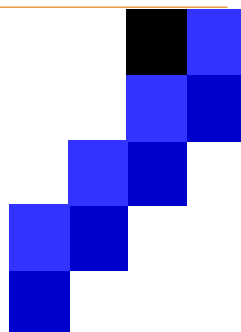


## 実証分析の基礎： 回帰分析とその応用(6) 重回帰分析

松井啓之  
京都大学 経営管理大学院  
2018年9月10日



## 重回帰分析(1)

- 一つの目的変数を、複数の説明変数で予測する事を考える。たとえば3つの独立変数がある場合、重回帰式は  $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$  となる。それぞれの独立変数にかかっている係数を「偏回帰係数」と呼ぶ。

$$\sum e_i^2 = \sum (y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i}))^2 \rightarrow \text{Min}$$

なので、 $\alpha, \beta_1, \beta_2, \beta_3$ で偏微分して0とおくことで正規方程式を作成し、それを解いて求めれば良い

- 単回帰分析と同様に、モデルの当てはまりの良さは(自由度修正済み)決定係数で確認可能
- ただし、単回帰と違い様々な要因が絡み難易度は上がる。
- 当然、偏回帰係数についての統計的検定が必要！

2



## 重回帰分析(2)

- $Y_i = \alpha + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + u_i$ 
  - $i$ は $i$ 番目、 $Y_i$ は被説明変数、 $X_{ji}$ は説明変数 $j$ 、 $u_i$ は誤差項
  - $i=1$ から $n$ まで、この式が成立するとすれば、

$$Y_1 = \alpha + \beta_1 X_{11} + \dots + \beta_k X_{k1} + u_1$$

⋮

$$Y_n = \alpha + \beta_1 X_{1n} + \dots + \beta_k X_{kn} + u_n$$

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & X_{11} & \dots & X_{k1} \\ \vdots & \vdots & & \vdots \\ 1 & X_{1n} & \dots & X_{kn} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \alpha \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$$

とすれば、この重回帰モデルは

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

と表現できる

3



## 重回帰分析(3)

$\boldsymbol{\beta}$ の通常の最小二乗推定量を $\hat{\boldsymbol{\beta}}$ とすると、モデルからの $\mathbf{y}$ の推定値は

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

であり、 $\mathbf{y}$ と $\hat{\mathbf{y}}$ との差は残差

$$\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}}$$

である。これを整理すると最小二乗法の正規方程式

$$(\mathbf{X}'\mathbf{X})\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$$

を得るので、この正規方程式を $\hat{\boldsymbol{\beta}}$ について解くと

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

が得られる。また残差 $\hat{\mathbf{u}}$ は以下の性を持つ。

$$(1)\mathbf{X}'\hat{\mathbf{u}} = \mathbf{0}, (2)\hat{\mathbf{u}}'\hat{\mathbf{y}} = 0, (3)E(\hat{\mathbf{u}}) = \mathbf{0}, (4)\text{Var}(\hat{\mathbf{u}}) = \sigma^2\mathbf{M}$$

$$\text{ただし、}\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

4



## 重回帰分析(4)

### ■ 線形回帰モデルの古典的仮定

- 仮定0: 真のモデルは線形
- 仮定1: 説明変数は非確率変数(=  $X_1, \dots, X_n$  は定数)
- 仮定2: 誤差項の期待値は(全ての  $i$  について) 0
  - $E(u) = 0$
- 仮定3: 分散均一性=誤差項の分散は全ての  $i$  について等しい
  - $E(uu') = \sigma^2 I$  ( $I$  は  $n \times n$  の単位行列)
- 仮定4: 誤差項は互いに独立  $\Rightarrow$  相関は存在しない。
- 仮定5: 誤差項の確率分布は正規分布に従う
  - この仮定と仮定1、仮定2、仮定3をあわせると  
 $u_i \sim N(0, \sigma^2)$  i.i.d.: 誤差項は互いに独立で、同一の正規分布  $N(0, \sigma^2)$  に従う

5



## 重回帰分析(5)

- これまでと同じ議論から、個々の観測  $Y_i$  は正規分布に従う
 
$$Y_i \sim N(\alpha + \beta_1 X_{1i} + \dots + \beta_k X_{ki}, \sigma^2)$$

$$E(Y_i) = \alpha + \beta_1 X_{1i} + \dots + \beta_k X_{ki}$$

$$Var(Y_i) = \sigma^2$$

### ■ ガウス・マルコフの定理

- 古典的仮定1~4が成立するならば、最小二乗推定の期待値と分散は

$$E(\hat{\beta}_j) = \beta_j, \text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{S_{jj}(1 - R_j^2)}, j = 1, 2, \dots, k$$

- ここで、 $S_{jj}$  は、第  $j$  説明変数  $X_{ji}$  の偏差2乗和、 $R_j^2$  は、 $X_{ji}$  をそれ以外の  $k-1$  個の説明変数に重回帰した際の決定係数。上式の分散は、線形不偏推定量の中で最小となる。

6



## 重回帰分析の検定(1)

### ■ 母分散 $\sigma^2$ の不偏推定量

$$s^2 = \frac{1}{n - (k + 1)} \sum \hat{u}_i^2, E(s^2) = \sigma^2$$

### ■ $s^2$ より、各 $b$ の標準誤差

$$s.e.(\hat{\beta}_j) = \frac{s}{\sqrt{S_{jj}(1 - R_j^2)}}, j = 1, 2, \dots, k$$

を得る。さらに、標準誤差を使えば、t統計量となる。

$$t_j = \frac{\hat{\beta}_j - \beta_j}{s / \sqrt{S_{jj}(1 - R_j^2)}} = \frac{\hat{\beta}_j - \beta_j}{s.e.(\hat{\beta}_j)} \sim T(n - (1 + k))$$

- 係数  $\beta_j$  に関する仮説検定  $H_0: \beta_j = \beta_{j*}$  が可能に

7



## 重回帰分析の検定(2)

- $Y_i = \alpha + \beta_1 X_{1i} + \dots + \beta_k X_{ki}$  という回帰関係を仮定し、個々の偏回帰係数  $\beta_j = 0$  という仮説を検定する。

- 回帰係数  $\beta_j$  (説明変数  $x_j$ ) の有意性検定

$$H_0: \beta_j = 0$$

- 帰無仮説が真  $\Rightarrow$  偏回帰係数に意味がない

- 帰無仮説を棄却

$\Rightarrow$  対立仮説を採択:  $x$  と  $y$  の間に統計的な関係がある

「 $\beta_j$  は ( $x_j$  は) 統計的に有意である」

「 $\beta_j$  は有意にゼロと異なる」と表現される

- 統計ソフトで最小二乗検定を実施する場合に、t検定 (t値が計算) されるのは、優位性検定を実施しているため

8



## 重回帰分析の検定(3)

- 回帰モデル全体の評価を考える。
- $Y_i = \alpha + \beta_1 X_{1i} + \dots + \beta_k X_{ki}$  という回帰関係を仮定し、全ての偏回帰係数  $\beta_1, \dots, \beta_k$  が0という仮説を検定する。  
 帰無仮説  $H_0: \beta_1 = \dots = \beta_k = 0$   
 対立仮説  $H_1: H_0$  でない
- 帰無仮説が真  $\Rightarrow$  偏回帰係数に意味がない
- 最小二乗法で推定した決定係数  $R^2$  を用いて  $F$  値を求めると  

$$F = \frac{R^2/K - 1}{(1 - R^2)/(n - K)} \sim F(K - 1, n - K)$$
- $F$  値が棄却域に入らなければ、帰無仮説を棄却できず、説明変数は全て説明力を持たない = モデルが意味を持たない。

9



## 自由度修正済み決定係数

- 決定係数  $R^2$  の問題
  - 一般の回帰モデル(重回帰モデル)では複数の説明変数がいられる  $\Rightarrow$  説明変数が増えると  $R^2$  は必ず値が大きくなる。  
 $\Leftrightarrow$  モデルを複雑にすればするほど(説明変数が多いモデルほど)説明力の高いモデル
- モデルの複雑度に左右されないためには  
 $\Rightarrow$  自由度修正済み決定係数  $\bar{R}^2$

$$\bar{R}^2 = 1 - \frac{RSS/n - (k+1)}{TSS/n - 1} = 1 - \frac{\sum_{i=1}^n e_i^2 / (n - (k+1))}{\sum_{i=1}^n (y_i - \bar{y})^2 / (n - 1)} = 1 - (1 - R^2) \frac{n-1}{n-k}$$

なお、 $k$  は変数の個数

10



## 最適なモデルの決定

- $F$  検定
  - nested model の場合
- Adjusted  $R^2$  を用いる方法
- AIC 基準 (Akaike Information Criteria)  
 $AIC = -2\ln(L) + 2k$   
 $\ln(L)$ : 対数尤度,  $k$ : パラメータの数(説明変数の数)  
 AIC を最小にするようなモデルを選ぶ
- BIC 基準 ベイジアン情報量規準  
 $BIC = T \cdot \log(s^2) + K \cdot \log(T)$   
 $T$  は標本の大きさ,  $s^2$  はモデルの誤差項の分散推定量で、 $K$  はモデルに含まれる係数の数  
 BIC を最小にするようなモデルを選ぶ
- 変数増減法(stepwise regression)



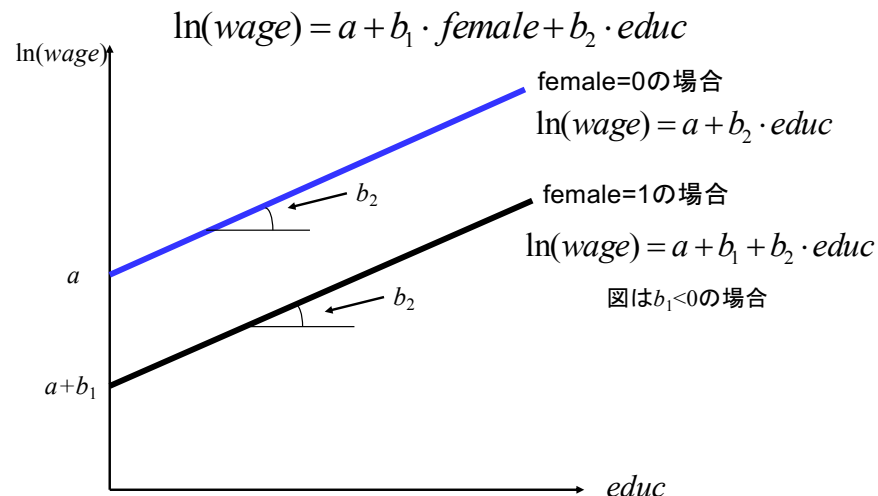
## ダミー変数(1)

- 2つのグループの所属を区別するための質的変数
  - $D_i = \begin{cases} 0 & i \text{ がグループに該当しない場合} \\ 1 & i \text{ がグループに該当する場合} \end{cases}$
  - 複数のグループの場合は、「グループ数-1」のダミー変数を用意
    - 中卒、高卒、大卒の3区分  $\rightarrow D_{\text{高卒}}, D_{\text{大卒}}$  の2つのダミー変数
      - 中卒:  $D_{\text{高卒}}=0, D_{\text{大卒}}=0$  高卒:  $D_{\text{高卒}}=1, D_{\text{大卒}}=0$
      - 大卒:  $D_{\text{高卒}}=1, D_{\text{大卒}}=1$
- $Y = \alpha + \beta X$ 
  - $Y = \alpha + \beta_1 X + \beta_2 D \leftarrow$  切片パラメータのみに影響
  - $Y = \alpha + \beta_1 X + \beta_2 (X \times D) \leftarrow$  交差項: 傾き(回帰)パラメータにも影響

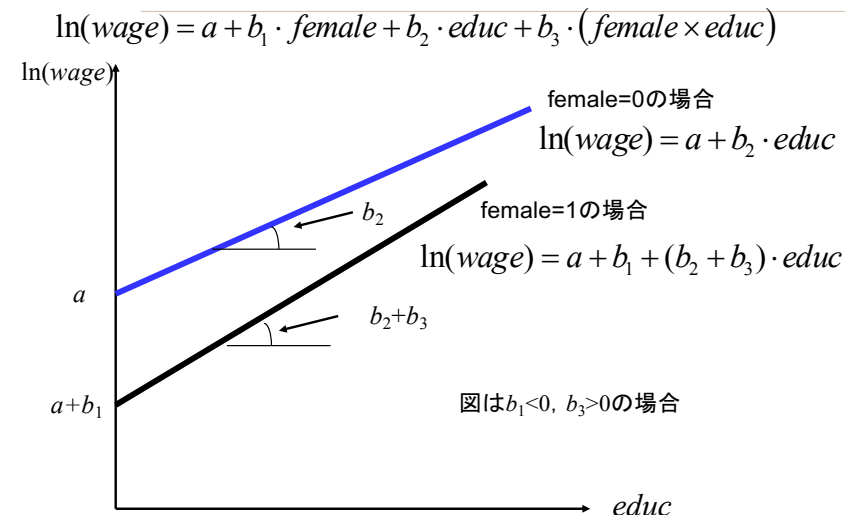
12



## 定数項ダミー



## 傾きのダミー



## ダミー変数(2)

- 重回帰モデルの構造変化が生じる
  - ⇔グループ間での重回帰パラメータ全てが異なる
  - ⇔交差項まで含めたF検定＝大変！
- チョウ(Chow)検定
  - 全体の残差二乗和:  $SSR$
  - 2つのグループの残差二乗和:  $SSR_1, SSR_2$
$$F = \frac{SSR - (SSR_1 + SSR_2)}{SSR_1 + SSR_2} \times \frac{n - 2(k + 1)}{k + 1} \sim F(2, n - 2k - 2)$$
  - F値が棄却域に入れば、帰無仮説「グループ間で違いがない」を棄却する。＝構造変化が起った
    - 時系列データでも利用。その場合には、何時生じたかも確認



## 多重共線性

- 説明変数間で線形関係があるとき、正規方程式で解を一意に定めることが出来ない状態となる。変数間に高い相関がある場合も同様。
  - 指定されたパラメータの符号が理論と合わない。説明変数のt値が小さい反面、決定係数が高い。データ数が増えると、推定値が大きく変動する。新しい説明変数を追加すると、パラメータの符号や値が大きく変わる。
- 相関の高い場合には、どちらかの変数をモデルから、取り除く必要がある。
  - VIF (Variance Inflation Factor) の値が10以上
  - ⇔相関係数が0.95以上
- 変数の変更の際には、適切なモデル化どうか、自由度済み決定係数やAIC、BICを用いて確認。



## 多重共線性の検出

- OLSにおいて説明変数 $x_j$ の係数の分散は次の通りになる。

$$\text{var}(b_j) = \frac{\sigma^2}{S_j(1 - R_j^2)}$$

$\sigma^2$ : 誤差項の分散,  $S_j$ : 説明変数 $x_j$ の平均値の回りの平方和,  $R_j^2$ : 説明変数 $x_j$ を他の説明変数に回帰した場合の $R^2$ (決定係数)

多重共線性 $\rightarrow R_j^2$ が高い $\rightarrow b_j$ の分散が大きくなる

- ・ VIF (Variance inflation factor 分散増幅因子)

$$\text{VIF}(b_j) = \frac{1}{1 - R_j^2}$$



## 均一分散が成立しない場合の対応

- 古典的仮定の定理3が成立しない  
 $\Rightarrow \text{Var}(\hat{\beta})$ を推定できない $\rightarrow$ 標準誤差を求められない
- 不均一分散を確認するために  
 $\Rightarrow$ ホワイト検定、ブルーシュ・ペーガン検定など
  - 補助回帰(残差を二乗した系列で回帰)の決定係数( $nR^2$ )がカイ二乗分布に従うことを利用
  - 帰無仮説: 不均一分散は生じていない
- 不均一分散への対応  
 $\Rightarrow$ (ホワイトの)頑健標準誤差を用いる
  - 一般に、分析ソフトのOLSのオプションで用意されている  
 $\Rightarrow$ 分散の構造が分かっている場合ウェイト付き回帰(WLS)を用いる $\Rightarrow$ 推定の精度を向上させる

18



## 分散不均一性の検出

- Whiteのテスト
  - 残差の平方  $e^2$  を被説明変数
  - 説明変数:  $x_j$ ,  $x_j$ の平方,  $x_j$ と $x_h$ の交差項
  - これらの説明変数の係数が全て0という仮説を検定する

- Breusch and Paganのテスト

estimate:  $y_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i} + u_i$

save:  $e_i = y_i - \alpha - b_1 x_{1,i} - b_2 x_{2,i} - \dots - b_k x_{k,i}$

compute  $e_i^2$

estimate:  $e_i^2 = \delta_0 + \delta_1 x_{1,i} + \delta_2 x_{2,i} + \dots + \delta_k x_{k,i} + v_i$

test  $H_0: \delta_1 = \delta_2 = \dots = \delta_k = 0$

$$\frac{(RSS - TSS) / k}{RSS / (n - (k + 1))} = \frac{ESS / k}{RSS / (n - (k + 1))} \sim F(k, n - (k + 1))$$



## 加重最小二乗法 Weighted Least Square

- 不均一性のテストは検出のみ
  - どのような方法で対処すべきかは教えてくれない
  - 推定する方程式の関数型を変えることで解決する場合もある
- 誤差項の分散がある変数に比例していることがわかっている場合  
 $\rightarrow$  Weighted Least Square 加重最小二乗法
- WLS: 次の式を最小化するように係数を決定

$$\sum_{i=1}^n w_i (y_i - a - b_1 x_{1,i} - \dots - b_k x_{k,i})^2$$

$w_i$ : weight



## Weighted Least Square

次のモデルを考える

$$y_i = \alpha + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i} + u_i \quad (1)$$

ただし,  $\text{var}(u_i) = \sigma_i^2 = h_i \sigma^2$  (誤差項の分散が変数 $h$ に比例している→ 分散不均一性)。このとき次のように式変換すれば

$$\frac{y_i}{\sqrt{h_i}} = \alpha \frac{1}{\sqrt{h_i}} + \beta_1 \frac{x_{1,i}}{\sqrt{h_i}} + \dots + \beta_k \frac{x_{k,i}}{\sqrt{h_i}} + \frac{u_i}{\sqrt{h_i}} \quad (2)$$

$$\text{var}\left(\frac{u_i}{\sqrt{h_i}}\right) = \sigma^2 \rightarrow \text{分散は均一}$$



## Weighted Least Square (2)

(2)式をもとに係数を推定→次の式の最小化

$$\begin{aligned} & \sum_{i=1}^n \left[ \frac{y_i}{\sqrt{h_i}} - \alpha \frac{1}{\sqrt{h_i}} - \beta_1 \frac{x_{1,i}}{\sqrt{h_i}} - \dots - \beta_k \frac{x_{k,i}}{\sqrt{h_i}} \right]^2 \\ &= \sum_{i=1}^n \frac{1}{h_i} [y_i - \alpha - \beta_1 x_{1,i} - \dots - \beta_k x_{k,i}]^2 \\ &= \sum_{i=1}^n w_i (y_i - \alpha - \beta_1 x_{1,i} - \dots - \beta_k x_{k,i})^2 \end{aligned}$$

元のモデルの誤差項の分散が $h$ に比例する→weight 変数を $1/h$ にする



## 非線形の回帰分析(1)

- 線形回帰モデル: 非説明変数と説明変数に線形線形を仮定している。ただし、線形回帰モデルのOLS推定が通用するための要件は、 $Y_i$ と回帰係数の線形性であり、 $X_i$ に関しては、適切な変数変換で対応可能
- 例: 多項式モデル
  - $Y = \alpha + \beta_1 X + \beta_2 X^2 + \dots + \beta_k X^k, X \rightarrow X_1, X^2 \rightarrow X_2, \dots, X^k \rightarrow X_k$  と変数変換することで対応可能
- 例: 対数線形モデル
  - 対数変換することで、掛け算→足し算へ変換可能
  - $Y = \alpha + \beta X_1 \cdot X_2 \cdot \dots \cdot X_k$   
 $\rightarrow \text{Log}(Y) = \text{Log}(\alpha) + \beta \{ \text{Log}(X_1) + \text{Log}(X_2) + \dots + \text{Log}(X_k) \}$



## 対数変換したモデルの解釈

応答変数	説明変数	係数 $b$ の意味
無変換	無変換	説明変数が1単位増えると、応答変数は $b$ だけ増える
無変換	自然対数	説明変数が1%増えると、応答変数が $b$ だけ増える
自然対数	無変換	説明変数が1単位増えると、応答変数が100 <b>b</b> %増える
自然対数	自然対数	説明変数が1%増えると、応答変数が100 <b>b</b> %増える (弾力性)

注:  $b$  が0に近くないときは  $\exp(b) - 1$  を計算する必要がある



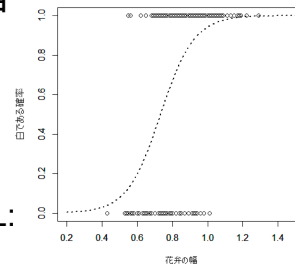
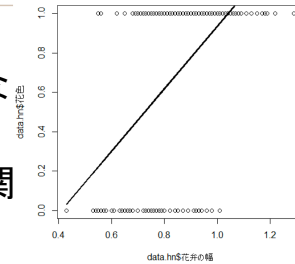


## 非線形の回帰分(2)

- 被説明変数が「0-1」になるような変数の場合に、線形回帰ではなく条件付き確率として割り当てる
- プロビット: 正規分布の累積分布関数を割り当てる。
- ロジット: ロジステック関数を割り当てる。

$$Y = \frac{1}{1 + e^{-(\alpha + \beta_l X_{li} + \dots + \beta_k X_{ki})}}$$

- 均一分散を満たさないので、OLS推定は出来ない。最尤推定量 (ML: Maximum Likelihood estimator) を求める。→ ML推定。



25



## スコアリング

- 個々の見込み客が持つ、自社への将来的な「価値」を予測し、その価値に準じて順位をつけること
- 具体的には、ロジステック回帰分析で推定したロジステック関数を用いて確率を計算する。
  - 確率の予測は、医療分野では疾病リスク、金融分野では貸し倒れリスク、マーケティング分野では見込み客が財・サービスを利用する確率といった形で、多くの分野で応用
  - 例: ある妊婦が未熟児を生まれる確率を求めたい
    - birthwt 「体重2.5kg未満の新生児が出生される原因となる因子」を探索するために調査されたデータ(1986年マサチューセッツ州スプリングフィールドにある湾岸州医療センター)
    - low: 2.5kg未満の出生体重の指標。age: 母親の年齢。lwt: 最後の月経時の母親の体重(ポンド)。race: 母親の人種(1=白、2=黒、3=他)。smoke: 妊娠中の喫煙状態。ptl: 早産の経験。ht: 高血圧の病歴。ui: 子宮過敏性の存在。ftv: 最初の妊娠中の医師の診察数。bwt: 出生体重(グラム)。

26