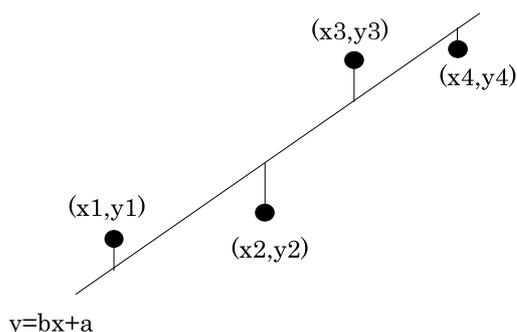


回帰直線の計算法：検量線を例にして

回帰直線は、2つの変数間の関係を散布図として描いたときに、両者の間に直線的な関係があるようにみえることを前提として、縦軸の変数 (y) のばらつきを、横軸の変数 (x) のばらつきによって説明しようとする考え方である。その意味で、縦軸の変数を目的変数（あるいは従属変数）、横軸の変数を説明変数（あるいは独立変数）と呼ぶ。なるべく偏りがないように、測定点にみられる直線的な関係をうまく表すような直線の式を推定するため、最小二乗法を用いるのが普通である。



図のような測定点 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ が得られたときに、検量線 $y = bx + a$ を推定するには、図に示した線分の二乗和が最小になるように a と b を設定すればよい、というのが最小二乗法の考え方である。つまり、

$$\begin{aligned} f(a, b) &= \sum_{i=1}^n \{y_i - (bx_i + a)\}^2 \\ &= b^2 \sum_{i=1}^n x_i^2 - 2b \sum_{i=1}^n x_i y_i + 2ab \sum_{i=1}^n x_i - 2a \sum_{i=1}^n y_i + na^2 + \sum_{i=1}^n y_i^2 \end{aligned}$$

を最小にするような a と b を推定すればよい。

通常、 a と b で偏微分した値がそれぞれ 0 となることを利用して計算すると簡単である。つまり、

$$\frac{\partial f(a, b)}{\partial a} = 2na + 2(b \sum_{i=1}^n x_i - \sum_{i=1}^n y_i) = 0$$

$$i.e. \quad na = \sum_{i=1}^n y_i - b \sum_{i=1}^n x_i$$

$$\frac{\partial f(a, b)}{\partial b} = 2b \sum_{i=1}^n x_i^2 + 2(a \sum_{i=1}^n x_i - \sum_{i=1}^n x_i y_i) = 0$$

$$i.e. \quad b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i$$

を連立方程式として a と b について解けばよい。これを解くと、

$$b = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

が得られ、これを上の式に代入すれば a も得られる。検量線に限らず、一般の回帰直線でも、計算方法は原則として同じである。名称の説明をしておくと、一般に、 $y = bx + a$ という回帰直線について、 b を回帰係数 (regression coefficient)、 a を切片 (intercept) と呼ぶ。

データから得た回帰直線は、 $pV = nRT$ のような物理法則と違って、完璧にデータに乗ることはない。そこで、回帰直線の当てはまりのよさを評価する必要が出てくる。いま、 $z_i = a + bx_i$ とおいたときに、 $e_i = y_i - z_i$ を残差 (residual) と呼ぶ。残差は、 y_i のばらつきのうち、回帰直線では説明できなかった残りに該当する。つまり、残差が大きいくほど、回帰直線の当てはまりは悪いと考えられる。残差にはプラスもマイナスもあるので、二乗和をとって、

$$Q = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - z_i)^2 = \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i\right)^2/n - \frac{\left(n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i\right)^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} / n$$

が、回帰直線の当てはまりの悪さを示す尺度となる。この Q を「残差平方和」と呼び、それを n で割った Q/n を残差分散という。この残差分散 $var(e)$ と Y の分散 $var(Y)$ とピアソンの相関係数 r の間には、 $var(e) = var(Y)(1 - r^2)$ という関係が常に成り立つので、 $r^2 = 1 - var(e)/var(Y)$ となる。このことから r^2 が 1 に近いほど回帰直線の当てはまりがよいことになる。その意味で、 r^2 を「決定係数」と呼ぶ。また、決定係数は、 Y のばらつきがどの程度 X のばらつきによって説明されるかを意味するので、 X の「寄与率」と呼ぶこともある。

回帰直線は最小二乗法でもっとも残差平方和が小さくなるように選ぶわけだが、データの配置によっては、何通りもの回帰直線の残差平方和が大差ないという状況がありうる。例えば、独立変数と従属変数（として選んだ変数）が実はまったく無関係であった場合は、データの重心を通るどのような傾きの線を引きても残差平方和はほとんど同じになってしまう。その意味で、回帰直線のパラメータ（回帰係数 b と切片 a ）の推定値の安定性を評価することが大事である。そのためには、 t 値というものが使われている。いま、 Y と X の関係が $Y = a_0 + b_0 X + e$ というモデルで表されるとして、誤差項 e が平均 0、分散 σ^2 の正規分布に従うものとするれば、回帰係数の推定値 a も、平均 a_0 、分散 $\sigma^2/n(1 + M^2/V)$ （ただし M と V は x の平均と分散）の正規分布に従い、残差平方和 Q を誤差分散 σ^2 で割った Q/σ^2 が自由度 $(n - 2)$ のカイ二乗分布に従うことから、

$$t_0(a_0) = \frac{\sqrt{n(n-2)}(a - a_0)}{\sqrt{(1 + M^2/V)Q}}$$

が自由度 $(n - 2)$ の t 分布に従うことになる。しかしこの値は a_0 がわからないと計算できない。 a_0 が 0 に近ければこの式で $a_0 = 0$ と置いた値（つまり $t_0(0)$ 。これを切片に関する t 値と呼ぶ）を観測データから計算した値が $t_0(a_0)$ とほぼ一致し、自由度 $(n - 2)$ の t 分布に従うはずなので、その絶対値は 95% の確率で t 分布の 97.5% 点（サンプルサイズが大きければ約 2 である）よりも小さくなる。つまり、データから計算された t 値がそれより大きければ、切片は 0 でない可能性が高いことになる。言い換えると、「切片が 0 である」という帰無仮説に対する有意確率が計算できることになる。回帰係数についても同様に、

$$t_0(b) = \frac{\sqrt{n(n-2)V}b}{\sqrt{Q}}$$

が自由度 $(n - 2)$ の t 分布に従うことを利用して、「回帰係数が 0」であるという帰無仮説に対する有意確率が計算できる。