

応用統計学 2017 第2回 最小二乗法, 重回帰分析

2017年10月4日(水)

清智也 sei@mist.i.u-tokyo.ac.jp

<http://ur0.pw/yTzt>

- 単回帰：目的変数, 説明変数, 最小二乗法, 回帰式, 回帰係数, 予測値, 残差¹。
- 重回帰：計画行列, 正規方程式, 決定係数, 重相関係数, ダミー変数²。
- 後日扱う内容：回帰モデル, 変数選択, 一般化線形モデル。

演習問題

前回までの演習問題(の一部)の略解は, 講義用 web ページからダウンロードできる。
以下では行列の転置を'で表す。

問題 2-1. 目的変数 y を説明変数 x で説明する回帰式 $\hat{y}(x) = \hat{a} + \hat{b}x$ の回帰係数は

$$\hat{a} = \bar{y} - \hat{b}\bar{x}, \quad \hat{b} = r_{xy}(s_y/s_x)$$

で与えられることを示せ。ただし \bar{x}, \bar{y} は平均, s_x, s_y は標準偏差, r_{xy} は相関係数を表す。

問題 2-2. 講義の web ページにある「LPGA データ」は 2017 年女子オープンゴルフ選手権競技の 3 日目までのスコア x と 4 日目のスコア y からなる CSV ファイルである³。これをダウンロードし, y を x で説明する回帰式を求めよ。また散布図を描き, 回帰直線を重ね描きせよ。

問題 2-3. 2016 年度の東京の日最高気温 (°C) の月平均値 y_t ($t = 1, \dots, 12$) は

10.6, 12.2, 14.9, 20.3, 25.2, 26.3, 29.7, 31.6, 27.7, 22.6, 15.5, 13.8

であった。目的変数を y_t , 説明変数を $x_{t1} = \cos(2\pi t/12)$, $x_{t2} = \sin(2\pi t/12)$ として回帰式を求めよ。

なお, 上のデータは気象庁のサイト <http://www.jma.go.jp/> から「各種データ・資料」に進み, 地点を「東京都 東京」, 年月日を「2016 年」, データの種類を「2016 年の月ごとの値を表示」とすれば得られる。

¹simple regression: response variable, explanatory variable, least squares method, regression equation, regression coefficient, predicted value, residual.

²multiple regression: design matrix, normal equation, coefficient of determination, multiple correlation coefficient, dummy variable.

³日本女子プロゴルフ協会ホームページ (<https://www.lpga.or.jp/>) より作成。

問題 2-4. 目的変数 y を p 個の説明変数 x_1, \dots, x_p で説明する回帰式

$$\hat{y}(x_1, \dots, x_p) = \hat{a} + \hat{b}_1 x_1 + \dots + \hat{b}_p x_p, \quad (x_1, \dots, x_p) \in \mathbb{R}^p,$$

は点 $(\bar{x}_1, \dots, \bar{x}_p, \bar{y})$ を通ることを示せ。また, $\mathbf{1}_n, \mathbf{x}_1, \dots, \mathbf{x}_p$ が互いに直交しているとき, $\hat{a}, \hat{b}_1, \dots, \hat{b}_p$ はどのように表されるか。

問題 2-5 (Simpson's paradox revisited). $y, \mathbf{x}_1, \mathbf{x}_2$ がそれぞれ

$$\mathbf{y} = (4, 5, 6, 1, 2, 3)', \quad \mathbf{x}_1 = (1, 2, 3, 4, 5, 6)', \quad \mathbf{x}_2 = (1, 1, 1, -1, -1, -1)'$$

であるとする。 y を \mathbf{x}_1 と \mathbf{x}_2 で説明する場合の回帰係数と, \mathbf{x}_1 だけで説明する場合の回帰係数を比較せよ。ただし回帰式には定数項も含めるものとする。

問題 2-6. ランク p の行列 $\mathbf{X} \in \mathbb{R}^{n \times p}$ に対し, $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ と定義される行列 \mathbf{P} は直交射影行列になること, すなわち $\mathbf{P}^2 = \mathbf{P}$ かつ $\mathbf{P}' = \mathbf{P}$ が満たされることを示せ。また \mathbf{P} はどのような部分空間への射影になっているか説明せよ。

問題 2-7. 計画行列 \mathbf{X} の QR 分解 $\mathbf{X} = \mathbf{Q}\mathbf{R}$ が得られているとき, 回帰係数を数値的に求める上で有利な点を説明せよ。ただし $\mathbf{X} = \mathbf{Q}\mathbf{R}$ が \mathbf{X} の QR 分解であるとは, \mathbf{Q} が列直交 ($\mathbf{Q}'\mathbf{Q}$ が単位行列) で, \mathbf{R} が正の対角成分を持つ上三角行列となることをいう。

宿題

問題 2-8. アンケート結果を講義の web ページからダウンロードし, 睡眠時間を通学時間で説明する回帰式を求めよ。