

回帰分析I

15. 不均一分散

1

イントロダクション

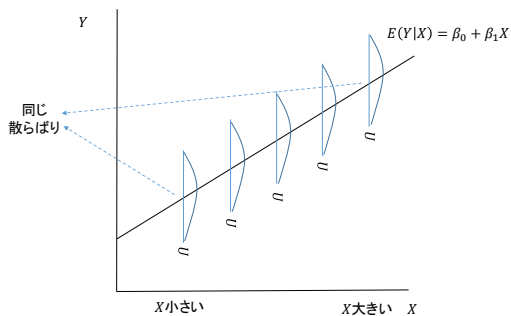
- この講義ノートで扱うトピックは「誤差項の分散」です。
- これまで誤差項の分散について、「均一分散」仮定を置いてきました（講義ノート9、10参照）。
- すなわち、説明変数がいかなる値でも、誤差項の分散は一定：

$$VAR(U|X_1, \dots, X_k) = \sigma^2$$

どういう意味？

2

均一分散というのは、ざっくり言うと、
説明変数 (X) の値に関わらず、誤差項の散らばりが同じということ



3

- この均一分散の仮定は、最小二乗推定量の不偏性や一致性には必要ありませんでした。
- それでは、この仮定は何のためにしたのでしょか？
- この仮定は最小二乗推定量 $\hat{\beta}_1$ の分散を導出する際に使いました。
- 話を簡単にするために単純回帰モデルを考えます。以下は講義ノート9の復習です。

$$Y = \beta_0 + \beta_1 X + U$$

- 最小二乗推定量 $\hat{\beta}_1$ は

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

と表すことができます。

- これを使って最小二乗推定量 $\hat{\beta}_1$ の分散を導出します。

4

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N)$$

$$= VAR\left(\beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 VAR\left(\sum_{i=1}^N (X_i - \bar{X}) U_i \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 VAR(U_i | X_1, \dots, X_N)\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 \sigma^2\right)$$

ここで誤差項の均一分散の仮定を使いました

$$= \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

ちなみに、これは中間試験の問題になってました

5

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N)$$

$$= VAR\left(\beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 VAR\left(\sum_{i=1}^N (X_i - \bar{X}) U_i \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 VAR(U_i | X_1, \dots, X_N)\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 \sigma^2\right)$$

ここで誤差項の均一分散の仮定を使いました

$$= \sigma^2 \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2\right)$$

でも、もし誤差項が均一分散ではなかったら？

$$= \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

6

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N)$$

$$= VAR\left(\beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 VAR\left(\sum_{i=1}^N (X_i - \bar{X}) U_i \middle| X_1, \dots, X_N\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 VAR(U_i | X_1, \dots, X_N)\right)$$

$$= \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2 \sigma^2\right)$$

ここで誤差項の均一分散の仮定を使いました

$$= \sigma^2 \left(\frac{1}{\sum_{i=1}^N (X_i - \bar{X})^2}\right)^2 \left(\sum_{i=1}^N (X_i - \bar{X})^2\right)$$

でも、もし誤差項が均一分散ではなかったら？

$$= \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

最小二乗推定量 $\hat{\beta}_1$ の分散は $\sigma^2 / \sum_{i=1}^N (X_i - \bar{X})^2$ ではない！

均一分散の仮定が正しくなかったら？

- もし均一分散の仮定が正しくなかったら、

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

またこれに平方根をとった

$$SD(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sigma}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2}}$$

は正しくないということが分かりました。

- では、これの何が問題なのでしょう？

8

- $SD(\hat{\beta}_1|X_1, \dots, X_N) = \frac{\sigma}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2}}$ この結果、何に使いましたっけ？
- SD の σ を(残差を使って推定した) $\hat{\sigma}$ で置き換えたものが、最小二乗推定量 $\hat{\beta}_1$ の標準誤差です。

$$SE(\hat{\beta}_1|X_1, \dots, X_N) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2}}$$

- これどこで使いましたっけ？
- t 値ですね。

$$t = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1|X_1, \dots, X_N)}$$

9

Source	SS	df	MS	Number of obs = 1000000
Model	359065.326	1	359065.326	F(1, 999999) =
Residual	750369.903999999	.750371404		Prob > F = 0.0000
Total	1109435.239999999	1.10943634		R-squared = 0.3236
				Adj R-squared = 0.3236
				Root MSE = .86624

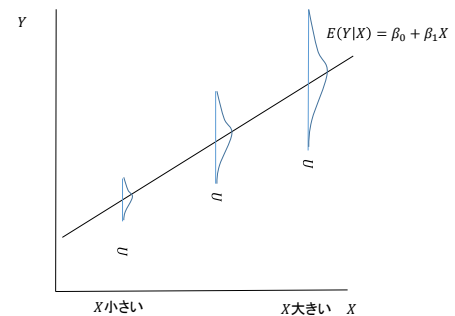
- $VAR(\hat{\beta}_1|X_1, \dots, X_N) = \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$ が正しくないなら、それに基づいて計算される最小二乗推定量 $\hat{\beta}_1$ の標準誤差 $SE(\hat{\beta}_1|X_1, \dots, X_N) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2}}$ も正しくない。
- すると、 $t = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1|X_1, \dots, X_N)}$ も正しい値ではない。
- なので p 値も正しくない → 仮説検定の結果は信用できない!
- 標準誤差が正しくなければ、信頼区間も正しくない!ですね。

10

不均一分散

- 均一分散は、誤差項の条件分散が X の値に依存しないで一定。
- それ以外の場合は、誤差項は「不均一分散」となります。
- 次のページの図は、不均一分散の例。
- ここでは、 X の値が大きいに、誤差項の散らばりが大きくなるとして図を描いています。

11



12

不均一分散の具体例

- 具体例を使って、不均一分散が何を意味するかを説明します。
- 次のモデルを考えましょう。

$$WAGE_i = \beta_0 + \beta_1 MALE_i + U_i$$

- このモデルは

$$WAGE_i = \beta_0 + U_i \quad (\text{女性})$$

$$WAGE_i = \beta_0 + \beta_1 + U_i \quad (\text{男性})$$

ということですね。

- 女性のとき: U_i は、女性 i の賃金が、母集団の女性の平均賃金 β_0 とどれだけ異なるのかを捉える。
- 男性のとき: U_i は、男性 i の賃金が、母集団の男性の平均賃金 $\beta_0 + \beta_1$ とどれだけ異なるのかを捉える。

13

- 従って、、、
- 「 U_i の分散が $MALE_i$ に依存しない」とは、「賃金の分散が男性と女性で等しい」ということと同じです。
- 言い換えると、この例では、母集団における賃金の分散が、男女で等しければ誤差項は均一分散。
- 等しくなければ誤差項は不均一分散ということになります。

14

均一分散と不均一分散はどちらがより現実的？

- 賃金のモデルについて考えてみましょう。
- 過去のことを言えば、トップクラスの高給の仕事についていたのは、男性がほとんどで女性は少なかった。
- 程度の差はあれ、現在でもある程度はそうと言えるでしょう。
- 一方で、低賃金の仕事に就く男性はいつの時代にも存在します。
- 従って、女性の賃金式の誤差項の分散の方が、男性の賃金式の誤差項の分散よりおそらく小さいと考えられます。
- 従って、不均一分散であると考えるのが妥当でしょう。

15

- 一般的なことをいうと、誤差項は均一分散に従うはず、なんてことを示唆する理論は普通はありません。
 - 誤差項が均一分散に従うかどうかは、実際のデータ次第ですね。
- そのため、最近の実証分析では、「誤差項は不均一分散している」ことを前提にするのが普通です。
- 前提にして何？
- 問題なのは、「誤差項は均一分散している」という仮定に基づいて計算される最小二乗推定量の標準誤差でした。
- よって、「誤差項が不均一分散しているかもしれないことを前提にし、それを考慮に入れて最小二乗推定量の標準誤差を計算するのが最近の実証分析の常識となっています。
- 以下では、その方法について説明します。

16

不均一分散に対して頑健な標準誤差

- ここでは話を簡単にするために単純回帰モデルを考えます。

$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

- 添え字 i は以下のことを強調するためにつけました。

- 誤差項は不均一分散を仮定します。

$$VAR(U_i | X_i) = \sigma_i^2$$

- このポイントは、 σ^2 についている添え字の i です。

- これは、誤差項の分散が、 X_i の値に依存し一定ではないことを意味します。

17

- 最小二乗推定量 $\hat{\beta}_1$ の分散は

$$\begin{aligned} VAR(\hat{\beta}_1 | X_1, \dots, X_N) &= VAR\left(\beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \middle| X_1, \dots, X_N\right) \\ &= \frac{1}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2} VAR\left(\sum_{i=1}^N (X_i - \bar{X}) U_i \middle| X_1, \dots, X_N\right) \\ &= \frac{1}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2} \left(\sum_{i=1}^N (X_i - \bar{X})^2 VAR(U_i | X_1, \dots, X_N)\right) \end{aligned}$$

ここまでは均一分散のときと同じです。

18

- 最小二乗推定量 $\hat{\beta}_1$ の分散は

$$\begin{aligned} VAR(\hat{\beta}_1 | X_1, \dots, X_N) &= VAR\left(\beta_1 + \frac{\sum_{i=1}^N (X_i - \bar{X}) U_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \middle| X_1, \dots, X_N\right) \\ &= \frac{1}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2} VAR\left(\sum_{i=1}^N (X_i - \bar{X}) U_i \middle| X_1, \dots, X_N\right) \\ &= \frac{1}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2} \left(\sum_{i=1}^N (X_i - \bar{X})^2 VAR(U_i | X_1, \dots, X_N)\right) \end{aligned}$$

ここが均一分散の場合とは違う点です。

$$= \frac{1}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2} \left(\sum_{i=1}^N (X_i - \bar{X})^2 \sigma_i^2\right)$$

19

- すなわち不均一分散 $VAR(U_i | X_i) = \sigma_i^2$ の仮定の下では、最小二乗推定量 $\hat{\beta}_1$ の分散は

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sum_{i=1}^N (X_i - \bar{X})^2 \sigma_i^2}{\left[\sum_{i=1}^N (X_i - \bar{X})^2\right]^2}$$

- ちなみに均一分散の仮定の下では、

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sigma^2}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

- 不均一分散の時は、 σ_i^2 が定数ではないので、 Σ の前に出ない。

- これにより、均一分散と不均一分散のときの最小二乗推定量 $\hat{\beta}_1$ の分散が一般的には異なることになります。

- 従って、もし誤差項が不均一分散なら、誤差項が均一分散の仮定の下で導出された分散の式は正しくありません。

20

$$VAR(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sum_{i=1}^N (X_i - \bar{X})^2 \sigma_i^2}{\left[\sum_{i=1}^N (X_i - \bar{X})^2 \right]}$$

- この分散の推定量として、次の式が使われます:

$$\widehat{VAR}(\hat{\beta}_1 | X_1, \dots, X_N) = \frac{\sum_{i=1}^N (X_i - \bar{X})^2 \hat{\sigma}_i^2}{\left[\sum_{i=1}^N (X_i - \bar{X})^2 \right]}$$

残差の二乗

- この推定量はWhite(1980)が開発したもので、サンプルサイズが十分に大きい場合、いかなる形の不均一分散にも妥当な推定量であることが証明されています。

White, H. (1980) "A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity," *Econometrica*, 48, 817-838.

21

- これに平方根をとったものを、「**不均一分散に対して頑健(ロバスト)な標準誤差**」と呼びます。単に「**ロバスト標準誤差**」と呼ぶこともあります。

$$Robust\ SE(\hat{\beta}_1 | X_1, \dots, X_N) = \sqrt{\widehat{VAR}(\hat{\beta}_1 | X_1, \dots, X_N)}$$

- 誤差項の均一分散は、誤差項の不均一分散の特殊ケース。
- 従って、たとえ誤差項が均一分散であったとしても、サンプルサイズが十分に大きい場合、ロバスト標準誤差は妥当な推定量です。
- 不均一分散に対してロバストなt検定**は、この標準誤差を用いて行われます。
- 例えば、 $H_0: \beta_1 = \alpha, H_1: \beta_1 \neq \alpha$ なら、t統計量は

$$t = \frac{\hat{\beta}_1 - \alpha}{Robust\ SE}$$
- これを「**不均一分散に対してロバストなt統計量**」と呼びます。
- 検定の仕方はこれまでと全く同じです。

22

自分で実証論文を書くとき

- たとえ誤差項が均一分散であったとしても、サンプルサイズが十分に大きい場合、ロバスト標準誤差は妥当な推定量。
- そのため最近の実証分析では、ロバスト標準誤差、そしてそれに基づいた検定結果を報告するのが普通です。
- ほとんどの論文で普通の標準誤差は報告されていません。
- 普通の標準誤差「のみ」を報告すると、論文の読み手が推定結果を全く信用しないこととなります。
- 従って、皆さんが実証論文を書く場合は、ロバスト標準誤差を使うのが無難だと思います。

23

使ってみましょう

- データ: hprice1_practice.dta
 - Wooldridge (2006) *Introductory Econometrics: A Modern Approach*, Thomson, South-Western からのデータで一部改変したものの。
- 変数: price (家の価格、単位\$1000)
 lotsize (土地面積、単位square feet)
 sqft (建物面積、単位square feet)

24

```
. reg pricemcy lotsizetb sqrrttb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
				R-squared	=	0.6631
				Adj R-squared	=	0.6552
Total	91785450.8	87	1055005.18	Root MSE	=	603.12

pricemcy	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	-.306522 1.2589
sqrrttb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- ロバスト標準誤差を使う場合は一番最後に robust とつけるだけでいいです。

```
. reg pricemcy lotsizetb sqrrttb, robust
```

Linear regression

Number of obs	=	88
F(2, 85)	=	32.81
Prob > F	=	0.0000
R-squared	=	0.6631
Root MSE	=	603.12

pricemcy	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrrttb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

25

```
. reg pricemcy lotsizetb sqrrttb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
				R-squared	=	0.6631
				Adj R-squared	=	0.6552
Total	91785450.8	87	1055005.18	Root MSE	=	603.12

pricemcy	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	-.306522 1.2589
sqrrttb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- ロバスト標準誤差を使う場合は一番最後に robust とつけるだけでいいです。

推定された係数は変わりません。
いずれも最小二乗法
で推定されていますから。

```
. reg pricemcy lotsizetb sqrrttb, robust
```

Linear regression

Number of obs	=	88
F(2, 85)	=	32.81
Prob > F	=	0.0000
R-squared	=	0.6631
Root MSE	=	603.12

pricemcy	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrrttb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

26

```
. reg pricemcy lotsizetb sqrrttb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
				R-squared	=	0.6631
				Adj R-squared	=	0.6552
Total	91785450.8	87	1055005.18	Root MSE	=	603.12

pricemcy	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	-.306522 1.2589
sqrrttb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- ロバスト標準誤差を使う場合は一番最後に robust とつけるだけでいいです。

よって回帰直線の標準誤差や決定係数
なども変わりません。

```
. reg pricemcy lotsizetb sqrrttb, robust
```

Linear regression

Number of obs	=	88
F(2, 85)	=	32.81
Prob > F	=	0.0000
R-squared	=	0.6631
Root MSE	=	603.12

pricemcy	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrrttb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

27

```
. reg pricemcy lotsizetb sqrrttb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
				R-squared	=	0.6631
				Adj R-squared	=	0.6552
Total	91785450.8	87	1055005.18	Root MSE	=	603.12

pricemcy	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	-.306522 1.2589
sqrrttb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- ロバスト標準誤差を使う場合は一番最後に robust とつけるだけでいいです。

変わるのは標準誤差です。

```
. reg pricemcy lotsizetb sqrrttb, robust
```

Linear regression

Number of obs	=	88
F(2, 85)	=	32.81
Prob > F	=	0.0000
R-squared	=	0.6631
Root MSE	=	603.12

pricemcy	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrrttb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

28

```
. reg pricemy lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs =	88
Model	60866850	2	30433425	F(2, 85)	= 83.67
Residual	30918600.8	85	363748.245	Prob > F	= 0.0000
				R-squared	= 0.6631
				Adj R-squared	= 0.6552
Total	91785450.8	87	1055005.18	Root MSE	= 603.12

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	.306522 1.2589
sqrfthb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- ロバスト標準誤差を使う場合は一番最後に robust とつけるだけでいいです。

```
. reg pricemy lotsizetb sqrfthb, robust
```

Linear regression

変わるのは標準誤差です。

それに伴い、t値、p値、信頼区間も変わります。

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrfthb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

29

```
. reg pricemy lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs =	88
Model	60866850	2	30433425	F(2, 85)	= 83.67
Residual	30918600.8	85	363748.245	Prob > F	= 0.0000
				R-squared	= 0.6631
				Adj R-squared	= 0.6552
Total	91785450.8	87	1055005.18	Root MSE	= 603.12

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	.306522 1.2589
sqrfthb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
._cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

この例では、土地面積の有意性に変化があります。

普通の標準誤差を使うと土地面積は1%水準でも有意ですが、ロバスト標準誤差を使った場合は5%水準で有意ではありません。(10%水準で有意ですが、10%水準だとエビデンスとしては弱い)

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.4457812	1.76	0.083	-.1035566 1.669108
sqrfthb	49.39334	6.234889	7.92	0.000	36.99671 61.78997
._cons	59.32415	336.7108	0.18	0.861	-610.1472 728.7955

30

```
. reg pricemy lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs =	88
Model	60866850	2	30433425	F(2, 85)	= 83.67
Residual	30918600.8	85	363748.245	Prob > F	= 0.0000
				R-squared	= 0.6631
				Adj R-squared	= 0.6552
Total	91785450.8	87	1055005.18	Root MSE	= 603.12

F値も違ってきます。

不均一分散が存在する場合、通常のF統計量は妥当では無くなります。

robustのオプションを使うと、Stataは不均一分散に対してロバストなF統計量を計算します。

です。

```
. reg pricemy lotsizetb sqrfthb, robust
```

Linear regression

	Number of obs =	88
F(2, 85)	=	32.87
Prob > F	=	0.0000
R-squared	=	0.6631
Adj R-squared	=	0.6552
Root MSE	=	603.12

不均一分散に対してロバストなF統計量については、Wooldridge, J.M. (2010) *Econometric Analysis of Cross Section and Panel Data*, Cambridge, MA: MIT Press を参照してください。

31

- 今の例のように、均一分散を仮定して標準誤差を計算した場合と、不均一分散の可能性を考慮に入れて標準誤差を計算した場合で、検定の結果が変わってしまうことが有ります。

- 必ずではないのですが、後者の標準誤差の方が大きくなる事が多いです。

- よって、変数Aの係数が、

- 普通の標準誤差を使って検定すると、統計的に有意にゼロと異なる。
- しかし、ロバスト標準誤差を使って検定すると、そうは無い。

というケースがよく起こります。

- このような場合、ロバスト標準誤差を使って検定するのが一般的ですから、「変数Aは従属変数に影響を与えない」と結論付けるのが普通ですね。

32

不均一分散の検定

- 「誤差項が不均一分散している」かどうかを検定することができます
 - ロバスト標準誤差を使うのがスタンダードになったため、最近ではこの検定はあまりされなくなりました。
 - とはいえ、いろいろな文脈で誤差項の不均一分散の検定はされることがあります（このコースの範囲外の文脈ですが）。
 - また、回帰分析における有名な検定の一つですので、詳しく解説することになります。

- 次のモデルを考えます：

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + U_i, \quad E(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = 0$$

- 帰無仮説は均一分散です：

$$H_0: \text{VAR}(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = \sigma^2$$

33

- $E(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = 0$ が成り立っているため、

$$\text{VAR}(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = E(U_i^2 | X_{1i}, X_{2i}, \dots, X_{ki})$$

- よって均一分散であるという帰無仮説は

$$H_0: E(U_i^2 | X_{1i}, X_{2i}, \dots, X_{ki}) = \sigma^2$$

と等しくなります。

- また繰り返し期待値の法則より、

$$E(U_i^2) = E(E(U_i^2 | X_{1i}, X_{2i}, \dots, X_{ki})) = E(\sigma^2) = \sigma^2$$

なので

$$H_0: E(U_i^2 | X_{1i}, X_{2i}, \dots, X_{ki}) = E(U_i^2) = \sigma^2$$

ですね。

34

$$\delta H_0: E(U_i^2 | X_{1i}, X_{2i}, \dots, X_{ki}) = E(U_i^2) = \sigma^2$$

- これは、均一分散を検定するためには、「 U_i^2 が説明変数と関係していない」、ということを検定すればよいことを示しています。
- いくつかアプローチがありますが、ここでは最も単純なものを取り上げます。
- まず線形関数を仮定します。

$$U_i^2 = \delta_0 + \delta_1 X_{1i} + \delta_2 X_{2i} + \cdots + \delta_k X_{ki} + V_i$$

ここで V_i は誤差項です。

- この回帰の左辺は、ちょっと変な感じがしますが、もとのモデルの「誤差項の二乗」であることに注意して下さい。
- 帰無仮説（均一分散）は、

$$H_0: \delta_1 = \delta_2 = \cdots = \delta_k = 0$$

35

- $\delta_1 = \delta_2 = \cdots = \delta_k = 0$ であれば、 U_i^2 がいかなる説明変数とも関係していないわけですから、分散は一定（均一分散）ということですね。

- $H_0: \delta_1 = \delta_2 = \cdots = \delta_k = 0$ は、結合帰無仮説ですから、F検定をすればよいです。

- $U_i^2 = \delta_0 + \delta_1 X_{1i} + \delta_2 X_{2i} + \cdots + \delta_k X_{ki} + V_i$

- 左辺は元の式の「誤差項の二乗」です。

- 誤差項は見えませんが、データとしてはありません。

- 実際の検定では、ここを「残差の二乗」で置き換えます。

36

不均一分散の検定のまとめ

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i, E(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = 0$$

$$H_0: \text{VAR}(U_i | X_{1i}, X_{2i}, \dots, X_{ki}) = \sigma^2$$

- ステップ1: モデルを最小二乗法で推定し残差 \hat{U}_i を出す。そして残差の二乗 \hat{U}_i^2 を計算する。
- ステップ2: \hat{U}_i^2 を従属変数、 $X_{1i}, X_{2i}, \dots, X_{ki}$ を説明変数としたモデルを最小二乗法を使って推定する。

$$\hat{U}_i^2 = \delta_0 + \delta_1 X_{1i} + \delta_2 X_{2i} + \dots + \delta_k X_{ki} + V_i$$

- ステップ3: $H_0: \delta_1 = \delta_2 = \dots = \delta_k = 0$ をF検定する。
 - F検定は通常のものでいい。不均一分散に対してロバストにする必要はない。
 - H_0 を棄却できない → 不均一分散を棄却できるだけの十分な証拠はない(均一分散)。
 - H_0 を棄却できる → 均一分散を棄却。不均一分散と結論付ける。

37

やってみましょう

```
. reg price my lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
Total	91785450.8	87	1055005.18	R-squared	=	0.6631
				Adj R-squared	=	0.6552
				Root MSE	=	603.12

price my	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	.3066522 1.2589
sqrfthb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
_cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

- 残差を計算します。predictコマンドを使えばいいです。
- predict 変数名(自分でつけます), r で残差を計算します。
- ここでは残差の変数名をresidとします。predict resid, r
- これの二乗を計算します。変数名はresidsqとします。(ステップ1終わり)

38

- residsqを従属変数、lotsizetbとsqrfthbを説明変数にしたモデルを推定します(ステップ2)。

```
. reg residsq lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs	=	88
Model	7.0715e+12	2	3.5358e+12	F(2, 85)	=	7.93
Residual	3.7888e+13	85	4.4574e+11	Prob > F	=	0.0007
Total	4.4959e+13	87	5.1677e+11	R-squared	=	0.1573
				Adj R-squared	=	0.1375
				Root MSE	=	6.7e+05

residsq	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	729.9321	245.085	2.75	0.007	202.8721 1256.992
sqrfthb	10852.95	4672.65	2.32	0.023	1562.471 20143.43
_cons	-416489.1	260277.4	-1.60	0.113	-933990.3 101012

- lotsizetbとsqrfthbの係数がゼロをF検定(ステップ3)。これは
- 帰無仮説は1%の有意水準で棄却できます → 不均一分散あり、ですね。
- よって普通の標準誤差は適切ではない(ロバスト標準誤差が適切)、ということになります。

39

- 同じことをコマンド一発でできます。モデルを推定後、
estat hettest lotsizetb sqrfthb, fs
です。

```
. reg price my lotsizetb sqrfthb
```

Source	SS	df	MS	Number of obs	=	88
Model	60866850	2	30433425	F(2, 85)	=	83.67
Residual	30918600.8	85	363748.245	Prob > F	=	0.0000
Total	91785450.8	87	1055005.18	R-squared	=	0.6631
				Adj R-squared	=	0.6552
				Root MSE	=	603.12

price my	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lotsizetb	.7827759	.2394666	3.27	0.002	.3066522 1.2589
sqrfthb	49.39334	4.221075	11.70	0.000	41.00071 57.78597
_cons	59.32415	235.1236	0.25	0.801	-408.1645 526.8128

```
. estat hettest lotsizetb sqrfthb, fs
Breusch-Pagan / Cook-Weisberg test for heteroskedasticity
Ho: Constant variance
Variables: lotsizetb sqrfthb
```

F(2, 85)	=	7.93
Prob > F	=	0.0007

40

- これはBreusch-Pagan テストと呼ばれるもののF検定バージョンです。
 - コマンドラインの最後の fs っていうのが、「F検定バージョンをお願いします」と頼んでいるところです。
- アイデアは同じですが、別の検定統計量を使うバージョンもあります。
- Stataは、これら以外の不均一分散の検定も数多く取り揃えています。
- 興味のある人は、estate hettestで検索してみてください。
 - ヘルプ → 検索

41

大切なおまけ

- 不均一分散についての学習は以上で終わります。
- 実証分析上、大切なことは、
 - 「回帰分析するときは、reg の最後に robust オプションをつける」

そして、

「ロバスト標準偏差を使いました」

と論文で説明する、です。

- ただちょっとしたバズルが残りました。
- 「家の価格に土地面積が関係ない？」です。

42

- 「家の価格に土地面積が関係ない」なんてことはありません。
- この結果は、「モデルの関数形が適切なものではなかった」ことから生じたものだと考えられます。
- 家の価格、土地面積、建物面積、すべてに自然対数をとってみましょう。

```
. reg lpricemv llotsizeb lsqftfb, robust
```

```
Linear regression      Number of obs   =    88
                      F(2, 85)         =   58.51
                      Prob > F          =  0.0000
                      R-squared         =  0.6353
                      Root MSE       =   .18547
```

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
lpricemv					
llotsizeb	.1684569	.0382696	4.40	0.000	.0923668 .2445471
lsqftfb	.7623692	.0772034	9.87	0.000	.6088682 .9158703
_cons	4.024582	.3719377	10.82	0.000	3.28507 4.764094

- 土地面積は1%水準で有意ですね。

43

- このように不適切な関数形を使うことで、誤った結論に至ってしまうことが有ります。
- 従って、先行研究をチェックして、どのような関数形が使われているのかを知ることが大切です。
- ちなみに、今回の場合、対数-対数モデルだと、均一分散のようですね。

```
. reg lpricemv llotsizeb lsqftfb
```

```
Source          SS           df           MS       Number of obs   =    88
Model            5.097628         2       2.548814       F(2, 85)        =   74.04
Residual        2.92297433         85      .03439008       Prob > F         =  0.0000
Total           8.02060233         87      .092156348       R-squared        =  0.6353
                                   Root MSE       =   .18547
```

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lpricemv					
llotsizeb	.1684569	.0384596	4.38	0.000	.091989 .2449249
lsqftfb	.7623692	.0808863	9.43	0.000	.6015457 .9231928
_cons	4.024582	.3722772	10.80	0.000	3.283809 4.605356

```
. estat hettest llotsizeb lsqftfb, dx
```

```
Breusch-Pagan / Cook-Weisberg test for heteroskedasticity
Ho: Constant Variance
Variables: llotsizeb lsqftfb

F(2, 85)         =   0.78
Prob > F         =   0.4613
```

44

- もう一つおまけです。対数を取ると単位が関係なくなります。
- 変数のおおもとの単位は\$, square feetでした。これらの変数に対数を取って同じモデルを推定してみましょう。

```

generate lprice = ln(price)
generate llotsize = ln(lotsize)
generate lsqrft = ln(sqrft)

. reg lprice llotsize lsqrft

```

Source	SS	df	MS	Number of obs	=	88
Model	5.09362891	2	2.54681446	F(2, 85)	=	74.04
Residual	2.92397461	85	.034399701	Prob > F	=	0.0000
Total	8.01760352	87	.092156362	R-squared	=	0.6353
				Adj R-squared	=	0.6267
				Root MSE	=	.18547

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
llotsize	.1684569	.0384096	4.38	0.000	.0919889 .2449249
lsqrft	.7623693	.0808863	9.43	0.000	.6015457 .9231929
_cons	-1.640071	.6018805	-2.72	0.008	-2.836771 -.4433717

- 前のページの結果と定数項を除いてすべて一致しました。

45