



慶應義塾大学ビジネス・スクール

回帰分析シリーズ 1

— 単回帰分析 —

"... there is no one more dangerous than the unthinking user of a computer, who has no real understanding of the nature of the computations being processed inside the machine."

--- Johnston, J. "Econometric Method", 3rd.ed., 1984; pp.79. ---

単回帰分析は、すべての計量経済分析における基礎となるものである。したがってどのような複雑なモデルであっても、基本的には単回帰の延長線上にある。このノートでは、単回帰の概念や仮定、制約、計算方法を例題を用いて解説しているが、その目的は、回帰分析によって「何が出来るか」を知ることだけでなく、その限界を理解することにある。

一般に、関係式に2つの変数しか含まれない確率モデル

$$Y_i = \alpha + \beta X_i + u_i \quad i = 1, 2, \dots, n \quad (1)$$

は単純回帰と呼ばれる。 Y は被説明変数あるいは従属変数、 X は説明変数あるいは独立変数と呼ばれる。ここで α と β はそれぞれ、この関数の切片と勾配を表す未知のパラメータである。また u は Y の変動について X だけでは説明しきれない要因を表し、誤差項あるいは攪乱項と呼ばれる。

回帰モデルによって α と β を推定する場合、誤差項 u の確率分布にいくつかの仮定をおく必要がある。それは分布の平均値が0、分散が均一、自己相関しない、確率分布は正規分布にしたがうというものであり、それぞれ(2)から(4)のように表すことができる。

$$E(u_i) = 0 \quad i = 1, 2, \dots, n \quad (2)$$

$$E(u_i u_j) = \begin{cases} 0 & i \neq j \\ \sigma_u^2 & i = j \end{cases} \quad (3)$$

$$p(u_i) = N(0, \sigma^2) \quad (4)$$

これらの仮定は、一括して

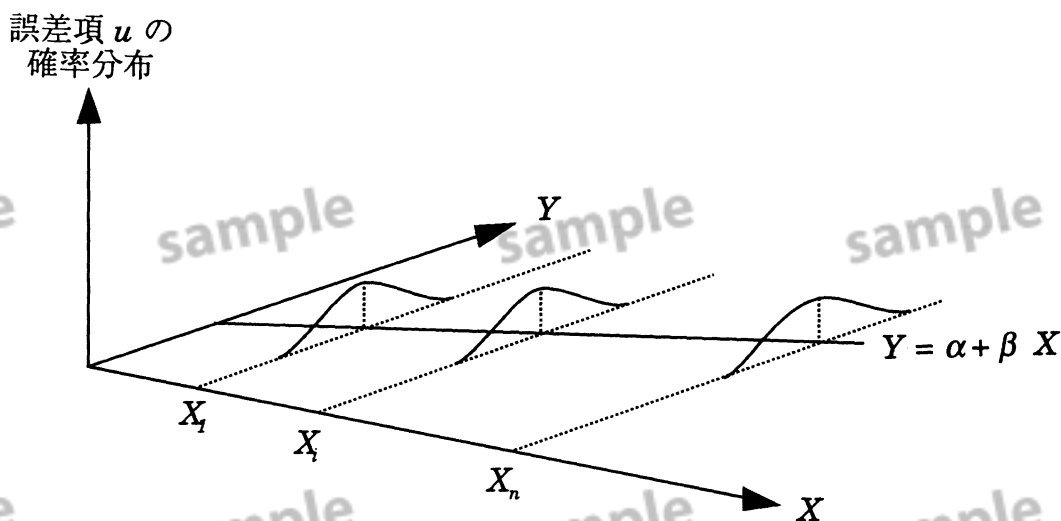
$$u = NID(0, \sigma_u^2)$$

によって示すこともできる。これは誤差項の平均が0で、正規に独立して分布してい

このノートは、慶應義塾大学ビジネス・スクールにおける補助教材として、同ビジネス・スクール教授矢作恒雄と博士課程磯辺剛彦が作成した。(1995年5月作成)

ることを意味する。特に分散の均一性については、図1に描かれているように、いずれの X_i においてもそのバラツキが同じであることを仮定する。

(図1) 分散の均一性



最小二乗法

まず n 個の観測データ、 (X_1, Y_1) 、 (X_2, Y_2) 、...、 (X_n, Y_n) が図2のようにプロットされたとしよう。単純回帰の目的は、この直線を推定することであり、 $\hat{Y} = \hat{\alpha} + \hat{\beta} X$ で表されることにする。ここで、

\hat{Y} : X を固定したときの推定値 Y

$\hat{\alpha}$ 、 $\hat{\beta}$: パラメータの推定値 (これ以降、 $\hat{\alpha} = a$ 、 $\hat{\beta} = b$ とする)

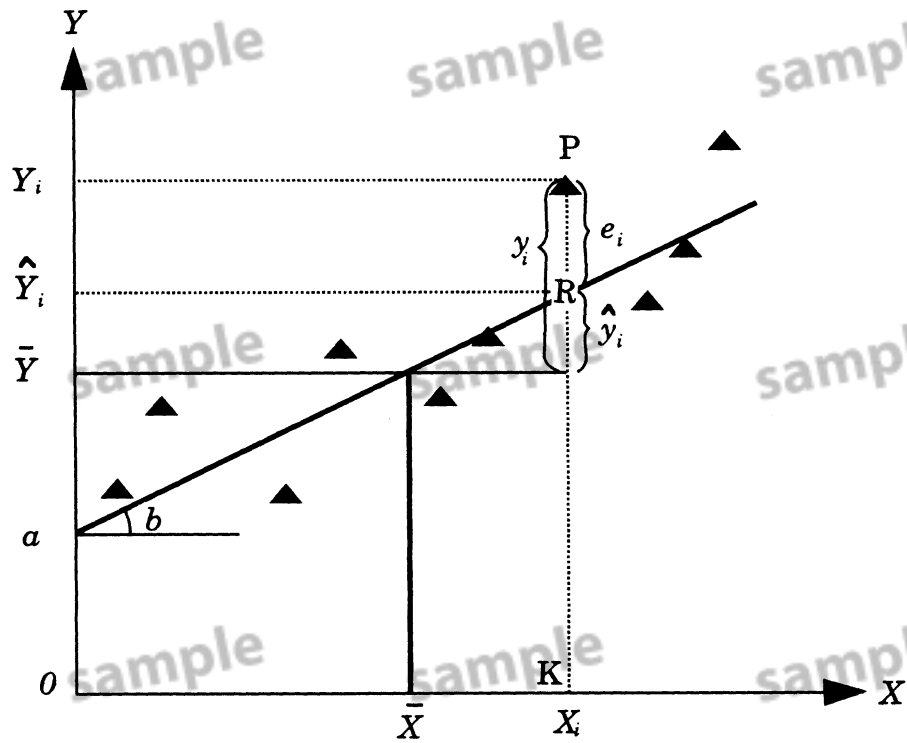
を示している。そして推定線を求めるには、ある公式を用いて a と b を観測データから導き出すことが必要になる。ここで任意の点 $P(X_i, Y_i)$ をとり、その X 軸への垂線を引く。垂線と推定線との交点を $R(X_i, \hat{Y}_i)$ 、 X 軸との交点を $K(X_i, 0)$ とする。このとき $OK = X_i$ 、 $KP = Y_i$ 、 $KR = \hat{Y}_i$ である。また X 軸の垂線に沿った点 P と点 R との差を

$$e_i = Y_i - \hat{Y}_i = Y_i - (a + b X_i) \quad (5)$$

と定義する。この推定線と観測値との差 (残差) はプラスであったりマイナスであったりするが、これらの残差を二乗して和をとれば決してマイナスにはならない。最小二乗法とは、この残差平方和の値を最小にするような a と b を選択するものである。そのための必要条件は、

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - a - b X_i)^2 \quad (6)$$

(図2) 観測データと推定線



の最小化であり、(6)式の残差平方和を a と b について偏微分したものが 0 となるような方程式を解けばよい¹⁾。つまり、

$$\frac{\partial(\Sigma e^2)}{\partial a} = \frac{\partial(\Sigma e^2)}{\partial b} = 0 \quad (7)$$

の条件が必要になる。その結果、

$$\frac{\partial(\Sigma e^2)}{\partial a} = -2 \Sigma (Y - a - b X) = 0 \quad (8)$$

$$\Sigma Y = na + b \Sigma X$$

さらに、

$$\frac{\partial(\Sigma e^2)}{\partial b} = -2 \Sigma X (Y - a - b X) = 0 \quad (9)$$

$$\Sigma XY = a \Sigma X + b \Sigma X^2$$

¹⁾ これ以降 $\sum_{i=1}^n$ を Σ に簡略化し、また各記号から i を省略する。たとえば、 $\sum_{i=1}^n e_i$ は Σe と表す。

という2つの方程式が得られる。これらの方程式は正規方程式と呼ばれる。

このようにして推定された式は、データの平均点 (\bar{X}, \bar{Y}) を通るという特性をもつ。これは(8)式を n で割ることによって証明される¹⁾。

$$\bar{Y} = a + b \bar{X} \quad (10)$$

さらに回帰係数 (regression coefficient) b は、 $x = X - \bar{X}$ 、 $y = Y - \bar{Y}$ という平均からの偏差を用いれば、

$$b = \frac{\sum xy}{\sum x^2} \quad (11)$$

で示すことができる。これは(8)と(9)式によって導き出される。

$$\begin{aligned} \sum XY &= a \sum X + b \sum X^2 \\ &= (\bar{Y} - b \bar{X}) \sum X + b \sum X^2 \\ &= \frac{1}{n} (\sum Y) (\sum X) - \frac{b}{n} (\sum X)^2 + b \sum X^2 \end{aligned} \quad (12)$$

したがって、

$$\sum XY - \frac{1}{n} (\sum Y) (\sum X) = b \left[\sum X^2 - \frac{1}{n} (\sum X)^2 \right] \quad (13)$$

ここで X と Y とをそれぞれの偏差 x と y とに置き換えれば、

$$\begin{aligned} \sum xy &= \sum (X - \bar{X})(Y - \bar{Y}) \\ &= \sum XY - \bar{Y} \sum X - \bar{X} \sum Y + n \bar{X} \bar{Y} \\ &= \sum XY - \frac{1}{n} (\sum X) (\sum Y) \end{aligned} \quad (14)$$

また、

$$\begin{aligned} \sum x^2 &= \sum (X - \bar{X})(X - \bar{X}) \\ &= \sum X^2 - \bar{X} \sum X - \bar{X} \sum X + n \bar{X}^2 \\ &= \sum X^2 - \frac{1}{n} (\sum X)^2 \end{aligned} \quad (15)$$

(14)式と(15)式とを(13)式に代入すれば(11)式が得られる。

また別解として、 $\sum e^2$ を最小化することによってもパラメータ b を求めることが

¹⁾ X と Y の算術平均は、それぞれ $\sum X / n$ と $\sum Y / n$ である。

可能である。(5)式により Σe^2 は次のように変換される。

$$\begin{aligned}\Sigma e^2 &= \Sigma (Y - \hat{Y})^2 \\ &= \Sigma [(Y - \bar{Y}) - (\hat{Y} - \bar{Y})]^2 \\ &= \Sigma (y - \hat{y})^2 \\ &= \Sigma (y - b x)^2\end{aligned}\tag{16}$$

ここで、

$$\hat{y} = b x\tag{17}$$

であることが分かる。なぜなら、

$$\begin{aligned}\hat{y} &= Y - \bar{Y} \\ &= a + b X - (a + b \bar{X}) \\ &= b (X - \bar{X})\end{aligned}\tag{18}$$

Σe^2 を最小化するために、(16)式を b について偏微分を行えば、

$$\frac{\partial(\Sigma e^2)}{\partial b} = 2 \Sigma x (y - b x) = 0$$

したがって、

$$\Sigma x y = b \Sigma x^2$$

また切片のパラメータ a を求めるには、 b を(10)式に代入すればよい。

$$a = \bar{Y} - \left(\frac{\Sigma x y}{\Sigma x^2} \right) \bar{X}\tag{19}$$

平方和の分解

次に、総平方和(TSS)、もしくは全変動と呼ばれる Σy^2 は、説明された平方和($ESS : \Sigma \hat{y}^2$)と残差平方和($RSS : \Sigma e^2$)とに分解することができる。 $y = \hat{y} + e$ は $y = b x + e$ に変換できるので、

$$\begin{aligned}
\Sigma y^2 &= \Sigma (bx + e)^2 \\
&= b^2 \Sigma x^2 + \Sigma e^2 + 2be \Sigma xe \\
&= b^2 \Sigma x^2 + \Sigma e^2
\end{aligned}
\tag{20}$$

したがって、

$$TSS = ESS + RSS \tag{21}$$

ここで、

$$\begin{aligned}
\Sigma xe &= \Sigma (X - \bar{X}) e \\
&= \Sigma Xe - \bar{X} \Sigma e \quad \therefore \Sigma Xe = \Sigma e = 0 \\
&= 0
\end{aligned}
\tag{22}$$

となる。これは (8) 式と (9) 式の偏微分が 0 であることを用いることによって簡単に証明できる。

$$\frac{\partial(\Sigma e^2)}{\partial a} = -2 \Sigma (Y - a - bX) = -2 \Sigma e = 0$$

さらに、

$$\frac{\partial(\Sigma e^2)}{\partial b} = -2 \Sigma X (Y - a - bX) = -2 \Sigma Xe = 0$$

また ESS は (17) 式により b を消去し、 x と y との関数で表すことができる。

$$\begin{aligned}
\Sigma \hat{y}^2 &= \Sigma (bx)^2 \\
&= b^2 \Sigma x^2 \\
&= \left(\Sigma xy / \Sigma x^2 \right)^2 \Sigma x^2 \\
&= (\Sigma xy)^2 / \Sigma x^2
\end{aligned}
\tag{23}$$

したがって総平方和の分解することにより、回帰によって説明できる割合は、

$$\begin{aligned}
\Sigma \hat{y}^2 / \Sigma y^2 &= ESS / TSS = (1 - RSS) / TSS \\
&= (\Sigma xy)^2 / \Sigma x^2 \Sigma y^2 \\
&= \left(\frac{\Sigma xy}{n} \right)^2 / \left(\frac{\Sigma x^2}{n} \cdot \frac{\Sigma y^2}{n} \right)
\end{aligned}
\tag{24}$$

通常これを決定係数と呼ぶ。決定係数は r_{xy}^2 によって表され、0と1の間の数値になる。また r_{xy} は相関係数 (correlation coefficient) と呼ばれる。相関係数とは X と Y との直線関係についての尺度であり、分子は共分散、分母は X と Y の標準偏差の積で示される。

$$r = \frac{\frac{\sum xy}{n}}{\sqrt{\frac{\sum x^2}{n}} \cdot \sqrt{\frac{\sum y^2}{n}}} = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)} \cdot \sqrt{\text{Var}(Y)}} \quad (25)$$

最小二乗法の特徴

最小二乗法が回帰分析にとって望ましい理由は、この方法によって推定されたパラメータが最良線形不偏推定値 (best linear unbiased estimator ; BLUE) になりうるからである。そしてこの手法を、通常最小二乗法 (Ordinary Least Squares : OLS) と呼ぶ。最良線形不偏推定値とは、主に次のような条件を満たすような推定値である。ここでは傾き b を例にすると、

- (A) 推定された b は線形関数であること
- (B) b の平均値が真のパラメータ β と等しいこと
- (C) b の分布である $f(b)$ の分散が他のいかなる推定値の分散よりも小さいこと

そして、この3つの条件を満たす図3の(b)が不偏推定値になる。

次に、(1)から(3)の仮定をおくことによって、未知数である切片 a と傾き b の平均と分散を測定することができる。まず傾き b について(11)式を用いて考察すると、

$$\begin{aligned} b &= \frac{\sum xy}{\sum x^2} \\ &= \frac{\sum x(Y - \bar{Y})}{\sum x^2} \\ &= \frac{\sum xY}{\sum x^2} \\ &= \sum w Y \end{aligned} \quad (26)$$

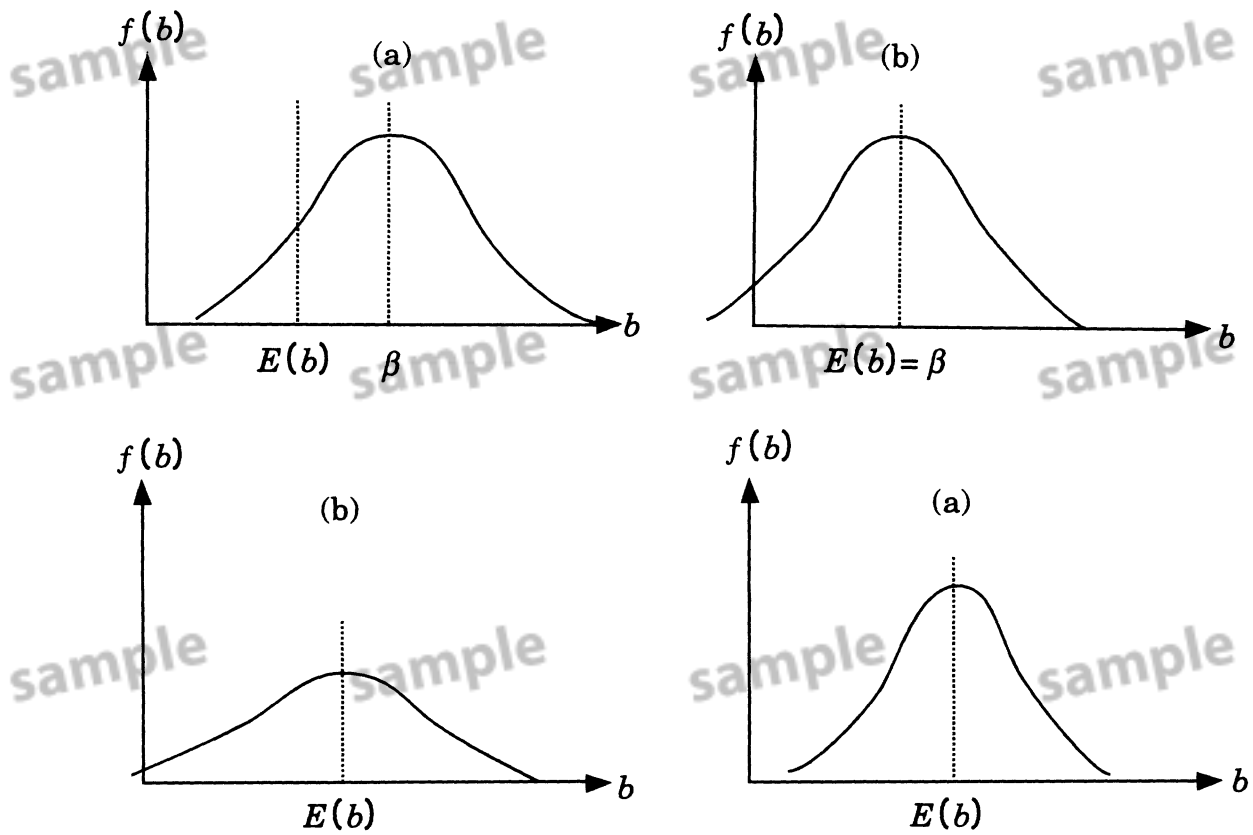
なぜなら、

$$\sum x \bar{Y} = \bar{Y} \sum (X - \bar{X}) = \bar{Y} (\sum X - \sum X) = 0 \quad (27)$$

ただし、

$$w = \frac{x}{\sum x^2} \quad (28)$$

(図3) 不偏推定値の特性



この(28)式から、次のような3つの条件が導き出される。

$$\sum w = 0 \quad \therefore \sum x = \sum (X - \bar{X}) = \sum X - \sum X \quad (29a)$$

$$\begin{aligned} \sum w X = 1 \quad \therefore \sum w x &= \frac{\sum x^2}{\sum x^2} = 1 \\ &= \sum w (X - \bar{X}) = \sum w X \end{aligned} \quad (29b)$$

$$\sum w^2 = \frac{\sum x^2}{(\sum x^2)^2} = 1 / \sum x^2 \quad (29c)$$

これらの条件と(1)式を組み合わせることによって傾き b は、

$$\begin{aligned} b &= \sum w Y = \sum w (\alpha + \beta X + u) \\ &= \beta + \sum w u \end{aligned} \quad (30)$$

そして(2)式の $E(u) = 0$ の仮定を用いると、この(30)式の期待値は、

$$E(b) = \beta + E(\sum w u) = \beta \quad (31)$$

つまり最小二乗法によって求められた傾き b は、真の傾き β の不偏推定値であることがわかる。同様に、傾き b の分散 (Var) については、

$$Var(b) = E[(b - \beta)^2] = E[(\sum w u)^2] \quad (32)$$

ここで、

$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$$

とすると、

$$\mathbf{w}'\mathbf{u} = w_1 u_1 + w_2 u_2 + \dots + w_n u_n$$

したがって、

$$\begin{aligned} (\sum w u)^2 &= w_1^2 u_1^2 + \dots + w_n^2 u_n^2 + 2w_1 u_1 w_2 u_2 + \dots + 2w_{n-1} u_{n-1} w_n u_n \\ &= \sum w^2 u^2 + 2 \sum_{i < j} w_i u_i w_j u_j \end{aligned}$$

ここで、(3) 式の仮定である $E(u^2) = \sigma_u^2$ と $E(u_i u_j) = 0$ 、さらに (29c) 式を用いれば、

$$Var(b) = E(\sum w u)^2 = \sigma_u^2 \sum w^2 = \sigma_u^2 / \sum x^2 \quad (33)$$

また切片 a については、(10) 式に (26) 式を代入することによって求めることができる。

$$\begin{aligned} a &= \bar{Y} - b \bar{X} \\ &= \bar{Y} - \bar{X} \sum w Y \quad \therefore (26) \text{ 式} \\ &= \sum \left(\frac{1}{n} - w \bar{X} \right) Y \\ &= \sum \left(\frac{1}{n} - w \bar{X} \right) (\alpha + \beta X + u) \\ &= \alpha + \sum \left(\frac{1}{n} - w \bar{X} \right) u \end{aligned} \quad (34)$$

(34) 式の期待値をとると、

$$E(a) = \alpha \quad \therefore E(u) = 0 \quad (35)$$

(33) 式を用いると切片 a の分散は、

$$\begin{aligned} \text{Var}(a) &= E[(a - \alpha)^2] \\ &= \Sigma \left(\frac{1}{n} - w \bar{X} \right)^2 u^2 \\ &= \sigma_u^2 \left(\frac{1}{n} + \bar{X}^2 \Sigma w^2 - 2 \bar{X} \Sigma w \right) \\ &= \sigma_u^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\Sigma x^2} \right) \end{aligned} \quad (36)$$

また傾き b と切片 a の共分散は、

$$\begin{aligned} \text{cov}(a, b) &= E[(a - \alpha)(b - \beta)] \\ &= E\{[\bar{u} - (b - \beta)\bar{X}](b - \beta)\} \\ &= -\bar{X} \cdot E[(b - \beta)^2] \\ &= -\bar{X} \sigma_u^2 / \Sigma x^2 \quad \therefore E(\bar{u}) = 0 \end{aligned} \quad (37)$$

以上のような最小二乗法による分散最小化は、仮定(1)式から(4)式に基づいたものである。仮定(1)式は、 a と b が u の線形関数であり、仮定(4)式は u が正規分布であることを示している。したがって a と b の平均と分散とが既知であれば、 a と b は次のような正規分布にしたがうことになる。

$$a \sim N\left(\alpha, \sigma_u^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\Sigma x^2} \right]\right) \quad (38)$$

$$b \sim N\left(\beta, \frac{\sigma_u^2}{\Sigma x^2}\right) \quad (39)$$

しかし u の分散である σ_u^2 は未知であるため、回帰分析を操作可能なものにするためには σ_u^2 の推定が必要になる。この推定には(5)式を用いる。

$$\begin{aligned} e &= Y - \hat{Y} \\ &= (\alpha + \beta X + u) - (a + b X) \\ &= u - (a - \alpha) - (b - \beta) X \end{aligned} \quad (40)$$

さらに、 $\bar{Y} = a + b \bar{X}$ より、

$$\begin{aligned} a &= (\alpha + \beta \bar{X} + u) - b \bar{X} \\ &= \alpha - (b - \beta) \bar{X} + \bar{u} \end{aligned} \quad (41)$$

したがって、

$$e = (u - \bar{u}) - (b - \beta) x \quad (42)$$

ただし $\bar{u} \neq 0$ である。(42)式の平方和は、

$$\begin{aligned} \sum e^2 &= \sum (u - \bar{u})^2 + (b - \beta)^2 \sum x^2 - 2(b - \beta) \sum (u - \bar{u}) x \\ &= \sum u^2 - n \bar{u}^2 + (b - \beta)^2 \sum x^2 - 2(b - \beta) \sum u x \end{aligned} \quad (43)$$

残差平方和である(43)式の期待値をとると、

$$\begin{aligned} \cdot E(\sum u^2) &= n \sigma_u^2 \\ \cdot E(\bar{u}^2) &= \text{Var}(\bar{u}) = \sigma_u^2 / n \\ E[(b - \beta)^2] &= \text{Var}(b) = \sigma_u^2 / \sum x^2 \\ \cdot E[(b - \beta)(\sum u x)] &= E[(\sum w u)(\sum u x)] \\ &= \sum u x / \sum x^2 \cdot \sum u x \\ &= \sigma_u^2 \end{aligned}$$

より、

$$\begin{aligned} E(\sum e^2) &= n \sigma_u^2 - \sigma_u^2 + \sigma_u^2 - 2 \sigma_u^2 \\ &= (n - 2) \sigma_u^2 \end{aligned} \quad (44)$$

つまり単回帰における誤差項の分散 σ_u^2 は、残差平方和をサンプル数から2を引いた数値によって割ったものになる。

$$E\left(\frac{\sum e^2}{n - 2}\right) = \sigma_u^2 \quad (45)$$

最小二乗法による推定

これまで見たように、仮定(1)式から(4)式のもとで、切片 a と傾き b の分布とを導き出してきた。さらに(38)式と(39)式とに含まれている未知の分散 σ_u^2 についても、その推定方法は既に述べたとおりである。次に考慮しなければならないのは、これらの仮定や帰結のもとで、どのように区間推定や仮説検定を行うべきかについてである。この推定・検定のプロセスについて解説する前に、さまざまな前提条件や命題が存在するため、これらについて一通り述べることにする。

定理1：確率変数 x の平均が μ であり、分散が σ^2 の正規分布であることが知られているとき、 $z = \frac{x - \mu}{\sigma}$ は平均 0、分散 1 の正規分布にしたがう。

$$z = \frac{x - \mu}{\sigma} \sim N(0, 1)$$

証明： x が確率変数であるとき、 $z = \frac{x - \mu}{\sigma}$ もまた確率変数である。ここで $f(x)$ を x の確率密度関数とすると、 z の平均は、

$$\begin{aligned}\mu_z &= \int_{-\infty}^{+\infty} \frac{x - \mu}{\sigma} f(x) dx \\ &= \frac{1}{\sigma} \left[\int_{-\infty}^{+\infty} x f(x) dx - \mu \int_{-\infty}^{+\infty} f(x) dx \right] \\ &= \frac{1}{\sigma} (\mu - \mu) = 0\end{aligned}$$

また分散については、

$$\begin{aligned}\sigma_z^2 &= E[(z - \mu_z)^2] \\ &= \int_{-\infty}^{+\infty} \left(\frac{x - \mu}{\sigma}\right)^2 f(x) dx - \mu_z^2 \\ &= \frac{1}{\sigma^2} \left[\int_{-\infty}^{+\infty} x^2 f(x) dx - 2\mu \int_{-\infty}^{+\infty} x f(x) dx + \mu^2 \int_{-\infty}^{+\infty} f(x) dx \right] \quad \because \mu_z = 0 \\ &= \frac{1}{\sigma^2} [(\sigma^2 + \mu^2) - 2\mu \cdot \mu + \mu^2] = 1\end{aligned}$$

定理2： z_1, z_2, \dots, z_n が互いに独立して標準正規分布 $N(0, 1)$ にしたがう確率変数ならば、 $x = z_1^2 + z_2^2 + \dots + z_n^2$ は自由度 n のカイ二乗分布にしたがう。

定理1により $z = \frac{x - \mu}{\sigma}$ は標準正規分布にしたがうことが分かっているので、
 (45)式から $\sum e^2 / \sigma_u^2$ は自由度 $(n - 2)$ のカイ二乗分布にしたがう。

$$\frac{e_1^2}{\sigma_u^2} + \frac{e_2^2}{\sigma_u^2} + \dots + \frac{e_n^2}{\sigma_u^2} = \frac{\sum e^2}{\sigma_u^2} \sim \chi(n-2)$$

定理3: z は標準正規分布にしたがう確率変数、 v は自由度 ϕ のカイ二乗分布にしたがう確率変数であり、さらに z と v とが互いに独立であるとき、 $\frac{z\sqrt{\phi}}{\sqrt{v}}$ は自由度 ϕ の t 分布にしたがう。

傾き b について、平均と分散はそれぞれ β と $\frac{\sigma_u^2}{\sum x^2}$ であることから、定理1により、

$$\frac{(b - \beta)}{\sigma_u / \sqrt{\sum x^2}} \sim N(0, 1)$$

また定理3より、未知の σ_u^2 を不偏推定値である s^2 に置き換えることによって正規分布から t 分布へと変わる（ただしサンプル数が増えるにつれて t 分布は正規分布へと近づく）。分母である $\frac{s}{\sqrt{\sum x^2}}$ は分散の平方根であり、一般に標準誤差と呼ばれる。

$$\begin{aligned} t = z \sqrt{\phi} / \sqrt{v} &= \frac{(b - \beta)}{\sigma_u / \sqrt{\sum x^2}} \cdot \frac{\sqrt{\sum e^2} / \sigma_u}{\sqrt{(n - 2)}} \\ &= \frac{(b - \beta)}{s / \sqrt{\sum x^2}} \sim t(n - 2) \end{aligned} \tag{46}$$

一般的な検定手続きは (46) 式によって行うことが可能である。まず傾き b の 95% の信頼区間は、

$$b \pm t_{0.025} \frac{s}{\sqrt{\sum x^2}} \tag{47}$$

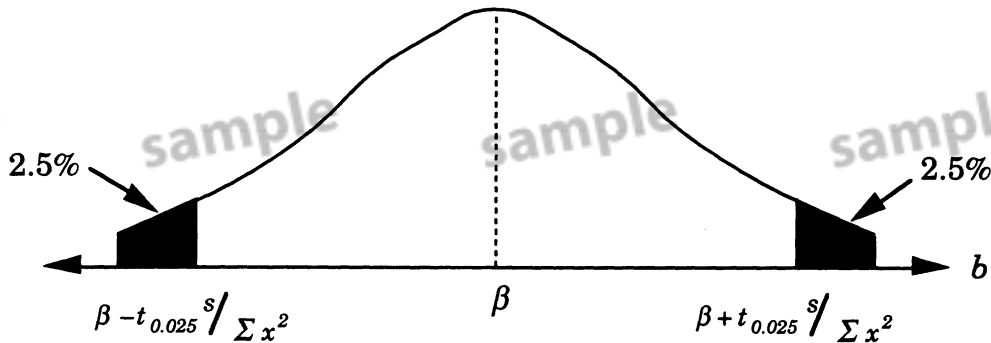
であり、図4のように描くことができる。この信頼区間とは、推定された傾き b が真の傾きである β と 95% の確率で等しい範囲を意味する。このように傾き b が β と等しいという仮説は帰無仮説 (null hypothesis) と呼ばれ、次のように記述される。

$$H_0: b = \beta$$

逆に b が β とは異なるという仮説は対立仮説と呼ばれる。

$$H_1: b \neq \beta$$

(図4) 信頼区間



この信頼区間は、

$$Prob \left[-t_{0.025} < \frac{(b - \beta)}{\sigma_u / \sqrt{\sum x^2}} < t_{0.025} \right] = 0.95$$

あるいは、

$$\left| \frac{(b - \beta)}{\sigma_u / \sqrt{\sum x^2}} \right| > t_{0.025}$$

と記述される。

多くの場合、帰無仮説は $\beta = 0$ であり、これを X の有意性の検定という。もし帰無仮説が有意であれば、 X は Y の決定において何ら意味を持たないことになる。逆に帰無仮説が棄却された場合には、5% の水準で有意であるという。つまり Y の決定において、 X の役割がゼロではないということを示すのものである。また有意性の検定には、5% や 1% といった水準が採用されるが、この数値自体には何ら根拠はなく、単に広く採用されているだけであることを認識しておくべきである。

次に切片 a についても、検定手続は傾き b と同様である。切片 a の分布は、

$$a \sim N \left(\alpha, \sigma_u^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2} \right] \right)$$

であることが分かっている。したがって、

$$t = \frac{a - \alpha}{s \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2}}} \sim t(n-2) \quad (48)$$

また 95% の信頼区間と帰無仮説は、それぞれ

$$a \pm t_{0.025} \cdot s \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2}}$$

$$H_0: a = \alpha$$

最小二乗法における分散分析

単回帰モデルに関連して、分散分析 (analysis of variance : ANOVA) が取り上げられることがある。これは (21) 式のように、全変動 (総平方和) が説明された平方和と残差平方和とに分解できる特性を利用するものである。表 1 は単回帰における総平方和分解したものである。

(表 1) 単回帰における分散分析

変動の原因	平方和	自由度	平均平方和
X	$ESS = \sum \hat{y}^2 = b^2 \sum x^2$ $= b \sum xy$	1	$ESS/1$
残差	$RSS = \sum e^2$	$n-2$	$RSS/(n-2)$
合計	$TSS = \sum y^2$	$n-1$	

(46) 式により

$$\frac{(b - \beta)}{\sigma_u / \sqrt{\sum x^2}} \sim N(0, 1)$$

したがってこれを二乗したものは、定理 2 により自由度 1 のカイ二乗分布にしたがう。

$$\frac{(b - \beta)^2}{\sigma_u^2 / \sum x^2} \sim \chi(1) \quad (49)$$

定理4： v_1, v_2 がそれぞれ自由度 ϕ_1, ϕ_2 のカイ二乗分布にしたがう独立した確率変数であるとき、 $g = \frac{v_1 / \phi_1}{v_2 / \phi_2}$ は自由度 (ϕ_1, ϕ_2) の F 分布にしたがう。

定理4により、

$$F = \frac{v_1 / \phi_1}{v_2 / \phi_2} = \frac{(b - \beta)^2 \Sigma x^2 / 1}{\Sigma e^2 / (n - 2)} \sim F(1, n - 2) \quad (50)$$

ここでもし帰無仮説が $\beta = 0$ であるならば、

$$\begin{aligned} F &= \frac{b^2 \Sigma x^2}{\Sigma e^2 / (n - 2)} = \frac{b^2 \Sigma x^2}{s^2} \\ &= \frac{ESS}{RSS / (n - 2)} \sim F(1, n - 2) \end{aligned} \quad (51)$$

この結果を t 検定の結果と比較すれば、

$$t = \frac{(b - \beta) \sqrt{\Sigma x^2}}{s} \sim t(n - 2)$$

であるため、単回帰モデルにおける F 統計量は t 統計量の二乗に等しくなる。

定理5： h が自由度 ϕ の t 分布にしたがう確率変数であるとき、 h^2 は自由度 $(1, \phi)$ の F 分布にしたがう。

ここで単回帰における決定係数 R^2 と t 値、 F 値との関係について簡単に述べておく。まず、 $R^2 = ESS / TSS$ および $TSS = ESS + RSS$ より、

$$\frac{ESS}{RSS} = \frac{R^2}{1 - R^2}$$

また、

$$F = \frac{ESS / 1}{RSS / (n - 2)} = \frac{R^2}{(1 - R^2) / (n - 2)}$$

であることから、

$$R^2 = \frac{F}{F + (n - 2)}$$

定理5により単回帰の場合には $t^2 = F$ なので、

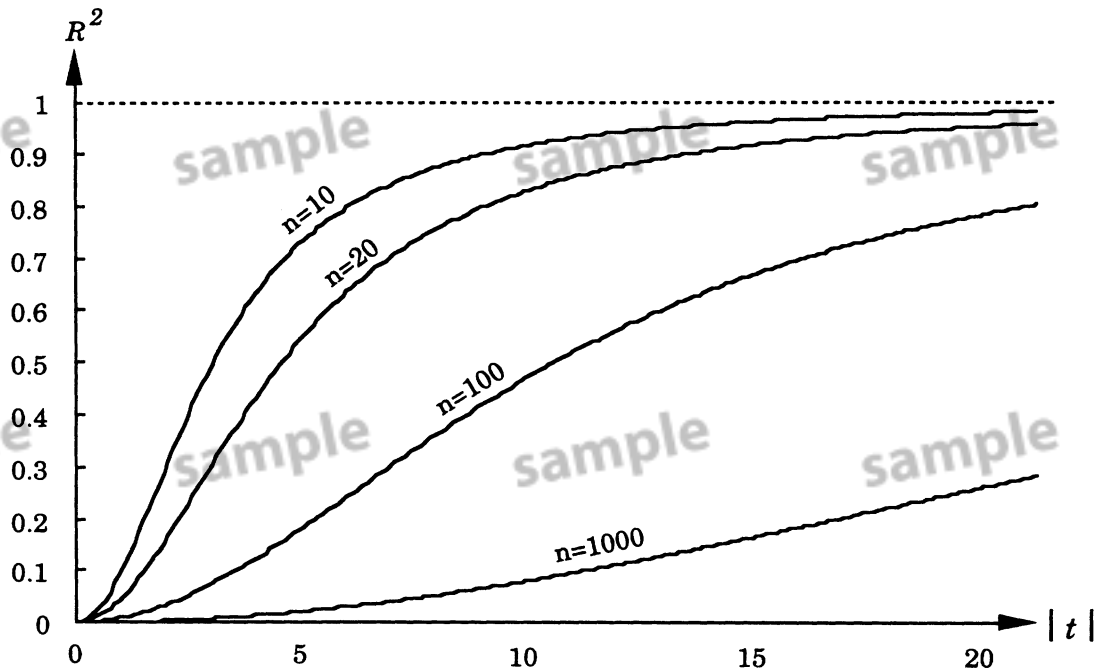
$$R^2 = \frac{t^2}{t^2 + (n - 2)}$$

ここでサンプル数ごとに決定係数 R^2 と t 値との関係を描いたのが図5である。この図によると、サンプル数が少ないほど、小さい t 値でも高い決定係数が得られることを示している。また両者の関係式を微分すると、

$$\frac{d(R^2)}{dt} = \frac{2(n-2)t}{[t^2 + (n-2)]^2} \quad \text{さらに} \quad \frac{d^2(R^2)}{dt^2} = \frac{2(n-2)[(n-2) - 3t^2]}{[t^2 + (n-2)]^3} \quad \text{より、}$$

$|t| < \sqrt{\frac{n-2}{3}}$ のとき R^2 は単純増加曲線、逆に $|t| > \sqrt{\frac{n-2}{3}}$ のとき R^2 は単純減少曲線であることも指摘される。

(図5) 決定係数 R^2 と t 値との関係



最小二乗法における予測

これまで最小二乗法による2変数データの分析手法を考察してきた。次のステップは、新たに説明変数が観察された場合に、導き出された回帰モデルがどの程度の精度をもつのか？、また X_0 を代入することによって計算された被説明変数はどのくらいの信頼区間をもつのか？、などに関連するものである。ここで任意の値(X_0)についての回帰式のパredictionを(\hat{Y}_0)とすると、

$$\hat{Y}_0 = a + b X_0 \quad (52)$$

また、

$$Y_0 = \alpha + \beta X_0 + u_0 \quad (53)$$

によって予測誤差(e_0)は、

$$\begin{aligned} e_0 &= Y_0 - \hat{Y}_0 \\ &= u_0 - (a - \alpha) - (b - \beta) X_0 \end{aligned} \quad (54)$$

また a と b はそれぞれ α と β の不偏推定値であるため、 $E(a - \alpha) = 0$ 、 $E(b - \beta) = 0$ である。さらに $E(u_0) = 0$ により、予測誤差の期待値は $E(e_0) = 0$ である。これは(45)式による予測値が不偏であることを意味する。また予測誤差の分散については、

$$\begin{aligned} \text{Var}(e_0) &= E(e_0^2) \\ &= E[u_0^2 + (a - \alpha)^2 + (b - \beta)^2 X_0^2 - 2(a - \alpha)u_0 \\ &\quad - 2(b - \beta)X_0 u_0 - 2(a - \alpha)(b - \beta)X_0] \\ &= \text{Var}(u_0^2) + \text{Var}(a) + X_0^2 \text{Var}(b) - 2X_0 \text{cov}(a, b) \\ &= \sigma_u^2 + \sigma_u^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2} \right] + X_0^2 \frac{\sigma_u^2}{\sum x^2} - 2X_0 \sigma_u^2 \frac{\bar{X}}{\sum x^2} \\ &= \sigma_u^2 \left[1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2} \right] \end{aligned} \quad (55)$$

つまり X_0 が \bar{X} に等しいとき、誤差の分散 $\sigma_u^2 \left[1 + \frac{1}{n} \right]$ は最小になる。また(54)式より e_0 は線形関数であり、さらに正規分布にしたがうので、

$$z = \frac{e_0}{\sigma} = \frac{e_0}{\sigma_u \sqrt{1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2}}} \sim N(0, 1) \quad (56)$$

ここで未知である σ_u^2 を推定値 $s^2 = \sum e^2 / (n-2)$ に置き換えることによって、 e_0 は自由度が $n-2$ の t 分布にしたがうことになる。

$$t = \frac{Y_0 - \hat{Y}_0}{s \sqrt{1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2}}} \sim t(n-2) \quad (57)$$

以上により、未知の数値である Y_0 の 95% 水準の信頼区間は、

$$(a + b X_0) \pm t_{0.025} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2}} \quad (58)$$

次に、もし Y_0 自体の数値よりも Y_0 の平均について興味があるのであれば、(53) 式の期待値をとればよい。

$$E(Y_0) = \alpha + \beta X_0 \quad (59)$$

したがって、

$$\begin{aligned} e_0 &= E(Y_0) - \hat{Y}_0 \\ &= -(a - \alpha) - (b - \beta) X_0 \end{aligned} \quad (60)$$

この分散は、

$$\begin{aligned} \text{Var}(e_0) &= E[(a - \alpha)^2 + X_0^2 (b - \beta)^2 - 2(a - \alpha)(b - \beta) X_0] \\ &= \text{Var}(a) + X_0^2 \text{Var}(b) - 2 X_0 \text{cov}(a, b) \\ &= \sigma_u^2 \left[\frac{1}{n} + \frac{(\bar{X} - X_0)^2}{\sum x^2} \right] \end{aligned} \quad (61)$$

つまり Y_0 の平均である $E(Y_0)$ の 95% 水準の信頼区間は、

$$(a + b X_0) \pm t_{0.025} \cdot s \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2}} \quad (62)$$

によって求めることができる。

例題

これまでの単回帰の仮定や手法を数学的に考察してきたが、例題を用いて実際に手計算することによって最小二乗法の仕組みを理解することが望ましい。表2は、ある機械の稼働時間(X)と固定費(Y)とを30ヶ月(n)にわたって測定したものである。ここでの興味は、固定費を決定する上で機械の稼働時間がどのような影響を及ぼすのか?、またそれは確率的に有意なのか?、などに関連するものである。

(表2) 機械の稼働時間と固定費

月数 (i)	固定費 (X)	機械の稼働時間 (Y)
1	76.67	1.77
2	73.68	1.82
3	80.14	1.63
4	61.99	1.01
5	72.69	1.38
6	87.68	1.96
7	78.45	1.56
8	70.63	1.46
9	63.42	1.55
10	56.06	1.12
11	67.45	1.38
12	72.10	1.32
13	68.53	1.26
14	69.08	1.34
15	85.55	1.80
16	58.20	1.02
17	61.63	1.51
18	80.69	1.79
19	58.26	1.15
20	55.34	1.16
21	85.11	1.85
22	76.49	1.83
23	67.78	1.14
24	56.40	1.14
25	66.62	1.33
26	63.49	1.36
27	89.42	1.88
28	67.52	1.17
29	73.68	1.64
30	66.13	1.38

(Kaplan, R.S & A.A. Atkinson "Advanced Management Accounting", Prentice-Hall, N.Y: 1989; pp.99.を修正)

$$\Sigma X = 43.71 \quad \Sigma Y = 2,110.88 \quad \Sigma XY = 3,144.157$$

$$\Sigma X^2 = 66.039$$

これらを(8)式と(9)式に代入すれば、

$$2,110.88 = 30a + b(43.71)$$

$$3,144.157 = a(43.71) + b(66.039)$$

この方程式を解くと、 $a = 27.891$ 、 $b = 29.150$ が求められる。この結果は、機械を稼働させなくとも固定費が 27.891 発生し、また稼働時間が 1 時間延びるにつき固定費が 29.150 発生するというものである。別解としては、(11) 式に

$$\sum x^2 = 2.354 \quad \sum xy = 68.605$$

を代入すれば切片 b が導かれる。

$$b = \frac{\sum xy}{\sum x^2} = \frac{68.605}{2.354} = 29.144 \quad (\text{A})$$

さらに切片 a は (10) 式により求めることができる。

$$a = \bar{Y} - b \bar{X} = 70.363 - (29.144)(1.457) = 27.900 \quad (\text{B})$$

$$\text{ここで、} \bar{X} = \frac{\sum X}{n} = 1.457 \quad \bar{Y} = \frac{\sum Y}{n} = 70.363$$

次に、この回帰モデルの総平方和 (TSS) を説明された平方和 (ESS) と残差平方和 (RSS) とに分解すると、

$$\text{TSS} = \sum y^2 = 2,708.273$$

$$\text{ESS} = \sum \hat{y}^2 = \frac{(\sum xy)^2}{\sum x^2} = \frac{(68.605)^2}{2.354} = 1,999.425 \quad (\text{C})$$

$$\text{RSS} = \text{TSS} - \text{ESS} = 2,708.273 - 1,999.425 = 708.849 \quad (\text{D})$$

残差平方和 (RSS) を検算する方法としては、サンプルごとに被説明変数の推定値を測定し、その平方和 $[\sum (Y - a - bX)^2]$ を計測すればよい。このように測定された 3 つの平方和により、決定係数と相関係数とが導き出される。まず (24) 式によって決定係数は、

$$r_{xy}^2 = \frac{\text{ESS}}{\text{TSS}} = \frac{1,999.425}{2,708.273} = 0.738 \quad (\text{E})$$

これは、機械の稼働時間によって固定費の全変動の 73.8% が説明されることを意味する。また機械の稼働時間と固定費との直線関係の尺度である相関係数は、

$$r_{xy} = \sqrt{0.738} = 0.859 \quad (\text{F})$$

であり、両変数の関係が究めて直線に近いことを示している。

以上の分析によって、切片 a と傾き b の平均値を測定したが、次の段階はこれらの分散についてである。まず (38) 式により切片 a の分散は、

$$\text{Var}(a) = \sigma_u^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2} \right) = \frac{\sum e^2}{n-2} \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2} \right) = \frac{708.849}{28} \left[\frac{1}{30} + \frac{(1.457)^2}{2.354} \right] = 23.674$$

また(39)式により傾き b の分散は、

$$\text{Var}(b) = \frac{\sigma_u^2}{\sum x^2} = \frac{\sum e^2}{n-2} \cdot \frac{1}{\sum x^2} = \frac{708.849}{28} \cdot \frac{1}{2.354} = 10.754$$

つまり切片 a と傾き b は、それぞれ次のような平均値と分散をもつ正規分布にしたがう。

$$a = N(27.900, 23.674) \quad b = N(29.144, 10.754)$$

またこれらの 95% の信頼区間は (34) 式と (35) 式によって求めることができる。付属資料 1 のスチューデントの t 分布表によると、自由度が 28 のときの $t_{0.025}$ は 2.048 である。したがってこの t 値を (34) 式と (35) 式とに代入すると、

$$\begin{aligned} b \pm t_{0.025} \frac{s}{\sqrt{\sum x^2}} &= b \pm t_{0.025} \sqrt{\text{Var}(b)} = 29.144 \pm 2.048 \sqrt{10.754} \\ &= [22.428, 35.860] \end{aligned}$$

$$\begin{aligned} a \pm t_{0.025} \frac{s}{\sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2}}} &= a \pm t_{0.025} \sqrt{\text{Var}(a)} = 27.900 \pm 2.048 \sqrt{23.674} \\ &= [17.935, 37.864] \end{aligned}$$

これは 95% の確率で、真の傾きがこの範囲に存在するというものである。しかし信頼期間は広く、意思決定には何ら有用ではない。

次に切片と傾きの t 値を計算すると、それぞれ

$$t = \frac{b}{s / \sqrt{\sum x^2}} = \frac{b}{\sqrt{\text{Var}(b)}} = \frac{29.144}{\sqrt{10.754}} = 8.887 \quad (\text{G})$$

$$t = \frac{a}{s / \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum x^2}}} = \frac{a}{\sqrt{\text{Var}(a)}} = \frac{27.900}{\sqrt{23.674}} = 5.734 \quad (\text{H})$$

この結果から理解されるように、切片 a と傾き b はともに $t_{0.025}$ である 2.048、さらに $t_{0.005}$ である 2.763 より大きい。したがって、回帰モデルによって導き出された切片と傾きがゼロであるという帰無仮説は、99% の水準で棄却されることになる。

つまり固定費の決定に機械の稼働率はプラスの影響を与えていることが否定できないことが検証されたのである。次に、分散分析による F 検定によってこの結果を確かめてみる。(51)式により、傾き b は自由度(1, 28)の F 値にしたがうことから、

$$F = \frac{ESS}{RSS/n-2} = \frac{1,999.425}{708.849/30-2} = 78.979 \sim F(1, 28) \quad (I)$$

この F 値の平方根は t 値と同じ8.887である。

次に、この回帰モデルを用いて予測を行う場合の留意点について述べることにする。もし経営者や固定費の管理責任者が、ある特定の稼働時間(X_0)がどの程度の固定費(Y_0)を発生させるかについて予測を行うものと仮定しよう。この場合の予測とは、区間予測と点(平均)予測の2つである。ここで想定される稼働時間(X_0)を1.50と2.00時間とする。前者は標本内予測であるが、後者の水準は観察データを越えたものであり、これを標本外予測と呼ぶ。まず稼働時間(X_0)が1.50の場合についての予測を考えてみると、予測誤差の分散は(55)式によって求めることができる。

$$\text{Var}(e_0) = \sigma_u^2 \left[1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2} \right] = \frac{708.849}{28} \left[1 + \frac{1}{30} + \frac{(1.50 - 1.457)^2}{2.354} \right] = 26.180$$

この分散の数値を(58)式に代入すると、次のような信頼区間が得られる。

$$\begin{aligned} (a + b X_0) \pm t_{0.025} \sqrt{\text{Var}(e_0)} &= [27.900 + (29.144)(1.50)] \pm 2.048 \sqrt{26.180} \\ &= [61.137, 82.095] \end{aligned}$$

また稼働時間(X_0)が2.00の場合については、[75.709, 96.667]の範囲が95%の信頼区間となる。このように稼働時間の平均に対して、より乖離が大きい2.00時間の予測の方が信頼区間の範囲が広がる。しかし予測者、もしくは分析者の関心は、信頼区間の計測よりも予測点や予測平均に向けられることが多い。この平均の予測についての95%の信頼区間は、(62)式により

$$(a + b X_0) \pm t_{0.025} s \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x^2}} = [69.711, 73.521]$$

および、 $[82.080, 90.296]$ である。この場合にも、平均より乖離するほどその範囲は広くなり、意思決定に際して不確定要素が大きくなることを意味している。

表3は、この例題の統計ソフトによる回帰分析の結果である。使用したパッケージは SYSTAT バージョン5.2.1 であり、この例題における結果を記号 (A)、(B)、(C) … (I) で示している。

(表3) SYSTAT バージョン 5.2.1 による分析結果

DEP VAR: Overhead		N: 30
MULTIPLE R: 0.859		F
SQUARED MULTIPLE R: 0.738		E
ADJUSTED SQUARED MULTIPLE R: 0.729		
STANDARD ERROR OF ESTIMATE: 5.032		

VARIABLE	COEFFICIENT	STD ERROR	STD COEF	T	P (2 TAIL)
CONSTANT	27.900	4.866	0.000	5.734	.38E-05
Machine-Hours	29.144	3.279	0.859	8.887	.12E-08

ANALYSIS OF VARIANCE					
SOURCE	SUM-OF-SQUARES	DF	MEAN-SQUARE	F-RATIO	P
REGRESSION	1999.413	1	1999.413	78.977	.121767E-08
RESIDUAL	708.859	28	25.316		

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

sample

不 許 複 製

慶應義塾大学ビジネス・スクール

Contents Works Inc.