

クレジット:

UTokyo Online Education 統計データ解析Ⅱ 2018 小池祐太

ライセンス:

利用者は、本講義資料を、教育的な目的に限ってページ単位で利用することができます。特に記載のない限り、本講義資料はページ単位でクリエイティブ・コモンズ 表示-非営利-改変禁止 ライセンスの下に提供されています。

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

本講義資料内には、東京大学が第三者より許諾を得て利用している画像等や、各種ライセンスによって提供されている画像等が含まれています。個々の画像等を本講義資料から切り離して利用することはできません。個々の画像等の利用については、それぞれの権利者の定めるところに従ってください。



統計データ解析 (II) 第 8 回

小池祐太

2018 年 5 月 30 日

- 1 重回帰分析: 復習
- 2 最小二乗法による線形回帰式の推定の幾何学的解釈
- 3 線形回帰式と標本平均
- 4 分析の評価
 - 残差
 - 標準誤差
 - t 値と p 値
 - 決定係数
 - F 値
- 5 予測

重回帰分析

- **回帰分析 (regression analysis)**

- ▶ ある 1 種類の変数/データを別の変数/データ (1 種類もしくは複数) によって説明もしくは予測するための関係式 (**回帰 (方程) 式 (regression equation)**) を構成することを目的とする分析法

- 説明される側のデータは, 目的変数, 被説明変数, 従属変数, 応答変数などと呼ばれる
- 説明する側のデータは, 説明変数, 独立変数, 共変量などと呼ばれる
- 説明変数が 1 種類の場合を**単回帰 (simple regression)**, 複数の場合を**重回帰 (multiple regression)**と呼ぶ

重回帰分析

- 目的変数を Y , 説明変数を X_1, \dots, X_p で表すことにし, 組 (Y, X_1, \dots, X_p) に対する n 個の観測データ

$$\{(y_i, x_{i1}, \dots, x_{ip})\}_{i=1}^n \quad (1)$$

が得られている状況を考える

- 観測データは次のモデルに従うとする:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n. \quad (2)$$

- ▶ $\beta_0, \beta_1, \dots, \beta_p$ は未知パラメーター (**回帰係数 (regression coefficients)**)
- ▶ $\epsilon_1, \dots, \epsilon_n$ (**誤差項 (error term)**): 独立な確率変数 (多くの場合それぞれ平均 0, 分散 σ^2 の正規分布に従うと仮定)

最小二乗法

- 回帰係数 $\beta := (\beta_0, \beta_1, \dots, \beta_p)^\top$ は**最小二乗法 (least squares)** によって決定 (推定) する
- すなわち, **残差平方和 (residual sum of squares)**

$$S(\beta) := \sum_{i=1}^n |y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})|^2$$

を最小化するベクトル $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)^\top$ を β の推定量とする
(**最小二乗推定量 (least squares estimator)**)

- R での実行: 関数 `lm()`

線形回帰式の行列による表現

- 最小二乗推定量の計算のために、モデル (2) を行列を用いて表現する

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix},$$
$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

とおくと、モデル (2) は

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

と表すことができる

線形回帰式の行列による表現

- 行列 \mathbf{X} は**デザイン行列 (design matrix)** と呼ばれることがある
- 最小二乗推定量がただ一つだけ存在するための必要十分条件は、**グラム行列 (Gram matrix) $\mathbf{X}^T \mathbf{X}$** が正則であることである (配布資料の定理 5.1)
- 以下では $\mathbf{X}^T \mathbf{X}$ の正則性を常に仮定する

最小二乗法による線形回帰式の推定の幾何学的解釈

- 一般に $\hat{\beta}$ を回帰係数の推定値とすると、目的変数の観測データ \mathbf{y} から観測誤差の影響を除いた目的変数の真の値が、

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$$

によって推定できる

- ▶ $\hat{\mathbf{y}}$ を **あてはめ値 (fitted values)** または **予測値 (predicted values)** と呼ぶ
- 特に $\hat{\beta}$ が最小二乗推定量のとき、デザイン行列 \mathbf{X} の列ベクトルたちで張られる \mathbb{R}^n の部分線形空間 (超平面) を $L[\mathbf{X}]$ と書くことにすれば、 $\hat{\mathbf{y}}$ はベクトル \mathbf{y} の $L[\mathbf{X}]$ への直交射影となることがわかる (資料の定理 5.3)
 - ▶ すなわち、 $\hat{\mathbf{y}}$ は、 $L[\mathbf{X}]$ に属するベクトル \mathbf{z} で $\mathbf{y} - \mathbf{z}$ が $L[\mathbf{X}]$ に属するすべてのベクトルに直交するような唯一のもの

最小二乗法による線形回帰式の推定の幾何学的解釈

- 幾何学的には, ベクトル \mathbf{y} からデザイン行列 \mathbf{X} の列ベクトルによって張られる (超) 平面に垂線を下ろした際の垂線の足が $\hat{\mathbf{y}}$ となる
- 次頁に $n = 3, p + 1 = 2$ の場合に最小二乗法による推定の様子を図示したものを示す

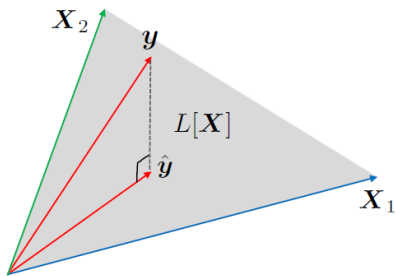


図 1: $n = 3, p + 1 = 2$ の場合の最小二乗法による推定の図示. $\mathbf{X}_1, \mathbf{X}_2$ はそれぞれデザイン行列 \mathbf{X} の第 1 列, 第 2 列からなるベクトルに対応している. グレーの平面が $L[\mathbf{X}]$ に対応.

UTokyo Online Education 統計データ解析 II 2018 小池祐太 CC BY-NC-ND

最小二乗法による線形回帰式の推定の幾何学的解釈

- 特に, ベクトル $\hat{\epsilon} := \mathbf{y} - \hat{\mathbf{y}}$ はあてはめ値 $\hat{\mathbf{y}}$ に直交する:

$$\hat{\epsilon} \cdot \hat{\mathbf{y}} = 0.$$

- $\hat{\epsilon}$ は**残差 (residuals)** と呼ばれ, 回帰式による目的変数のあてはめ値と実際の観測値とのずれを表す
- 実行例 `lse.r`

線形回帰式と標本平均

- 各 $i = 1, \dots, n$ について, 説明変数の i 番目の観測データに対応するベクトル $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ を導入する
- 説明変数および目的変数の標本平均を考える

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

- このとき, $\hat{\beta}$ を最小二乗推定量とすれば, $\bar{y} = (1, \bar{\mathbf{x}}^\top) \hat{\beta}$ が成り立つことが確認できる (資料参照)
- 幾何学的には, 方程式 $y = (1, \mathbf{x}^\top) \hat{\beta}$ によって定まる超平面は常に点 $(\bar{\mathbf{x}}^\top, \bar{y})$ を通るということである
- 実行例 `lse.r`

分析の評価

- 関数 `lm()` のアウトプットに関数 `summary()` を適用した際に表示される, 分析結果の評価をするための各種指標について解説する

残差

- 関数 `summary()` のアウトプットの “Residuals” の欄には残差 $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n$ に対する五数要約 (最小値, 第 1 分位点, メディアン, 第 3 分位点, 最大値) が表示される
- 残差は絶対値が小さいほど回帰式の観測データへのあてはまりがよいこととなるので, 残差のばらつきは小さいほどよい
- 実行例 `resid.r`

標準誤差

- モデル (2) において, 誤差項 $\epsilon_1, \dots, \epsilon_n$ は平均 0, 分散 σ^2 の正規乱数であると仮定する
- 最小二乗推定量 $\hat{\beta}$ は誤差項に依存して変化するため確率変数であるが, いまの仮定の下では平均 β , 共分散行列 $\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$ の $(p+1)$ 変量正規分布に従うことが知られている
- 特に, 行列 $(\mathbf{X}^\top \mathbf{X})^{-1}$ の対角成分を $\xi_0, \xi_1, \dots, \xi_p$ とすれば, 各 $j = 0, 1, \dots, p$ について, $\hat{\beta}_j$ は平均 β_j , 分散 $\sigma^2 \xi_j$ の正規分布に従う
- 正規分布の分散, もしくはその平方根である標準偏差は, 確率変数の値の平均からの離れやすさを表す指標だと解釈できる (大きいほど離れやすい)

標準誤差

- $\hat{\beta}_j$ の値は “真の” 回帰係数値 β_j に近ければ近いほどよい推定であるといえるから、その標準偏差 $\sigma\sqrt{\xi_j}$ は $\hat{\beta}_j$ の推定精度を評価するのに利用できる
- ただし、 σ は未知パラメーターだから、データから推定する必要がある
- σ の推定量としては、通常

$$\hat{\sigma} = \sqrt{\frac{1}{n-p-1} \sum_{i=1}^n \hat{\epsilon}_i^2}$$

が利用される (この推定量が利用される理由の直感的説明については資料の注意 5.1 を参照)

標準誤差

- こうして推定値 $\hat{\beta}_j$ の精度を評価するための指標

$$\hat{\sigma} \sqrt{\xi_j}$$

が得られるが、これを $\hat{\beta}_j$ の**標準誤差**と呼ぶ

- 標準誤差は、関数 `summary()` のアウトプットの “Coefficients” の欄の2列目 “Std. Error” で確認できる
- また、関数 `summary()` のアウトプットの “Residual standard error” の欄で $\hat{\sigma}$ の値を確認できる ($\hat{\sigma}$ は残差のばらつき具合を表す指標として利用できる)
- 実行例 `se.r`

t 値と p 値

- $(n - p - 1)\hat{\sigma}^2/\sigma^2$ は自由度 $n - p - 1$ のカイ二乗分布に従い、かつ $\hat{\beta}$ と独立であることが知られている
- $(\hat{\beta}_j - \beta_j)/(\sigma\sqrt{\xi_j})$ が標準正規分布に従うことに注意すれば、 $\hat{\beta}_j - \beta_j$ を $\hat{\beta}_j$ の標準誤差で割った量

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}\sqrt{\xi_j}}$$

は自由度 $n - p - 1$ の t 分布に従うことがわかる (配布資料 4 章 4.4.4 節参照)

t 値と p 値

- よって, もし $\beta_j = 0$ であったならば, 統計量

$$t = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{\xi_j}}$$

は自由度 $n - p - 1$ の t 分布に従う

- ▶ この t を $\hat{\beta}_j$ の t 値と呼ぶ
- また, 自由度 $n - p - 1$ の t 分布に従う確率変数の絶対値が $|t|$ を超える理論上の確率

$$2 \int_{|t|}^{\infty} f(x) dx, \quad (f(x) \text{ は自由度 } n - p - 1 \text{ の } t \text{ 分布の確率密度関数}) \quad (3)$$

を $\hat{\beta}_j$ の p 値と呼ぶ

t 値と p 値

- $\beta_j = 0$ が成り立つということは, j 番目の説明変数 X_j は回帰式に寄与していないこととなるから, 回帰式から除外しても問題無いことが示唆される
- 上で定義した t 値および p 値は, $\beta_j = 0$ であるか否かという仮説を検証するのに利用できる
- 実際, もし $\beta_j = 0$ という仮説が正しければ, 確率 (3) はそれほど小さくはならないはずなので, もし p 値が想定よりも小さい場合, はじめの仮説である $\beta_j = 0$ が誤っているという結論した方が自然である

t 値と p 値

- 統計の言葉で言うと, t 値及び p 値は, 仮説検定

$$H_0 : \beta_j = 0 \quad \text{vs} \quad H_1 : \beta_j \neq 0$$

に対する検定統計量の t 値と p 値となっている

- また, ここでいうはじめに想定する p 値の下限を**有意水準**といい, 通常 0.01 もしくは 0.05 とすることが多い
- p 値が有意水準より小さいような回帰係数の推定値をもつ説明変数は**有意**であるといわれる
 - ▶ 有意な説明変数は目的変数の変動を説明するのに有用であるといえる

t 値と p 値

- t 値および p 値は、関数 `summary()` のアウトプットの “Coefficients” の欄の 3-4 列目 “t value” および “Pr(>|t|)” で確認できる
- また、p 値の横についているアスタリスク等の記号は、p 値がどの程度小さいかを示している
- 例えば、‘*’ は p 値が 0.05 以下であることを示し、‘**’ は p 値が 0.01 以下であることを示す (記号の意味は “Coefficients” の欄の下部に書いてある)
- 実行例 `tpvalues.r`

決定係数

- **決定係数**は線形回帰分析のあてはまり具合を評価するためのもっとも代表的な指標である
- 決定係数は記号 R^2 で表され, 回帰モデルによる目的変数のあてはめ値 $\hat{y}_1, \dots, \hat{y}_n$ と実際の観測データ y_1, \dots, y_n の相関の2乗として定義される:

$$R^2 = \frac{(\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y}))^2}{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 \sum_{i=1}^n (y_i - \bar{y})^2}. \quad (4)$$

ここに, $\bar{\hat{y}}$ と \bar{y} はそれぞれ $\hat{y}_1, \dots, \hat{y}_n$ と y_1, \dots, y_n の平均を表す:

$$\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

決定係数

- あてはめ値と実際の観測データの変動が近いほどあてはまりが良いと考えられるので、決定係数は高ければ高いほどよい。
- 決定係数は以下のようにも書くことができる (資料参照, もしくは自分で確認):

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}. \quad (5)$$

- (5) の分子と分母をそれぞれ $n - 1$ で割ることで、決定係数はあてはめ値の分散を目的変数の観測データの分散で割ったものだとも解釈できる
- すなわち、目的変数の観測データの分散のうち何パーセントを回帰モデルが説明できているかを表す指標とも解釈できる

決定係数

- 決定係数はさらに以下のようにも書き直せる (資料参照, もしくは自分で確認):

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2}. \quad (6)$$

- (6) より, 決定係数は説明変数を付け加えるほど高くなることからわかるため, 決定係数は本来回帰式に不要である説明変数の効果を過剰に見積もっているおそれがある

決定係数

- この問題を解消するために、推定されて得られた未知パラメータの影響を考慮して以下のように決定係数を修正したものが**自由度調整済み決定係数**である:

$$\bar{R}^2 = 1 - \frac{\frac{1}{n-p-1} \sum_{i=1}^n \hat{\epsilon}_i^2}{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_i)^2}.$$

- なお、決定係数は**寄与率**とも呼ばれる
- 決定係数および自由度調整済み決定係数は、それぞれ関数 `summary()` のアウトプットの “Multiple R-squared” および “Adjusted R-squared” の欄で確認できる
- 実行例 `rsquared.r`

F 値

- t 値は個々の説明変数の要・不要を判断するための指標であったが、説明変数のうち 1 つでも目的変数の説明の役に立つものがあるか否かを判定するための指標に回帰モデルの **F 値**がある
- これは、現在の説明変数を用いて回帰分析を実行することに意味があるかどうかを検証するための指標ともいえる
- 回帰モデルの F 値は次式で定義される:

$$F = \frac{\frac{1}{p} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \frac{n-p-1}{p} \frac{R^2}{1-R^2}.$$

F 値

- もしすべての説明変数が不要, すなわち $\beta_1 = \dots = \beta_p = 0$ であったならば, F は自由度 $p, n - p - 1$ の F 分布に従うことが知られている
- したがって, 自由度 $p, n - p - 1$ の F 分布に従う確率変数が F を超える理論上の確率

$$\int_F^{\infty} f(x) dx, \quad f(x) \text{ は自由度 } p, n - p - 1 \text{ の } F \text{ 分布の確率密度関数} \quad (7)$$

はそれほど小さくはないはずなので, この確率が想定より小さければ回帰分析に意味があると結論付けられる

F 値

- 統計の言葉で言うと, F 値及び確率 (7) は, 仮説検定

$$H_0 : \beta_1 = \dots = \beta_p = 0$$

$$\text{vs } H_1 : \text{ある } j = 1, \dots, p \text{ に対して, } \beta_j \neq 0$$

に対する検定統計量の F 値と p 値となっている

- 回帰モデルの F 値および確率 (7) は, それぞれ関数 `summary()` のアウトプットの “F-statistic” およびその隣の “p-value” の欄で確認できる
- 実行例 `fstatistics.r`

予測

- 回帰分析の目的の1つは、説明変数の新規データが与えられたときに、そのデータに対応する目的変数の値を予測することであるが、これは関数 `predict()` で実行できる
- 実行例 `predict.r`