

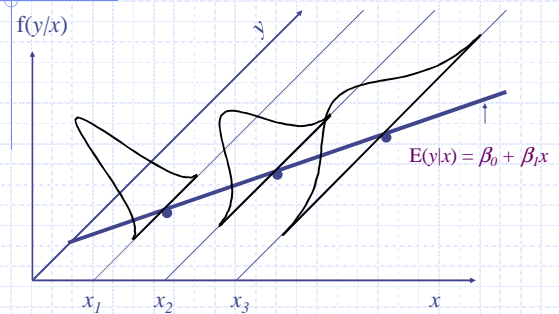
## Heterocedasticidade

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$$

*Outras formas de detectar heterocedasticidade*

1

## Exemplo de heterocedasticidade



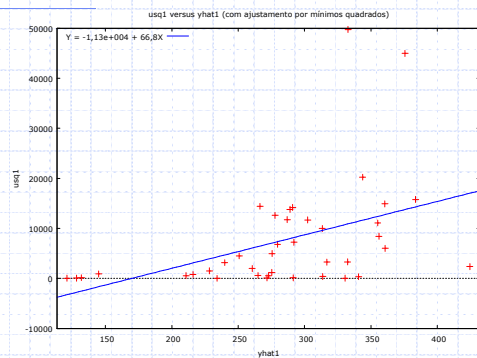
2

## Gráficos residuais

- Estimar o modelo usando MQO e fazer os gráficos dos resíduos relacionando-os com as variáveis explicativas ou com o valor estimado de  $y$ .
- Se erros são heterocedásticos – relação sistemática.

3

## Gráficos residuais



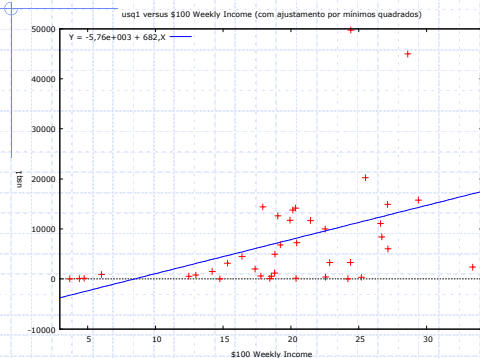
4

## Gráficos residuais

- No exemplo, o exame dos resíduos ao quadrado (proxies dos termos de erro), é feito quando eles são plotados junto com o  $y$  estimado.
- Verifica-se que o valor médio estimado de  $y$  se relaciona de forma sistemática com os resíduos elevados ao quadrado.

5

## Gráficos residuais



6

## Gráficos residuais

- ◆ No exemplo acima, o exame dos resíduos ao quadrado é feito quando eles são plotados junto com a variável explicativa  $x$ .
- ◆ Verifica-se que o valor médio estimado de  $y$  se relaciona de forma sistemática com os resíduos elevados ao quadrado.
- ◆ Para RLS, o gráfico é idêntico ao anterior. Para RLM, plota-se para cada uma das variáveis explicativas.

7

## Teste Goldfeld-Quandt

- ◆ Separar a amostra em duas subamostras aproximadamente iguais.
- ◆ Uma amostra com variância grande  $\sigma_1$  e outra com variância pequena  $\sigma_2$ .
- ◆ Calcular a variância estimada do erro para cada subamostra.
- ◆ Se a hipótese nula de variâncias iguais não é verdade, esperamos que a razão de  $\sigma_1/\sigma_2$  seja grande.

8

## Teste GQ

$$GQ = \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2}$$

$$H_0 : \sigma_t^2 = \sigma^2$$

$$H_1 : \sigma_t^2 = \sigma^2 x_t$$

- Rejeito a hipótese nula quando  $GQ > F_c$ .
- $F_c$  é o valor crítico da distribuição F com (T1-K) e (T2-K) graus de liberdade.

- **Exemplo:**  $GQ = 2285,9/682,46 = 3,35 > F_c = 2,22$ .
- Rejeitamos a hipótese nula e concluímos que há heterocedasticidade.

9

## Mínimos quadrados ponderados

- ◆ O estimador de *White* para os erros padrão ajuda a melhorar o problema de inferências incorretas das estimativas de MQO com heterocedasticidade.
- ◆ Podemos obter um estimador melhor do que o MQO, mas depende de como modelamos a heterocedasticidade.

10

## Mínimos quadrados ponderados (cont.)

- ◆ A idéia é minimizar a soma dos quadrados ponderados por  $1/h_i$ .
- ◆ Dá-se menos peso para as observações com maior variância.
- ◆ **No gretl temos que informar esta função que pondera as observações.**

11

## Exemplo

- ◆ Banco de dados: food.gdt
- ◆ Relacionar gastos em alimentos com a renda mensal
- ◆ Definir a variável peso:  $1/x$  (inverso da renda).
- ◆ Cada variável, inclusive a constante, é multiplicada pela raiz quadrada do peso.
- ◆ Modelo: outros modelos lineares/mínimos quadrados ponderados

12

## Exemplo

Estimativas WLS  
Variável dependente: y

	Modelo 1
const	78,68** (23,79)
x	10,45** (1,386)
n	40
R <sup>2</sup>	0,5994
lnL	-172,980

Erros padrão entre parênteses  
\* indica significância num nível de 10 por cento  
\*\* indica significância num nível de 5 por cento

13

## MQGF

- ◆ Heterocedasticidade tem uma forma específica.

- ◆ *Heterocedasticidade multiplicativa:*

$$\sigma_i^2 = \exp(\alpha_1 + \alpha_2 z_i)$$

- ◆ A variável explicativa  $z$  determina como a variância muda para cada observação. Defina  $z_i = \ln(x_i)$
- ◆ Estima os parâmetros por MQO e usa o valor predito para calcular o peso.

14

## MQGF (cont.)

- ◆ A estimativa de  $h$  é obtida por  $\hat{h} = \exp(\hat{g})$ ; o peso será o inverso dessa estimativa.
- ◆ Resumindo:
- ◆ Faça a regressão por MQO da equação original, salve os resíduos,  $\hat{u}$ , eleve-os ao quadrado e tire o log.
- ◆ Faça a regressão de  $\ln(\hat{u}^2)$  em todas as variáveis independentes e obtenha o valor ajustado  $\hat{g}$ .
- ◆ Faça a regressão por MQP utilizando  $1/\exp(\hat{g})$  como ponderador.

15

## Exemplo

$$\sigma_i^2 = \exp(\alpha_1 + \alpha_2 z_i)$$

$$z_i = \ln(x_i)$$

$$\ln(\sigma_i^2) = \alpha_1 + \alpha_2 z_i + v_i$$

$$\ln(\hat{e}_i^2) = \alpha_1 + \alpha_2 z_i + v_i$$

16

## Exemplo

Modelo 1: Mínimos Quadrados (OLS), usando as observações 1-40  
Variável dependente: y

	Coefficiente	Erro Padrão	razão-t	p-valor	*
const	83,416	43,4102	1,9216	0,06218	*
x	10,2096	2,09326	4,8774	0,00002	***
Média var. dependente	283,5735	D.P. var. dependente	112,6752		
Soma resid. quadrados	304505,2	E.P. da regressão	89,51700		
R-quadrado	0,385002	R-quadrado ajustado	0,368818		
F(1, 38)	23,78884	P-valor(F)	0,000019		
Log da verossimilhança	-235,5088	Critério de Akaike	475,0176		
Critério de Schwarz	478,3954	Critério Hannan-Quinn	476,2389		

Salva os resíduos ao quadrado e tira o log destes.

17

## Exemplo

$$\ln(\hat{e}_i^2) = \alpha_1 + \alpha_2 z_i + v_i$$

- ◆ Tira exponencial dos valores preditos (acha  $h$ ).
- ◆ Ache o peso ( $w=1/h$ )

Modelo 2: Mínimos Quadrados (OLS), usando as observações 1-40  
Variável dependente: l\_usq2

	Coefficiente	Erro Padrão	razão-t	p-valor	*
const	0,937796	1,58311	0,5924	0,55711	
l_x	2,32924	0,541336	4,3028	0,00011	***
Média var. dependente	7,648159	D.P. var. dependente	2,071519		
Soma resid. quadrados	112,5310	E.P. da regressão	1,720855		
R-quadrado	0,327597	R-quadrado ajustado	0,309903		
F(1, 38)	18,51375	P-valor(F)	0,000114		
Log da verossimilhança	-77,44452	Critério de Akaike	158,8890		
Critério de Schwarz	162,2668	Critério Hannan-Quinn	160,1103		

18

## Exemplo

$$\hat{y} = 76,0538 + 10,6335x$$

(9,7135)                      (0,97153)

$T = 40 \quad R^2 = 0,7529 \quad F(1,38) = 119,8 \quad \hat{\sigma} = 1,5467$

### Rode MQP

Modelo 3: WLS, usando as observações 1-40  
Variável dependente: y  
Variável usada como peso: w

	Coefficiente	Erro Padrão	razão-t	p-valor	
const	76,0538	9,71349	7,8297	<0,00001	***
x	10,6335	0,971514	10,9453	<0,00001	***

Estadísticas baseadas nos dados ponderados:

Soma resíd. quadrados	90,91135	E.P. da regressão	1,546740
R-quadrado	0,759187	R-quadrado ajustado	0,752850
F(1, 38)	119,7991	P-valor(F)	2,62e-13
Log da verossimilhança	-73,17765	Critério de Akaike	150,3553
Critério de Schwarz	153,7331	Critério Hannan-Quinn	151,5766

Estadísticas baseadas nos dados originais:

Média var. dependente	283,5735	D.P. var. dependente	112,6752
Soma resíd. quadrados	304869,6	E.P. da regressão	89,57055

19

## Exemplo

Variável dependente: y

	Modelo 1 WLS h=x	Modelo 2 WLS h=exp(g)	Modelo 3 Mínimos Quadrados (OLS)
const	78,68** (23,79)	76,05** (9,713)	83,42* (43,41)
x	10,45** (1,386)	10,63** (0,9715)	10,21** (2,093)
n	40	40	40
R <sup>2</sup>	0,5994	0,7592	0,3850
lnL	-172,980	-73,178	-235,509

Erros padrão entre parênteses

\* indica significância num nível de 10 por cento

\*\* indica significância num nível de 5 por cento

20

## Observações sobre MQP

- ◆ Lembre-se que utiliza-se MQP apenas por eficiência, pois MQO continua não tendencioso e consistente.
- ◆ As estimativas serão diferentes devido a erros amostrais, mas se forem muito diferentes, então alguma outra hipótese de Gauss-Markov também deve estar sendo violada.

21

## Multicolinearidade

22

## Multicolinearidade

- ◆ Quando existem relação linear exata entre as variáveis independentes será impossível calcular os estimadores de MQO.
- ◆ O procedimento MQO utilizado para estimação não será efetivado.
- ◆ Mensagem: "matriz quase singular" (uma matriz quase singular  $X'X$  não pode ser invertida) ou "a variável  $x_k$  dropped".
- ◆ **Relacionamento linear exato:** só quando os dados foram construídos pelo pesquisador, pe., no caso de inclusão de dummies.
- ◆ **Relacionamento linear aproximado** entre as variáveis independentes: comuns em economia.
- ◆ O procedimento de estimação não é rompido quando as variáveis são bastante correlacionadas, contudo, surgem problemas de estimação.

## Multicolinearidade

- ◆ Multicolinearidade: nome dado ao fenômeno de presença de relação linear aproximada entre os regressores.
- ◆ Problema de estimação causado por uma amostra particular. Não é um problema teórico.
- ◆ Multicolinearidade nos dados pode existir por diferentes motivos:
  - Regressores possuem a mesma tendência temporal.
  - Algumas variáveis caminham na mesma direção porque os dados não foram coletados de uma base grande.
  - Pode existir realmente algum tipo de relacionamento aproximado entre os regressores.

## Variância do estimador de MQO

A variância estimada de  $b_k$  é

$\text{Var}[b_k|X] =$

$$\frac{s^2}{(1-R_k^2) \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2} = \frac{s^2}{(1-R_k^2) S_{22}}$$

Quanto maior o fit da regressão de  $x_2$  em  $X_1$ , maior a variância. No limite, um ajuste perfeito produz uma variância infinita.

## Variância do estimador de MQO Forma mais geral

Defina a matriz  $X$  que contém uma constante e  $K-1$  variáveis explicativas

A variância estimada de  $b_k$  é

$\text{Var}[b_k|X] =$

$$\frac{s^2}{(1-R_k^2) \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}$$

**Ingrediente para existência de multicolinearidade:**

- Quanto maior a correlação entre  $x_k$  e as outras variáveis ( $R_k^2$ ).

## Consequências da Multicolinearidade

- ◆ O estimador de MQO permanece **não viesado** e **BLUE**.
- ◆ O grau de ajuste não é afetado.
- ◆ Problemas práticos:
  - Pequenas mudanças nos dados produzem grandes variações nas estimativas dos parâmetros.
  - Os coeficientes estimados apresentam erros padrão muito elevados e baixos níveis de significância, mesmo que sejam conjuntamente significativos e com o grau de ajuste da regressão elevado ( $R^2$ ).
  - Os coeficientes podem ter o sinal "errado" e magnitudes irreais.

## Consequências da Multicolinearidade

- ◆ Na presença de multicolinearidade, o procedimento de estimação MQO não recebe variação independente suficiente de uma variável para realizar o cálculo com confiança do efeito que esta tem sobre a variável dependente.
- ◆ Quando os regressores são altamente correlacionados, a maior parte da sua variação é comum às duas variáveis, deixando pouca variação exclusiva a cada variável.
- ◆ MQO tem pouca informação para usar ao fazer as estimativas do coeficiente (similar a um problema de amostra pequena ou que a variável não mudasse muito).

## Consequências da Multicolinearidade

- ◆ As variâncias dos estimadores MQO dos parâmetros são muito grandes – **Imprecisão dos estimadores dos parâmetros**.
- ◆ **Erros de especificação** : não sabemos qual variável é mais ou menos importante para explicar a variação da variável dependente.

## Como detectar?

- ◆ Controvérsia: muitos métodos inadequados.
- ◆ Sinais hipotéticos não são encontrados.
- ◆ Variáveis consideradas *a priori* importantes não são significativas individualmente, mas estatística F (significância coletiva) é alta.
- ◆ Resultados alterados quando uma variável independente é excluída ou quando uma observação é retirada.
- ◆ Matriz de correlação (0,8 a 0,9 são valores absolutos altos): detecta colinearidade de duas variáveis, mas não de mais de duas.

## Como detectar?

### Índice de condição dos dados (IC):

- Raiz quadrada da razão da maior para a menor raiz característica de  $X'X$

$$\gamma = \left[ \frac{\text{raizmáxima}}{\text{raizmínima}} \right]^{1/2}$$

- Medida de sensibilidade das estimativas a pequenas perturbações dos dados.
- Medida de proximidade de  $X'X$  da singularidade (multicolinearidade perfeita): quanto maior o IC maior dificuldade em inverter a matriz.
- Índice maior que 20 indica colinearidade forte: mudança de 1% nos dados faz surgir uma mudança de IC% nos estimadores.

## Como detectar?

### Inverso da matriz de correlação:

- Elementos na diagonal: Fatores de inflação da variância (VIF).

$$VIF = \frac{1}{(1 - R_k^2)}$$

$R^2$  da regressão da k-ésima variável independente em todas demais variáveis independentes.

- Quanto maior VIF, mais o  $R_k^2$  está perto da unidade.
- Medida da quantidade pela qual a variância da k-ésima estimativa do coeficiente é aumentada devido a associação linear com as outras variáveis explicativas.
- Se  $VIF > 10$ : presença de colinearidade

## No stata:

```
. reg ln_sal_hora filho idade idade2 sexo educa
Source | SS      df      MS              Number of obs = 14537
-----+-----+-----+-----+----- F( 5, 14531) = 1939.23
Model | 5434.055   5     1086.811       Prob > F       = 0.0000
Residual | 8143.68463 14531  560.435251     R-squared      = 0.4002
Total | 13577.7496 14536  934.077458     Adj R-squared  = 0.4000
                                           Root MSE     = .74862
```

ln_sal_hora	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
filho	-.208508	.0202922	-10.32	0.000	-.2492833	-.1697328
idade	.0604502	.0028366	21.31	0.000	.05469	.0660103
idade2	-.0006105	.0000332	-15.37	0.000	-.0005796	-.0004454
sexo	-.246604	.0129488	-26.77	0.000	-.3739854	-.1212227
educa	.1304724	.0014665	88.97	0.000	.1275979	.1333469
_cons	-.4814204	.061482	-7.83	0.000	-.601933	-.3609078

```
. vif
Variable |      VIF      1/VIF
-----+-----+-----
idade    | 33.37    0.029969
idade2   | 30.63    0.032650
filho    | 1.63     0.613227
educa    | 1.11     0.901243
sexo     | 1.04     0.961969
-----+-----+-----
Mean VIF | 13.56
```

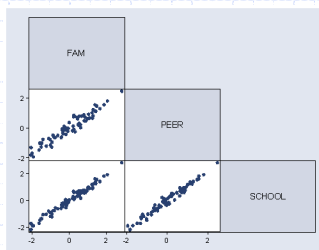
## No stata:

```
. collin idade idade2 sexo educa filho
Collinearity Diagnostics
-----+-----+-----+-----+-----
Variable | VIF      SQRT VIF      Tolerance      Squared
-----+-----+-----+-----+-----
idade    | 33.37    5.77    0.0300    0.9698
idade2   | 30.63    5.53    0.0300    0.9500
sexo     | 1.04     1.02    0.9943    0.0157
educa    | 1.11     1.05    0.6449    0.3551
filho    | 1.63     1.26    0.3218    0.6782
-----+-----+-----+-----+-----
Mean VIF | 10.66
```

Eigenval	Cond Index
1	4.3513
2	1.0883
3	0.3723
4	0.1404
5	0.0395
6	0.0063

Condition Number: 26.3514  
Eigenvalues & Cond Index computed from scaled raw sscp (w/ intercept).  
Det(corrrelation matrix) 0.0194

## No stata: graph matrix fam peer school



## Multicolinearidade

Não existe "cura" para a colinearidade.

1. Exclusão de variáveis: eliminar as variáveis que causam o problema – impor na regressão a hipótese de que a variável problemática não deve aparecer no modelo. **Possível problema de especificação.**
2. Obtenção de mais dados: dados adicionais e tamanho da amostra.
3. Formalizar os relacionamentos entre os regressores: equações simultâneas.
4. Especificar o relacionamento entre alguns parâmetros: dois parâmetros iguais ou que a soma das elasticidades deve ser igual a um, etc.
5. Análise componente principal: as variáveis colineares poderiam ser agrupadas para formar um índice composto capaz de representar este conjunto de variáveis. Variável só pode ser criada se tiver uma interpretação económica.