

# 第7回：単回帰モデルの係数の 検定

【教科書第5章】

北村 友宏

2020年11月13日

# 本日の内容

1. 標準誤差

2. 仮説検定

# gretl での単回帰分析

前回と同様に、「駅へのアクセスのよさがマンション価値に与える影響」を分析するためのモデル

$$price_i = \beta_0 + \beta_1 minutes_i + u_i$$

- ▶  $price_i$  : 中古マンション価格 (万円)
- ▶  $minutes_i$  : 最寄り駅までの所要時間 (分)
- ▶  $i$  : 中古マンション番号

を推定する.

➡ 「中古マンション価格」を「最寄り駅までの所要時間」に回帰する.

# 実習 1

1. gretl を起動.
2. 「ファイル」 → 「データを開く」 → 「ユーザー・ファイル」と操作.
3. setagayaapartment.gdt を選択し、「開く」をクリック.
4. gretl のメニューバーから「モデル」 → 「通常の最小二乗法」と操作.
5. 出てきたウィンドウ左側の変数リストにある price\_10th をクリックし、3つの矢印のうち上の青い右向き矢印をクリック.
  - ▶ 推定式の左辺の変数（被説明変数，従属変数）が price\_10th（万円単位の中古マンション価格）となる.

6. 「デフォルトとして設定」にチェック。
  - ▶ gretl を終了するまでの間、次回以降「通常の最小二乗法」での推定を行う際に、いま選択した変数が自動的に被説明変数（従属変数）に入力される。
7. ウィンドウ左側の変数リストにある minutes をクリックし、3つの矢印のうち真ん中の緑の右向き矢印をクリック。
  - ▶ 推定式の右辺の変数（説明変数，独立変数）が minutes（最寄り駅までの所要時間）となる。
  - ▶ 最初から説明変数リストに入っている const は推定式の切片（定数項）のこと。
8. 「OK」をクリックすると、結果が新しいウィンドウに表示される。

gretl: モデル

ファイル 編集(E) 検定(D) 保存(S) グラフ(G) 分析(A) LaTeX

モデル 1

モデル 1: 最小二乗法 (OLS), 観測: 1-194  
 従属変数: price\_10th

	係数	標準誤差	t値	p値	
const	3092.68	295.260	10.47	1.35e-020	***
minutes	74.5608	28.1685	2.647	0.0088	***
Mean dependent var	3782.577	S.D. dependent var	2150.961		
Sum squared resid	8.62e+08	S.E. of regression	2118.252		
R-squared	0.035207	Adjusted R-squared	0.030182		
F(1, 192)	7.006396	P-value(F)	0.008796		
Log-likelihood	-1759.988	Akaike criterion	3523.976		
Schwarz criterion	3530.512	Hannan-Quinn	3526.623		

このような画面が表示されれば成功。まだ作業があるので、「gretl: モデル」のウィンドウは**まだ閉じない!**

# 出力結果の見方

- ▶ 係数: 回帰係数推定値
- ▶ 標準誤差: 回帰係数の標準誤差
- ▶  $t$  値: 「回帰係数が 0」という帰無仮説の両側  $t$  検定における検定統計量の実現値 ( $t$  値)
- ▶  $p$  値: 両側  $p$  値
- ▶ R-squared: 決定係数

# 標準誤差

- ▶ 推定量の標準偏差の推定値を**標準誤差 (standard error)** という。
- ▶ 回帰係数の OLS 推定量  $\hat{\beta}_0$  と  $\hat{\beta}_1$  の (デフォルトの) 標準誤差は、それぞれ

$$\text{s.e.}(\hat{\beta}_0) = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-2} \cdot \left\{ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right\}},$$

$$\text{s.e.}(\hat{\beta}_1) = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-2} \cdot \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

⇒ これらの標準誤差は、任意の  $i$  について  $V(u_i | x_i)$  が一定 (均一分散) の場合のみ正しい。

# 頑健標準誤差

- ▶  $V(u_i | x_i)$  が一定でないことを（条件付き）不均一分散（heteroskedasticity）という。
- ▶ 不均一分散があっても厳密な標準誤差を求めするために、頑健標準誤差（robust standard error）が開発されている。
- ▶ gretl では、例えば White の頑健標準誤差などを出力できる。
  - ▶ 「gretl: モデル推定」ダイアログボックスの、「頑健標準誤差を使用する」をチェックすればよい。

- ▶ **経済学分野**の実証分析では、誤差項  $u_i$  に**不均一分散**があることを**前提**として**頑健標準誤差**を計算する**場合が多い**.
- ▶ 頑健標準誤差のほうがデフォルトの標準誤差より大きくなることもあれば、小さくなることもある.

# 仮説検定

- ▶  $y_i$  や  $x_i$  は様々な値をとり，観測される前はどのような値になるかが不確定（**確率変数**，**random variable**）.
- ▶  $y_i$  や  $x_i$  を用いて計算する  $\bar{y}$  や  $\bar{x}$  の値も不確定.
- ▶  $y_i, \bar{y}, x_i, \bar{x}$  を用いて計算する  $\hat{\beta}_0$  や  $\hat{\beta}_1$  の値も不確定.

⇒ 例えば回帰係数  $\beta_1$  の推定値として  $\hat{\beta}_1 = 74.5608$  という値が得られても，「推定値  $\hat{\beta}_1$  は真の  $\beta_1$  の値と必ずしも同じではなく，真の  $\beta_1$  は0で，その推定値  $\hat{\beta}_1$  は様々な値をとりうる中でたまたま74.5608になった」可能性もある.

⇒ **仮説検定 (hypothesis testing)** を行い，「真の  $\beta_0$  や  $\beta_1$  が0かどうか」を検証する.

gretl などの統計解析ソフトで線形回帰モデルを推定すると、各回帰係数  $\beta_j$  (単回帰の場合  $j = 0, 1$ ) について、

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

を検定するのに必要な情報が出力される。

- ▶ 回帰分析では、通常は両側検定を行う。

# 帰無仮説

- ▶ とりあえず「真」であると想定する仮説を**帰無仮説 (null hypothesis)** という.
  - ▶  $H_0$  と書くことが多い.
  - ▶ e.g.,  $H_0 : \beta_1 = 0$ .
  - ▶  $H_0$  は必ず「=」または「 $\leq$  や  $\geq$ 」を使った式.  
「 $\beta_1 < 0$ 」を  $H_0$  とする検定は不可能.
- ▶ まずは  $H_0$  が「真」であると仮定し、それを「偽」とするための証拠を探す.
  - ▶ 刑事裁判における推定無罪の原則と同様.

⇒ 具体的には、検定統計値を計算する.
- ▶ 標本の関数を**統計量 (statistic)** という.
  - ▶ e.g., 標本平均, 標本分散など
- ▶ 検定に用いる統計量を**検定統計量 (test statistic)** といい、その実現値を**検定統計値** という.

- ▶ 仮に  $H_0$  が真であれば，計算した検定統計値が5%や1%の**わずかな確率**でしか生じえない**値**になっている



それを証拠として  $H_0$  を偽と判断し， $H_0$  を**棄却する (reject)** .

- ▶ 仮に  $H_0$  が真であれば，計算した検定統計値が**小さすぎない確率**で生じうる**値**になっている



$H_0$  を偽とする証拠が不十分であり，偽とはいえないと判断し， $H_0$  を**採択する (accept)** .

- ▶ 15%や20%は「小さすぎない」.
- ▶ 「 $H_0$  は真」という判断ではない.

# 対立仮説

- ▶  $H_0$  が偽のときに代わりに採択する仮説を**対立仮説 (alternative hypothesis)** という.
  - ▶  $H_1$  と書くことが多い.
  - ▶ e.g.,  $H_1 : \beta_1 \neq 0$ .
  - ▶  $H_1$  は「 $\neq, <, >$ 」を使った式で設定できる.
- ▶ **両側検定 (two-sided test)** 問題の定式化 :

$$H_0 : \beta_1 = 0 \quad \text{vs} \quad H_1 : \beta_1 \neq 0$$

- ▶  $H_0$  の意味は、「その説明変数は被説明変数と相関していない」
- ▶  $H_1$  の意味は、「その説明変数は被説明変数と相関している」

# 有意水準

- ▶  $H_0$  が真なのに棄却することを第 1 種の誤り (type I error) という。
- ▶  $H_0$  を真としたときに、検定統計値が「わずかな確率でしか生じえない値」かの判断の基準となる確率、また、許容する第 1 種の誤りの確率を有意水準 (significance level) という。
  - ▶ 通常は 10%, 5%, 1% に設定。
  - ▶ e.g., 「有意水準 5% で  $H_0$  が棄却された」
    - ⇒ 仮に  $H_0$  が真であれば、そんな検定統計値が出てくる確率は 5% 以下に過ぎない ( $H_0$  を偽とする証拠) ので  $H_0$  を棄却。
    - ⇒ 言い換えると、 $H_0$  が真のとき、「そんな検定統計値」は 5% 以下の確率で出現しうる。
    - ⇒  $H_0$  を棄却する第 1 種の誤りを犯すことが、多くとも 5% の確率でありうる。

▶  $H_0$  (係数は 0) 棄却

↳ 「その回帰係数は統計的に有意に 0 と異なる」と判断.

- ▶ 「その説明変数は被説明変数と統計的に有意に相関している」と解釈.
- ▶ 定数項の検定の場合は「定数項は統計的に有意に 0 と異なる」と解釈.

▶  $H_0$  (係数は 0) 採択

↳ 「その回帰係数は 0 と異なるとは言えない」と判断.

- ▶ 「その説明変数は被説明変数と相関しているとは言えない」と解釈.
- ▶ 定数項の検定の場合は「定数項は統計的に有意に 0 と異なるとは言えない」と解釈.

# $p$ 値による判断

- ▶ 検定統計量（の絶対値）が実現値（検定統計値）を超える（以上になる）確率を  $p$  値という.
  - ▶  $p$  値が 0.1 以下（未満）：有意水準 10%で  $H_0$  を棄却.
  - ▶  $p$  値が 0.05 以下（未満）：有意水準 5%で  $H_0$  を棄却.
  - ▶  $p$  値が 0.01 以下（未満）：有意水準 1%で  $H_0$  を棄却.

⇒  $p$  値を見て，帰無仮説の採択・棄却を判断できる.

※検定統計量が連続型の確率分布（正規分布， $t$  分布，カイ二乗分布， $F$  分布など）に従う場合，「以上」と「超える」，「以下」と「未満」は区別しなくて良い.

gretl では，モデル推定結果の各説明変数の行の右端にアスタリスク（\*）が表示され，\*の個数を見れば，「有意水準何%で『回帰係数は0』の  $H_0$  を棄却できるか」が分かる．

- ▶ （アスタリスクなし）：有意水準 10%でも「係数は0」の  $H_0$  採択．
- ▶ \*：有意水準 10%で，「係数は0」の  $H_0$  棄却．
- ▶ \*\*：有意水準 5%で，「係数は0」の  $H_0$  棄却．
- ▶ \*\*\*：有意水準 1%で，「係数は0」の  $H_0$  棄却．

## $t$ 値による判断

定数項ありの単回帰の場合， $\beta_j = 0$  という  $H_0$  を検定するための  $t$  検定統計量は，

$$t = \frac{\hat{\beta}_j}{\text{s.e.}(\hat{\beta}_j)} \sim t(n-2).$$

- ▶ 観測値数が十分に大きいとき， $t$  値の絶対値がほぼ 2 を超えていれば， $H_0$  を棄却と判断（大雑把な判断）。
  - ↳ 「有意水準何%で  $H_0$  を棄却できるか」を厳密に判断するには， $t$  値ではなく  $p$  値を見る。

## 実習 2

1. gretl のメニューバーから「モデル」→「通常の最小二乗法」と操作. 説明変数（独立変数）は必ず前回の選択内容が記録されており，被説明変数（従属変数）は前回「デフォルトとして設定」にチェックしていれば前回の選択内容が記録されている.
2. 従属変数の入力ボックスに price\_10th が入力されていなければ，出てきたウィンドウ左側の変数リストにある price\_10th をクリックし，3つの矢印のうち上の青い右向き矢印をクリック.
  - ▶ 推定式の左辺の変数（被説明変数，従属変数）が price\_10th（万円単位の中古マンション価格）となる.

3. 「頑健標準誤差を使用する」にチェック.
  - ▶ 不均一分散に対して頑健な、White の標準誤差が計算され、推定式の誤差項  $u_i$  の分散に関する仮定が誤っていても、より厳密な分析ができるようになる.
4. 「OK」をクリックすると、結果が表示される.

gretl: モデル

ファイル 編集(E) 検定(D) 保存(S) グラフ(G) 分析(A) LaTeX

モデル 1 ✕ モデル 2 ✕

モデル 2: 最小二乗法 (OLS), 観測: 1-194  
 従属変数: price\_10th  
 不均一分散頑健標準誤差, バリエーション HC1

	係数	標準誤差	t値	p値	
const	3092.68	245.524	12.60	6.55e-027	***
minutes	74.5608	22.0194	3.386	0.0009	***
Mean dependent var	3782.577	S.D. dependent var	2150.961		
Sum squared resid	8.62e+08	S.E. of regression	2118.252		
R-squared	0.035207	Adjusted R-squared	0.030182		
F(1, 192)	11.46597	P-value(F)	0.000860		
Log-likelihood	-1759.988	Akaike criterion	3523.976		
Schwarz criterion	3530.512	Hannan-Quinn	3526.623		

このような画面が表示されれば成功。「gretl: モデル」のウィンドウは**まだ閉じない!**

# モデル推定結果

## ▶ 最寄り駅所要時間の係数

- ▶ 74.5608 (符号は正)
- ▶  $t$  値は 3.386,  $p$  値は 0.0009
  - ➡ 仮に「minutes の係数が 0」だとすると, 3.386 という  $t$  値は 0.09% の確率 (1% を下回る確率) でしか出てこない.
  - ➡ 有意水準 1% で, 係数ゼロの  $H_0$  棄却.
  - ➡ 最寄り駅までの所要時間はマンションの価格と統計的に有意に相関している.
  - ➡ 最寄り駅までの所要時間が 1 分長くなると, マンションの市場価値が 74.5608 万円 (745,608 円) 高くなる (?)

## ▶ 定数項

- ▶ 3092.68
- ▶  $t$  値は 12.60,  $p$  値は  $6.55 \times 10^{-27}$ 
  - ➡ 仮に「定数項が 0」だとすると, 12.6 という  $t$  値は  $6.55 \times 10^{-27}$ , つまりほぼ 0% の確率 (1% を下回る確率) でしか出てこない.
  - ➡ 有意水準 1% で, 係数ゼロの  $H_0$  棄却.
  - ➡ 定数項は統計的に有意に 0 と異なる.

## ▶ 決定係数

- ▶  $R^2 = 0.035207$ .
  - ➡ 「最寄り駅までの所要時間」の違いで, 「価格」のバラつきが約 3.5% のみ説明できる.

## ▶ 標準誤差の違い

- ▶ デフォルトの標準誤差を用いた推定結果と、Whiteの頑健標準誤差を用いた推定結果では、標準誤差や $t$ 値に差が生じているが、今回は両方で $t$ 検定の判断に違いは生じていない。

※使うデータ，推定するモデルによっては，両方で $t$ 検定の判断が異なる場合がある。

## 実習 3

1. 「モデル 2」が表示されている状態で、「gretl:モデル」のウィンドウのメニューバーから「ファイル」→「名前を付けて保存」と操作.
2. 「標準テキスト」を選び、「OK」をクリック.
3. results20201113.txt という名前で 2020microdatag フォルダに保存. すると、表示された推定結果をそのままテキストファイルで保存できる.

本日の作業はここまで.

今回は gretl のデータセットに変更を加えていないので、**gretl のデータセット (setagayaapartment.gdt) を上書き保存する必要はない.**