

# Métodos Quantitativos

## Aula 11. Regressão linear, parte 3

**Pedro H. G. Ferreira de Souza**

**[pedro.ferreira@ipea.gov.br](mailto:pedro.ferreira@ipea.gov.br)**

Mestrado Profissional em Políticas Públicas e Desenvolvimento

Instituto de Pesquisa Econômica Aplicada (Ipea)

12 dez. 2022

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

# Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

# Regressão multivariada

## O que vimos

1. Estimação no R e interpretação dos coeficientes
2. Viés de variável omitida
3. Inferência, testes de hipóteses e ICs
4. Pressupostos de OLS multivariada
5. Comparação e seleção de modelos

Recapitulação

## Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

## O que ficou de fora?

Até agora usamos só **Variáveis quantitativas**, assumindo linearidade e separabilidade dos efeitos parciais.

Mas o que acontece se...

- Quisermos usar **variáveis qualitativas ou discretas**?
- Os efeitos parciais não forem **lineares**?
- Os efeitos parciais dependerem de **interações** entre variáveis?

... e muito mais.

# Objetivos de hoje

1. OLS com variáveis qualitativas
  - Variáveis binárias e multinomiais
2. Transformações em variáveis
  - Logaritmos e polinômios
3. Interações entre variáveis
4. Mais sobre valores preditos e resíduos
5. Cuidados no mundo real

# Pacotes

```
# Pacotes de uso geral
```

```
library(tidyverse)
```

```
library(broom)
```

```
# Pacote para estatísticas descritivas
```

```
library(summarytools)
```

```
# Pacote para visualizacao de resultados de modelos
```

```
library(modelsummary)
```

```
# Pacote para ler arquivos csv
```

```
library(readr)
```

```
# Pacotes com bases de dados
```

```
library(HistData)
```

```
library(causaldata)
```



## Bases de dados #1 e #2

```
# Colera
```

```
colera <- Cholera
```

```
colera$water <-
```

```
  factor(colera$water,
```

```
        levels = c('New River', 'Kew', 'Battersea'))
```

```
# Restaurantes no Alasca
```

```
restaurantes <- restaurant_inspections
```

## Bases de dados #3

### Download

Baixem o arquivo `pnadc2021_limpa.csv` na [página da aula 11 no GitHub](#).

- Arquivo contém seleção de variáveis (renomeadas) da PNADC 2021, do IBGE, com filtro para manter apenas indivíduos ocupados que trabalham pelo menos 30 horas semanais.
- São 109,818 observações e 8 variáveis

### Leitura

```
local <- file.path("D:", "OneDrive", "work", "Incompletos",  
                  "mestrado_ipea", "11-regressao-linear-pt",  
                  "dados", "pnadc2021_limpa.csv")  
  
pnadc <- read_csv(local)
```

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

## Como fazer?

Para variáveis independentes com  $k$  categorias, incluímos no modelo  $k - 1$  **variáveis binárias** ou *dummy*, omitindo uma categoria de referência.

- O coeficiente de uma variável *dummy* oferece comparação direta entre aquela categoria e a categoria de referência

### Exemplo

Suponha uma variável  $z = \{0, 1, 2\}$  e as variáveis indicadoras  $z_1 = (z == 1)$  e  $z_2 = (z == 2)$ , com o modelo  $E(y) = \alpha + \beta_1 z_1 + \beta_2 z_2$ .

- Quando  $z_1 == 1$ ,  $E(y) = \alpha + \beta_1 \rightarrow$  média de  $y$  para esse grupo (e vice-versa para  $z_2$ )
- Teste t dos coeficientes avalia  $H_0$  de que média do grupo é igual à da categoria de referência
- Teste F da regressão avalia  $H_0$  conjunta de que todos os grupos têm a mesma média

## Exemplo com colera

```
mean(colera$cholera_drates)
```

```
freq(colera$water)
```

```
## [1] 66.13158
```

```
## Frequencies
```

```
## colera$water
```

```
## Type: Factor
```

```
##
```

```
##           Freq  % Valid  % Valid Cum.  % Total  %
```

```
## -----
```

```
-----
```

```
##      New River    20    52.63    52.63    52.63
```

```
##           Kew      6    15.79    68.42    15.79
```

```
##      Battersea   12    31.58   100.00    31.58
```

```
##           <NA>     0      0.00      0.00
```

```
##      Total     38   100.00   100.00   100.00
```

## Exemplo com colera

Como water já está formatada como variável factor com três categorias, basta incluir no comando da regressão:

```
col1 <- lm(cholera_drate ~ water, data = colera)
summary(col1)
```

Podemos comparar modelos também:

```
col2 <- lm(cholera_drate ~ water + elevation, data = colera)
col3 <- lm(cholera_drate ~ water + elevation + pop_dens,
           data = colera)
col4 <- lm(cholera_drate ~ water + elevation + pop_dens +
           house_valpp, data = colera)
msummary(list(col1, col2, col3, col4), output = "markdown",
          stars = TRUE, gof_omit = c('BIC|AIC|Log|Num. Obs|F'))
```

*(Output nos próximos slides)*

```
##
## Call:
## lm(formula = cholera_drate ~ water, data = colera)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -94.25 -16.80  -4.30   22.84   82.75
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    47.800     7.969   5.998 7.76e-07 ***
## waterKew       -32.800    16.589  -1.977  0.0559 .
## waterBattersea  74.450    13.013   5.721 1.80e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 35.64 on 35 degrees of freedom
## Multiple R-squared:  0.5752, Adjusted R-squared:  0.551
## F-statistic: 23.7 on 2 and 35 DF, p-value: 3.11e-07
```

	Model 1	Model 2	Model 3	Model 4
(Intercept)	47.800*** (7.969)	53.364*** (9.416)	46.863** (14.410)	64.503** (18.004)
waterKew	-32.800+ (16.589)	-24.403 (18.212)	-21.938 (18.839)	-8.259 (20.373)
waterBattersea	74.450*** (13.013)	69.595*** (13.703)	72.802*** (14.829)	68.532*** (14.758)
elevation		-0.133 (0.121)	-0.125 (0.123)	-0.159 (0.122)
pop_dens			0.045 (0.075)	0.050 (0.073)
house_valpp				-2.945 (1.869)
Num.Obs.	38	38	38	38
R2	0.575	0.590	0.594	0.623
R2 Adj.	0.551	0.554	0.545	0.565
RMSE	34.20	33.61	33.43	32.20



## Codificação de variáveis categóricas

Variáveis categóricas nem sempre estão codificadas como factor, mas isso não é um problema.

- Se a variável já for uma *dummy* numérica ou character, não há problema; caso contrário, recodificar com `mutate()` e `factor()` ou só pedir para o R considerar como factor na regressão.

### Restaurantes

```
glimpse(restaurantes)
```

```
## Rows: 27,178
```

```
## Columns: 5
```

```
## $ business_name      <chr> "MCGINLEYS PUB", "VILLAGE INN #1",
```

```
## $ inspection_score    <dbl> 94, 86, 80, 96, 83, 95, 94, 100, 92
```

```
## $ Year                <dbl> 2017, 2015, 2016, 2003, 2017, 2008,
```

```
## $ NumberofLocations   <dbl> 9, 66, 79, 86, 53, 89, 28, 37, 109,
```

```
## $ Weekend             <lgl> FALSE  FALSE  FALSE  FALSE  FALSE  FALSE
```

## Exemplo com restaurantes

```
table(restaurantes$Weekend)
print(c(min(restaurantes$Year), max(restaurantes$Year)))

##
## FALSE   TRUE
## 26968   210
## [1] 2000 2019
```

## Exemplo com restaurantes

```
table(restaurantes$Weekend)
print(c(min(restaurantes$Year), max(restaurantes$Year)))

##
## FALSE   TRUE
## 26968   210
## [1] 2000 2019

res1 <- lm(inspection_score ~ NumberofLocations + Weekend,
           data = restaurantes)
res2 <- lm(inspection_score ~ NumberofLocations + Weekend +
           as.factor(Year), data = restaurantes)
msummary(list(res1, res2), output = "markdown",
           stars = TRUE, statistic = NULL,
           gof_omit = c('BIC|AIC|Log|Num. Obs|F'))
```

	Model 1	Model 2
(Intercept)	94.851***	96.808***
NumberofLocations	-0.019***	-0.018***
WeekendTRUE	1.498***	1.328**
as.factor(Year)2001		0.578*
as.factor(Year)2002		0.522+
as.factor(Year)2003		1.091***
as.factor(Year)2004		-1.506***
as.factor(Year)2005		-1.980***
as.factor(Year)2006		-3.507***
as.factor(Year)2007		-2.873***
as.factor(Year)2008		-2.882***
as.factor(Year)2009		-3.470***
as.factor(Year)2015		-0.896***
as.factor(Year)2016		-2.520***
as.factor(Year)2017		-2.121***
as.factor(Year)2018		-2.654***
as.factor(Year)2019		-1.101***
Num.Obs.	27178	27178
R2	0.065	0.105
R2 Adj.	0.065	0.105

## Exemplo com a PNADC

```
freq(pnadc$regiao, headings = FALSE, report.nas = FALSE)
```

```
##
##              Freq              %      % Cum.
## -----
##           1    14427      13.14     13.14
##           2    26322      23.97     37.11
##           3    32844      29.91     67.01
##           4    23030      20.97     87.98
##           5    13195      12.02    100.00
##        Total  109818     100.00    100.00
```

```
summary( lm(renda_trabalho ~ factor(regiao), data = pnadc))
```

```
##
## Call:
## lm(formula = renda_trabalho ~ factor(regiao), data = pnadc)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2796  -1346   -750     70  297354
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2020.00      28.36   71.221  <2e-16 ***
## factor(regiao)2  -352.36      35.29   -9.985  <2e-16 ***
## factor(regiao)3   625.56      34.03   18.385  <2e-16 ***
## factor(regiao)4   729.73      36.17   20.174  <2e-16 ***
## factor(regiao)5   778.26      41.04   18.965  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3407 on 109813 degrees of freedom
## Multiple R-squared:  0.01772    Adjusted R-squared:  0.01766
```

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

**Transformações em variáveis**

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

# Por que transformar variáveis?

## Linearidade

Em MQO, a linearidade diz respeito aos parâmetros estimados: efeitos parciais dos coeficientes são lineares, ou seja, aumento de uma unidade em  $x$  aumenta a média condicional predita em  $\beta$ .

MQO é muito flexível porque pode incorporar **não linearidades** a partir de transformações em variáveis.

As duas transformações mais comuns são:

- Aplicação de logaritmos no  $y$  ou em algum  $x$
- Inclusão de polinômios em algum  $x$



# Logaritmos (i)

Por que aplicar?

- Às vezes a teoria prevê efeitos multiplicativos;
- Em muitos casos, logaritmos geram melhor aderência aos pressupostos do modelo clássico;
- Logaritmos são menos sensíveis a *outliers*
- etc

Uma boa **regra de bolso** é aplicar logs em variáveis monetárias (quando todos os valores são positivos).

Interpretação de logs

- Se  $y$  ou algum  $x$  está em log, interpretação passa a ser em termos de mudanças percentuais (ver próximos slides)

# Logaritmos (ii)

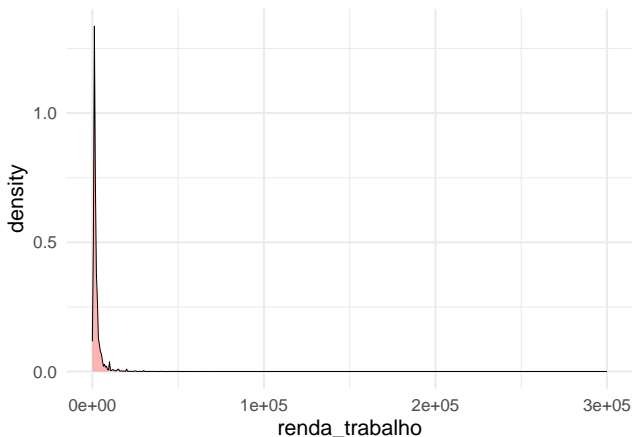
## Interpretação

Modelo	Dependente	Independente	Interpretação de $\beta_1$
Nível-nível	$y$	$x$	$\Delta y = \beta_1$ para $\Delta x = 1$
Nível-log	$y$	$\log(x)$	aprox: $\Delta y = \beta_1/100$ para $\Delta x = 1\%$ , exato: $\Delta y = \beta_1 \log(1 + c)$ para $\Delta x = c\%$
Log-nível	$\log(y)$	$x$	aprox: $\Delta y = 100\beta_1\%$ para $\Delta x = 1$ , exato: $\Delta y = 100(e^{\beta_1} - 1)\%$ para $\Delta x = 1$
Log-Log	$\log(y)$	$\log(x)$	aprox: $\Delta y = \beta_1\%$ para $\Delta x = 1\%$ , exato: $\Delta y = 100((1 + c)^{\beta_1} - 1)\%$ para $\Delta x = c\%$

*(Na prova vou pedir para vocês usarem as interpretações aproximadas!)*

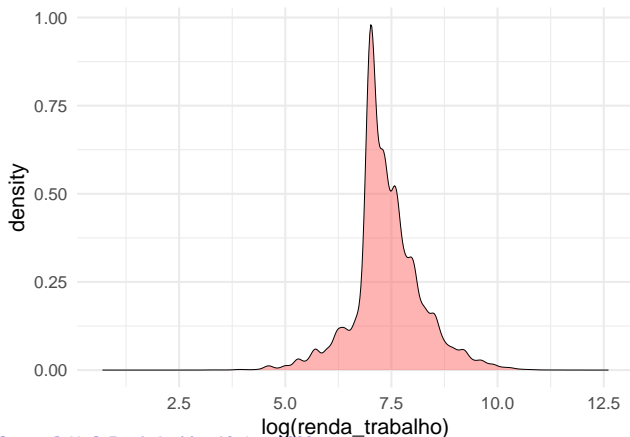
## Renda do trabalho na PNADC

```
ggplot(data = pnadc, aes(x = renda_trabalho)) +  
  geom_density(bw = .1, alpha = .5, fill = 'indianred1') +  
  theme_minimal(base_size = 24)
```



## Log Renda do trabalho na PNADC

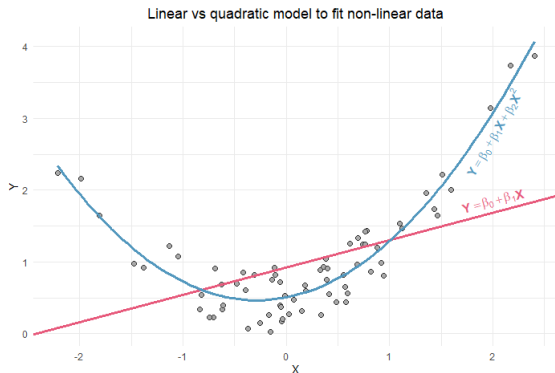
```
ggplot(data = pnadc, aes(x = log(renda_trabalho))) +  
  geom_density(bw = .1, alpha = .5, fill = 'indianred1') +  
  theme_minimal(base_size = 24)
```



# Polinômios

É comum utilizarmos polinômios - principalmente termos quadráticos - para incorporar não linearidades nas variáveis. Por exemplo, podemos estimar  $y = \alpha + \beta_1 x + \beta_2 x^2 + \epsilon$ .

- Efeito marginal sobre  $y$  varia conforme o valor de  $x$ :  $\Delta y \approx (\beta_1 + 2\beta_2 x)\Delta x$



## Exemplo: PNADC

```
bra1 <- lm(log(renda_trabalho) ~ factor(regiao) + homem + anos_estudo +
           poly(idade, degree = 2), data = pnadc)
msummary(bra1, output = "markdown", stars = TRUE, statistic = NULL)
```

	Model 1
(Intercept)	6.015***
factor(regiao)2	-0.166***
factor(regiao)3	0.201***
factor(regiao)4	0.327***
factor(regiao)5	0.293***
homem	0.274***
anos_estudo	0.100***
poly(idade, degree = 2)1	67.605***
poly(idade, degree = 2)2	-26.723***
Num.Obs.	109818
R2	0.348
R2 Adj.	0.348
AIC	1841859.0

## Exemplo: cólera

```
col5 <- lm(log(cholera_deaths) ~ water + log(popn),  
           data = colera)  
col6 <- lm(log(cholera_deaths) ~ water + log(popn) +  
           poly(pop_dens, degree = 2), data = colera)  
msummary(list(col5, col6),  
          output = "markdown", stars = TRUE,  
          gof_omit = c('BIC|AIC|Log|Num. Obs|F'))
```

*(Output no próximo slide).*

## Exemplo: cólera

	Model 1	Model 2
(Intercept)	-5.890** (1.791)	-5.840** (1.923)
waterKew	-1.192*** (0.259)	-1.128*** (0.276)
waterBattersea	0.923*** (0.204)	0.966*** (0.232)
log(popn)	1.040*** (0.163)	1.033*** (0.174)
poly(pop_dens, degree = 2)1		0.514 (0.606)
poly(pop_dens, degree = 2)2		-0.239 (0.599)
Num.Obs.	38	38
R2	0.755	0.762
R2 Adj	0.733	0.725



Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

**Interações entre variáveis**

Mais sobre valores preditos e resíduos

Cuidados no mundo real

Encerramento do semestre

## Do que se trata?

O modelo que vimos até agora foi:

$$E(y) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

Por definição, esse modelo pressupõe que os efeitos parciais de cada regressor **não dependem** dos valores dos outros regressores.

Às vezes, esse modelo é inadequado por ser simples demais. Podemos enriquecê-lo incluindo **interações** entre os regressores, de modo que a relação entre um dado  $x$  e o  $y$  podem variar conforme o nível de outra variável independente

- Por exemplo, o retorno salarial por anos de estudo pode variar conforme o gênero, idade ou região

## Interação entre variáveis quantitativas

Abordagem mais comum é via inclusão do **produto cruzado** para a interação entre variáveis:

$$E(y) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

Equivale a criar uma variável “artificial”  $x_3 = x_1 x_2$ . Observem que podemos reescrever:

- $E(y) = (\alpha + \beta_2 x_2) + (\beta_1 + \beta_3 x_2) x_1 = \alpha' + \beta' x_1$
- $\alpha' = \alpha + \beta_2 x_2$
- $\beta' = \beta_1 + \beta_3 x_2$

A média de  $y$  continua uma função linear de  $x_1$ , mas a inclinação do efeito parcial varia conforme  $x_2$  muda  $\rightarrow \beta_1$  equivale ao efeito parcial de  $x_1$  somente quando  $x_2 = 0$

## Outras interações

É ainda mais comum usarmos interações entre variáveis categóricas e variáveis contínuas:

$$E(y) = \alpha + \beta_1 d + \beta_2 x_1 + \beta_3 dx_1$$

No modelo, o coeficiente  $\beta_1$  desloca o intercepto para o grupo em que  $d = 1$ , enquanto o coeficiente  $\beta_3$  permite que a inclinação da reta (efeito parcial) também varie entre grupos.

- Na prática, o teste de hipóteses para  $H_0 : \beta_3 = 0$  avalia se o efeito parcial de fato varia ou não

## Exemplo: PNADC

```
bra2 <- lm(log(renda_trabalho) ~ factor(regiao) + homem +  
           branco + anos_estudo + poly(idade, degree = 2),  
           data = pnadc)  
bra3 <- lm(log(renda_trabalho) ~ factor(regiao) + branco +  
           homem * anos_estudo + poly(idade, degree = 2),  
           data = pnadc)  
bra4 <- lm(log(renda_trabalho) ~ factor(regiao) + homem +  
           branco * anos_estudo + poly(idade, degree = 2),  
           data = pnadc)  
bra5 <- lm(log(renda_trabalho) ~ factor(regiao) + anos_estudo +  
           branco * homem + poly(idade, degree = 2),  
           data = pnadc)  
msummary(list(bra2, bra3, bra4, bra5), output = "markdown",  
          stars = TRUE)
```

	Model 1	Model 2	Model 3	Model 4
(Intercept)	6.013*** (0.008)	5.878*** (0.012)	6.080*** (0.009)	6.010*** (0.009)
factor(regiao)2	-0.174*** (0.007)	-0.174*** (0.007)	-0.176*** (0.007)	-0.174*** (0.007)
factor(regiao)3	0.166*** (0.007)	0.170*** (0.007)	0.164*** (0.007)	0.166*** (0.007)
factor(regiao)4	0.252*** (0.007)	0.256*** (0.007)	0.252*** (0.007)	0.252*** (0.007)
factor(regiao)5	0.275*** (0.008)	0.277*** (0.008)	0.275*** (0.008)	0.275*** (0.008)
homem	0.275*** (0.004)	0.459*** (0.013)	0.275*** (0.004)	0.279*** (0.006)
branco	0.135*** (0.004)	0.134*** (0.004)	-0.035** (0.012)	0.140*** (0.007)
anos_estudo	0.097*** (0.0005)	0.108*** (0.0009)	0.091*** (0.0007)	0.097*** (0.0005)
poly(idade, degree = 2)1	65.436*** (0.680)	65.188*** (0.679)	65.324*** (0.679)	65.454*** (0.680)
poly(idade, degree = 2)2	-27.321*** (0.653)	-27.364*** (0.653)	-27.040*** (0.653)	-27.312*** (0.653)
homem × anos_estudo		-0.016*** (0.001)		
branco × anos_estudo			0.015*** (0.001)	
branco × homem				-0.008 (0.008)
Num.Obs.	109818	109818	109818	109818
R2	0.354	0.355	0.355	0.354

## Mais sobre interações

1. Sempre incluam no modelo os efeitos principais quando forem usar interações
  - Ou seja, o modelo deve ser  $E(y) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$ , e **não**  $E(y) = \alpha + \beta x_1 x_2$
2. Na prática, termos quadráticos equivalem à interação de uma variável com ela mesma
3. Nada impede a inclusão de interações entre três (ou mais) variáveis, mas a interpretação fica muito complicada
  - É tentador testar todas as interações possíveis, mas não se esqueçam de que modelos são úteis justamente por simplificarem o mundo

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

**Mais sobre valores preditos e resíduos**

Cuidados no mundo real

Encerramento do semestre



## Salvando valores preditos

A “reta” da regressão gera valores preditos (médias condicionais) para cada observação do banco de dados...

$$\hat{y} = a + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

... e os resíduos para cada observação são dados por  $y - \hat{y}$ .

No R, o jeito mais fácil de criar os valores preditos para o seu banco de dados é com a função `augment`, do pacote `broom`:

```
colera_predicted <- augment(col4)
```

`.fitted` são os valores preditos, `.resid` são os resíduos

## Função `predict()`

```
bra6 <- lm(log(renda_trabalho) ~ homem + branco + anos_estudo +  
           poly(idade, degree = 2), data = pnadc)  
yhat <- predict(bra6, data.frame(homem = 1, branco = 0,  
                                 anos_estudo = 5, idade = 30))  
  
print(yhat)  
  
##           1  
## 6.644839
```

## Função predict()

```
bra6 <- lm(log(renda_trabalho) ~ homem + branco + anos_estudo +  
           poly(idade, degree = 2), data = pnadc)  
yhat <- predict(bra6, data.frame(homem = 1, branco = 0,  
                                 anos_estudo = 5, idade = 30))  
  
print(yhat)  
  
##           1  
## 6.644839
```

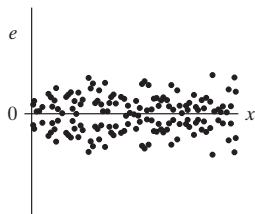
O valor predito está em log porque o  $y$  foi  $\log(\text{renda\_trabalho})$ ! Para virar nível, fazer  $e^{s_e^2/2} e^{\hat{y}}$

```
exp(yhat) * exp(summary(bra6)$sigma^2/2)  
  
##           1  
## 961.8995
```

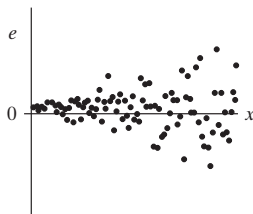
# Análise de resíduos

Muitos testes para diagnóstico de problemas podem ser feitos com os resíduos (ou resíduos padronizados) dos modelos:

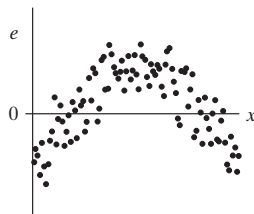
1. Histograma para verificar normalidade dos resíduos padronizados, bem como potenciais *outliers*
2. Scatterplot entre resíduo vs. valores preditos para verificar linearidade e homoscedasticidade



(a) Assumptions satisfied



(b) Nonconstant standard deviation



(c) Nonlinear term needed

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

**Cuidados no mundo real**

Encerramento do semestre

## Pequena lista...

1. Séries temporais e dados em painel normalmente possuem erros autocorrelacionados, o que fere os pressupostos clássicos e exige técnicas específicas
2. A maior parte dos bancos de dados amostrais não é coletada por amostra aleatória simples, o que exige o uso de pesos amostrais e desenhos amostrais complexos para cálculo dos erros padrão
3. Homoscedasticidade também é um pressuposto frequentemente violado, mas temos tanto testes para detectar a violação (e.g., Breusch-Pagan, White) quanto alternativas para saná-la (MQO com erros padrão robustos, GLS)
4. Como o viés de variável omitida é uma ameaça incontornável em dados observacionais, afirmações causais em geral dependem de experimentos ou quase experimentos

## Pequena lista...

5. Erros de medida são muito prejudiciais, e estão por toda parte.
  - Mesmo se o erro for aleatório em algum  $x$ , haverá viés de atenuação em todos os coeficientes para variáveis correlacionadas com o  $x$  com erro. Se o erro estiver em  $y$ , o problema é só de eficiência.
  - Se o erro não for aleatório, é mais complicado ainda
  - Se há suspeita de erros sérios, pode-se usar modelos de variáveis instrumentais
6. MQO tem limitações sérias quando o  $y$  é uma variável binária ou multinomial, entre outros casos
  - Aproximação linear só funciona bem em casos específicos
  - Para variáveis categóricas, pode ser melhor usar modelos logit ou probit

Recapitulação

Introdução

Variáveis qualitativas ou categóricas

Transformações em variáveis

Interações entre variáveis

Mais sobre valores preditos e resíduos

Cuidados no mundo real

**Encerramento do semestre**



# Datas importantes

26 de dezembro

Data limite para entrega da **Atividade #6** e da **prova final**

6 de janeiro

Entrega das notas.

Recuperação

Caso haja alunos com nota inferior ao mínimo de corte, podemos fazer atividades de recuperação em janeiro.