

# 回帰分析 - 予測と発展的なモデル

数理科学統論J

(Press ? for help, n and p for next and previous slide)

村田 昇

2019.11.01

# 講義の予定

- 第1日: 回帰モデルの考え方と推定
- 第2日: モデルの評価
- **第3日: モデルによる予測と発展的なモデル**

# 回帰分析の復習

# 線形回帰モデル

- **目的変数** を **説明変数** で説明する関係式を構成:
  - 説明変数:  $x_1, \dots, x_p$  (p次元)
  - 目的変数:  $y$  (1次元)
- **回帰係数**  $\beta_0, \beta_1, \dots, \beta_p$  を用いた一次式:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- **誤差項** を含む確率モデルで観測データを表現:

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

# 問題設定

- 確率モデル:

$$y = X\beta + \epsilon$$

- 式の評価: **残差平方和** の最小による推定

$$S(\beta) = (y - X\beta)^\top (y - X\beta)$$

# 解

- 解の条件: **正規方程式**

$$X^{\top} X \beta = X^{\top} y$$

- 解の一意性: **Gram 行列**  $X^{\top} X$  が正則

$$\hat{\beta} = (X^{\top} X)^{-1} X^{\top} y$$

# 寄与率

- **決定係数 (R-squared):**

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

- **自由度調整済み決定係数 (adjusted R-squared):**

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

## $t$ -統計量による検定

- 回帰係数  $\beta_j$  が回帰式に寄与するか否かを検定する
  - 帰無仮説:  $\beta_j = 0$
  - 対立仮説:  $\beta_j \neq 0$  ( $\beta_j$  は役に立つ)
- $t$ -統計量: 各係数ごと,  $\xi$  は  $(X^\top X)^{-1}$  の対角成分

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

- $p$ -値: 自由度  $n-p-1$  の  $t$  分布を用いて計算



# $F$ -統計量による検定

- 説明変数のうち1つでも役に立つか否かを検定する
  - 帰無仮説:  $\beta_1 = \dots = \beta_p = 0$
  - 対立仮説:  $\exists j \beta_j \neq 0$  (少なくとも1つは役に立つ)
- $F$ -統計量: 決定係数(または残差)を用いて計算

$$F = \frac{n-p-1}{p} \frac{R^2}{1-R^2}$$

- $p$ -値: 自由度  $p, n-p-1$  の  $F$  分布で計算

# 演習: 回帰分析とその評価

- 06-wine.r を確認してみよう

# 回帰モデルによる予測

# 予測

- 新しいデータ (説明変数)  $x$  に対する予測値

$$\hat{y} = (1, x^\top) \hat{\beta}, \quad \hat{\beta} = (X^\top X)^{-1} X^\top y$$

- 予測値は元データの目的変数の重み付け線形和

$$\hat{y} = w(x)^\top y$$

- 重みは元データと新規データの説明変数で決定

$$w(x)^\top = (1, x^\top) (X^\top X)^{-1} X^\top$$



# 予測値の分布

- 推定量は以下の性質をもつ多変量正規分布

$$\begin{aligned}\mathbb{E}[\hat{\beta}] &= \beta \\ \text{Cov}(\hat{\beta}) &= \sigma^2 (X^\top X)^{-1}\end{aligned}$$

- この性質を利用して以下の3つの値の違いを評価

$$\hat{y} = (1, x^\top) \hat{\beta} \quad (\text{回帰式による予測値})$$

$$\tilde{y} = (1, x^\top) \beta \quad (\text{最適な予測値})$$

$$y = (1, x^\top) \beta + \epsilon \quad (\text{観測値})$$

( $\hat{y}$  と  $y$  は独立な正規分布に従うことに注意)



# 最適な予測値との差

- 差の分布は以下の平均・分散の正規分布

$$\begin{aligned}\mathbb{E}[\tilde{y} - \hat{y}] &= \mathbf{x}^\top \boldsymbol{\beta} - \mathbf{x}^\top \mathbb{E}[\hat{\boldsymbol{\beta}}] = 0 \\ \text{Var}(\tilde{y} - \hat{y}) &= \underbrace{\sigma^2 \mathbf{x}^\top (X^\top X)^{-1} \mathbf{x}}_{\hat{\boldsymbol{\beta}} \text{の推定誤差の分散}} = \sigma^2 \gamma_c(\mathbf{x})^2\end{aligned}$$

- 正規化による表現

$$\frac{\tilde{y} - \hat{y}}{\sigma \gamma_c(\mathbf{x})} \sim \mathcal{N}(0, 1)$$





# 信頼区間

- 未知の分散を不偏分散で推定

$$Z = \frac{\tilde{y} - \hat{y}}{\hat{\sigma}_{\gamma_c}(\mathbf{x})} \sim \mathcal{T}(n-p-1) \quad (t\text{-分布})$$

- 確率  $\alpha$  の信頼区間 (最適な予測値  $\tilde{y}$  が入ることが期待される区間)

$$(\hat{y} - C_\alpha \hat{\sigma}_{\gamma_c}(\mathbf{x}), \hat{y} + C_\alpha \hat{\sigma}_{\gamma_c}(\mathbf{x}))$$

ただし  $C_\alpha$  は以下を満たす定数

$$P(|Z| < C_\alpha | Z \sim \mathcal{T}(n-p-1)) = \alpha$$



# 観測値との差

- 差の分布は以下の平均・分散の正規分布

$$\mathbb{E}[y - \hat{y}] = \mathbf{x}^\top \boldsymbol{\beta} - \mathbf{x}^\top \mathbb{E}[\hat{\boldsymbol{\beta}}] = 0$$

$$\text{Var}(y - \hat{y}) = \underbrace{\sigma^2 \mathbf{x}^\top (X^\top X)^{-1} \mathbf{x}}_{\hat{\boldsymbol{\beta}} \text{ の推定誤差の分散}} + \underbrace{\sigma^2}_{\text{誤差の分散}} = \sigma^2 \gamma_p(\mathbf{x})^2$$

- 正規化による表現

$$\frac{y - \hat{y}}{\sigma \gamma_p(\mathbf{x})} \sim \mathcal{N}(0, 1)$$



# 予測区間

- 未知の分散を不偏分散で推定

$$Z = \frac{y - \hat{y}}{\hat{\sigma}\gamma_p(\mathbf{x})} \sim \mathcal{T}(n-p-1) \quad (t\text{-分布})$$

- 確率  $\alpha$  の予測区間 (観測値  $y$  が入ることが期待される区間)

$$(\hat{y} - C_\alpha \hat{\sigma}\gamma_p(\mathbf{x}), \hat{y} + C_\alpha \hat{\sigma}\gamma_p(\mathbf{x}))$$

ただし  $C_\alpha$  は以下を満たす定数

$$P(|Z| < C_\alpha | Z \sim \mathcal{T}(n-p-1)) = \alpha$$

- $\gamma_p > \gamma_c$  なので信頼区間より広くなる

# 演習: 信頼区間と予測区間

- [06-interval.r](#)を確認してみよう



# 演習: モデルの更新

- 06-model.r  を確認してみよう

# 発展的なモデル

# 非線形な関係のモデル化

- 目的変数  $Y$
- 説明変数  $X_1, \dots, X_p$
- 説明変数の追加で対応可能
  - 交互作用 (交差項):  $X_i X_j$  のような説明変数の積
  - 非線形変換:  $\log(X_k)$  のような関数による変換

# R: 線形でないモデル式の書き方

- 交互作用を記述するためには特殊な記法がある
- 非線形変換はそのまま関数を記述すればよい
- 1つの変数の多項式は関数 `I()` を用いる

```
## 目的変数 y, 説明変数 x1,x2,x3
## 交互作用を含む式 (formula) の書き方
Y ~ X1 + X1:X2          # X1 + X1*X2
Y ~ X1 * X2              # X1 + X2 + X1*X2
Y ~ (X1 + X2 + X3)^2     # X1 + X2 + X3 + X1*X2 + X2*X3 + X3*X1
## 非線形変換を含む式 (formula) の書き方
Y ~ f(X1)                # f(X1)
Y ~ X1 + I(X1^2)         # X1 + X1^2
```

# 演習: 交互作用

- [06-cross.r](#) を確認してみよう

# カテゴリデータ

- 悪性良性や血液型などの数値ではないデータ
- 適切な方法で数値に変換して対応:
  - 2値の場合は0,1を割り当てる
    - 悪性:1
    - 良性:0
  - 3値以上の場合は **ダミー変数** を利用する (カテゴリ数-1個)
    - A型: (1,0,0)
    - B型: (0,1,0)
    - O型: (0,0,1)
    - AB型: (0,0,0)

# R: カテゴリデータの取り扱い

- 通常は適切に対応してくれる
- カテゴリデータとして扱いたい場合は `factor()` を利用する

```
## データフレーム mydat1, mydat2, mydat3
## 変数名 x, y, z
mydat2 <- transform(mydat1, Y=as.factor(X))
mydat3 <- transform(mydat1, Z=as.factor(X > 0))
```

## 演習: ダミー変数

- [06-dummy.r](#)を確認してみよう



# car package の紹介

# さまざまな評価

- 回帰モデルの評価
    - 与えられたデータの再現
    - 新しいデータの予測
  - モデルの再構築のための視覚化
    - **residual plots**: 説明変数・予測値と残差の関係
    - **marginal-model plots**: 説明変数と目的変数・モデルの関係
    - **added-variable plots**: 説明変数・目的変数をその他の変数で回帰したときの残差の関係
    - **component+residual plots**: 説明変数とそれ以外の説明変数による残差の関係
- などが用意されている

# 演習: car package の使用例

- [06-car.r](#) を確認してみよう

# 演習

- これまでに用いたデータでモデルを更新して評価してみよう
  - 変数間の線形回帰の関係について仮説を立てる
  - モデルのあてはめを行い評価する
    - 説明力があるのか? ( $F$ -統計量,  $t$ -統計量, 決定係数)
    - 残差に偏りはないか? (様々な診断プロット)
    - 変数間の線形関係は妥当か? (様々な診断プロット)
  - 検討結果を踏まえてモデルを更新する (評価の繰り返し)