ を棄却し、対立仮説 $H_1$ : 回帰直線は予測に役立つを採択する。つまり、この回帰直線は予測に役立つとする。