

# Probabilidade e Estatística

---

Eduardo Horta

# Correlação

---

# Coefficiente de correlação

- Se  $X$  e  $Y$  são duas variáveis aleatórias em uma mesma população (ou seja, definidas em um mesmo espaço amostral), então definimos o **coeficiente de correlação entre  $X$  e  $Y$**  como sendo o parâmetro

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y},$$

onde  $\sigma_X$  é o desvio padrão de  $X$ , e  $\sigma_Y$  é o desvio padrão de  $Y$ .

# Coefficiente de correlação

- Lembrando:
  - $\text{Cov}(X, Y) = \mathbb{E}\{(X - \mathbb{E}X)(Y - \mathbb{E}Y)\} = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$ ;
  - $\sigma_X^2 = \text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}^2(X)$ , etc.
- O *parâmetro*  $\rho$  é uma medida do grau de associação *linear* entre as variáveis  $X$  e  $Y$ .
- Em alguns casos é conveniente explicitar, na notação, as variáveis às quais o coeficiente se refere: escrevemos  $\rho(X, Y)$  ou  $\rho_{X,Y}$  ao invés de apenas  $\rho$ .

# Coefficiente de correlação

- O coeficiente de correlação é sempre um número entre  $-1$  e  $1$ :  
 $-1 \leq \rho \leq 1$ .
- De fato, como

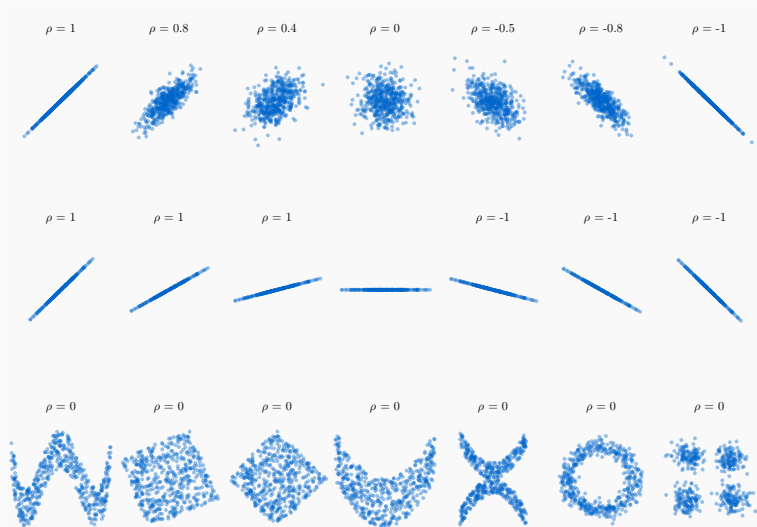
$$\left( \frac{X - \mathbb{E}X}{\sigma_X} \pm \frac{Y - \mathbb{E}Y}{\sigma_Y} \right)^2 \geq 0,$$

tomando o valor esperado em ambos os lados da desigualdade acima e expandindo o quadrado, obtemos a quota desejada  $|\rho| \leq 1$ .

# Coefficiente de correlação

- O valor de  $\rho$  nos permite interpretar o tipo de relação entre as variáveis  $X$  e  $Y$ , da seguinte forma:
  - $\rho = 1$ : associação linear *positiva* perfeita entre  $X$  e  $Y$ .
  - $0 < \rho < 1$ : associação linear *positiva* entre  $X$  e  $Y$  (quanto mais próximo de 1, mais forte é a associação linear entre essas variáveis. Quanto mais próximo de zero, mais fraca é tal associação);
  - $\rho = 0$ : ausência de associação *linear* entre  $X$  e  $Y$  (embora possa haver dependência *não linear* entre elas);
  - $-1 < \rho < 0$ : associação linear *negativa* entre  $X$  e  $Y$  (quanto mais próximo de  $-1$ , mais forte é a associação linear entre essas variáveis. Quanto mais próximo de zero, mais fraca é tal associação);
  - $\rho = -1$ : associação linear *negativa* perfeita entre  $X$  e  $Y$ .

- Na figura abaixo temos o gráfico de dispersão de amostras **simuladas**  $(X_1, Y_1), \dots, (X_n, Y_n)$ , de acordo com diferentes modelos populacionais.



## Estimação pontual de $\rho$

---



## Estimação pontual de $\rho$

- Dada uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$  de uma população descrita por uma função de distribuição  $F_{X,Y}$  com coeficiente de correlação  $\rho = \rho(X, Y)$ , um estimador natural para  $\rho$  é dado por

$$\hat{\rho} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sqrt{(\sum_{i=1}^n X_i^2 - n\bar{X}^2) \cdot (\sum_{i=1}^n Y_i^2 - n\bar{Y}^2)}}$$

- Pode-se mostrar que, nessas condições, o estimador  $\hat{\rho}$  definido acima é consistente para  $\rho$  (embora seja viesado).

## Testes de hipóteses para $\rho$

---

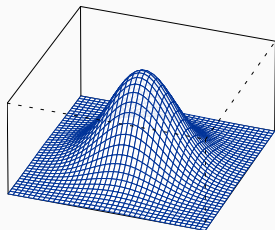
## Testes de hipóteses para $\rho$

- A estatística  $\sqrt{n}(\hat{\rho} - \rho)$  tem distribuição aproximadamente normal (se o tamanho da amostra,  $n$ , for suficientemente grande).
- Contudo, para explicitar a variância dessa estatística são necessárias suposições adicionais sobre a distribuição conjunta (modelo teórico) das variáveis  $X$  e  $Y$ .
- Por exemplo, se  $(X, Y)$  tem distribuição normal bivariada, então

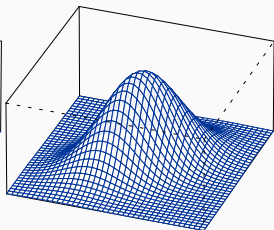
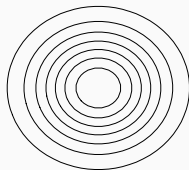
$$\sqrt{n}(\hat{\rho} - \rho) \stackrel{\text{aprox.}}{\sim} N(0, (1 - \rho^2)^2),$$

e com isso podemos testar hipóteses sobre  $\rho$  (ver mais abaixo).

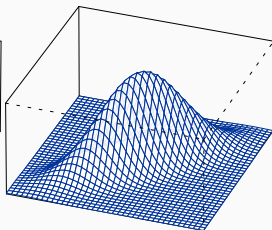
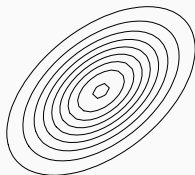
- Lembrando: distribuição *normal bivariada*



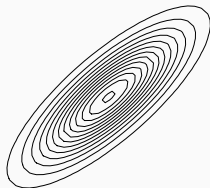
$\rho = 0$



$\rho = 0.5$



$\rho = 0.8$



## Testes de hipóteses para $\rho$

- A estatística de teste que usaremos para testar hipóteses sobre o coeficiente de correlação vai depender da hipótese nula que propusermos.
- Em qualquer caso, a validade do teste de hipóteses se baseia na suposição inicial de que o vetor aleatório  $(X, Y)$  tem distribuição normal bivariada.

- **Caso 1.** “ $H_0 : \rho = 0$ ”. Nesse caso, a estatística de teste a ser usada é

$$T = \hat{\rho} \sqrt{\frac{n-2}{1-\hat{\rho}^2}},$$

a qual tem distribuição  $t$  de Student com  $n - 2$  graus de liberdade (em particular, para valores de  $n$  razoavelmente grandes, essa distribuição é aproximadamente normal padrão).

- **Caso 2.** “ $H_0 : \rho = \rho_0$ ” (onde  $\rho_0$  é um valor numérico fixo diferente de 0). Nesse caso, a estatística de teste a ser usada é

$$\xi = \frac{\sqrt{n-3}}{2} \left( \ln\left(\frac{1+\hat{\rho}}{1-\hat{\rho}}\right) - \ln\left(\frac{1+\rho_0}{1-\rho_0}\right) \right),$$

a qual tem distribuição aproximadamente Normal padrão (isto é, com média igual a zero e variância igual a 1).

## **Coeficiente de correlação: alguns comentários**

---

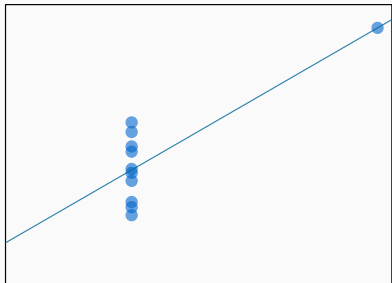
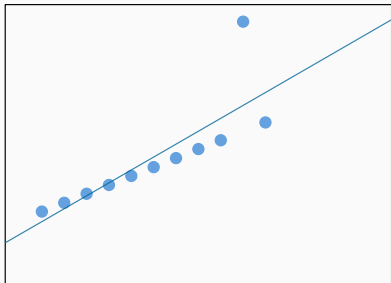
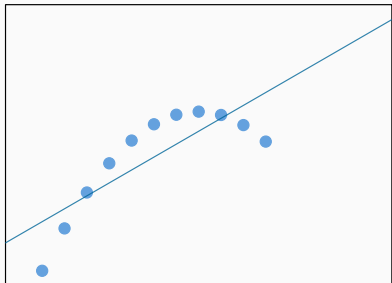
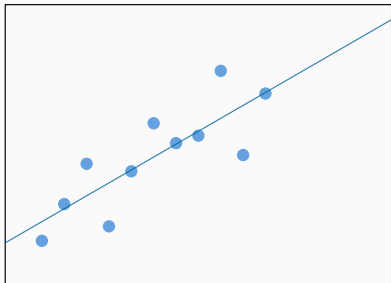


## Correlação: alguns comentários

- Conforme mencionado anteriormente, o coeficiente de correlação captura apenas o grau de associação *linear* entre duas variáveis aleatórias.
- Para dependências de tipo não-linear, o coeficiente de correlação deve ser interpretado com cautela.
- Vejamos mais alguns exemplos:

## Exemplo

- Os quatro conjuntos de dados na figura abaixo, embora bastante distintos, possuem o mesmo coeficiente de correlação (amostral):  
 $\hat{\rho} = 0.8165214$



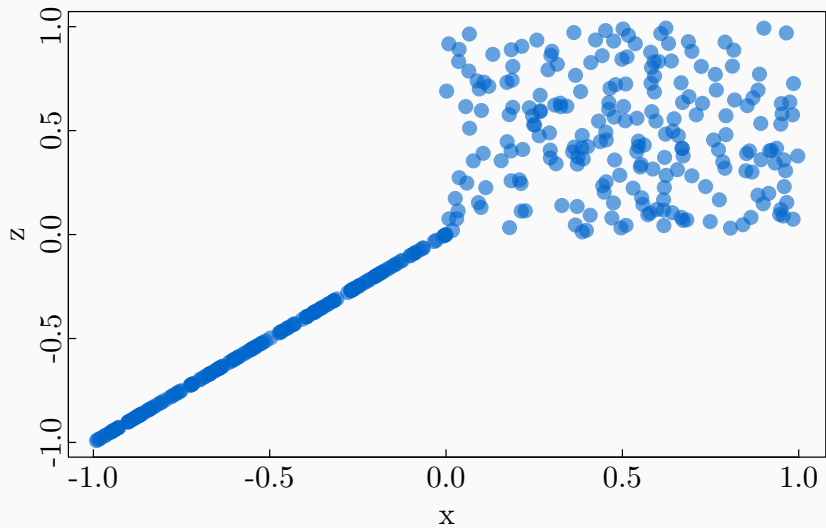
## Exemplo

- Considere variáveis aleatórias  $X \sim \text{Uniforme}[-1, 1]$  e  $Y \sim \text{Uniforme}[0, 1]$ .
- Seja  $Z$  a variável aleatória

$$Z = \begin{cases} X, & \text{se } X \leq 0; \\ Y, & \text{se } X > 0. \end{cases}$$

- Pode-se mostrar que o coeficiente de correlação entre  $X$  e  $Z$  é maior do que 0.87, algo que poderia ser interpretado como uma associação muito forte entre essas variáveis.
- Abaixo está o gráfico de dispersão de uma amostra simulada, de tamanho  $n = 400$ , dessas variáveis.

# Exemplo



## Coefficiente de correlação vs causalidade

- **Importante:** *correlação* não significa *causalidade*. Duas variáveis podem ser fortemente correlacionadas sem que uma delas seja causa (parcial) da outra.
- Por exemplo: numerosos estudos epidemiológicos mostraram que mulheres submetidas a terapia de reposição hormonal (TRH) também apresentavam uma incidência abaixo da média de doenças coronárias; isto é, esses estudos apontam para uma correlação negativa entre exposição a TRH e incidência de doenças coronárias.

## Coefficiente de correlação vs causalidade

- Levantou-se a hipótese de que a TRH seria preventiva de doenças coronárias; todavia, uma análise mais cuidadosa dos dados mostrou que mulheres submetidas a TRH provinham principalmente de classes econômicas ABC, as quais apresentam dietas e práticas de atividades físicas melhores do que a média populacional.
- Conclui-se que não havia uma relação causal entre exposição a TRH e redução do risco cardíaco, mas sim um terceiro fator causando ambas (no caso, boa alimentação e prática de exercícios, mais prevalentes nas classes ABC que também tendem a ser mais expostas a TRH).

# Coeficiente de correlação vs causalidade

- Ver também: <https://www.tylervigen.com/spurious-correlations>

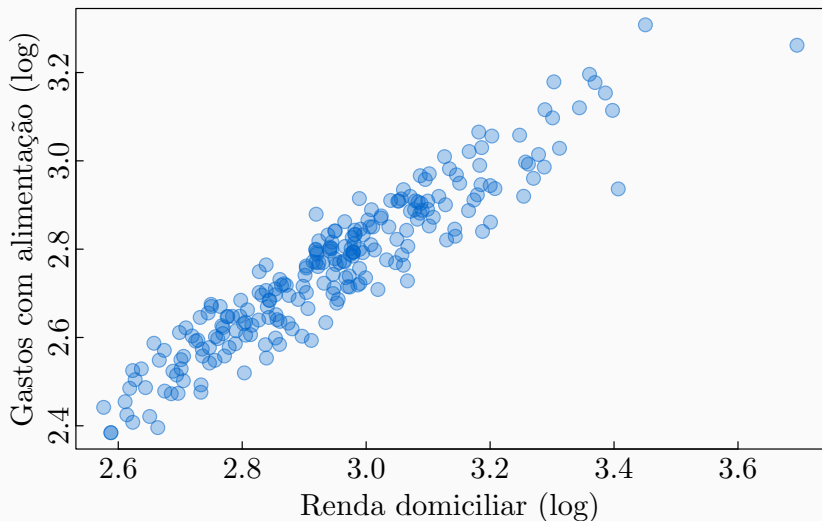


# Modelos de regressão linear

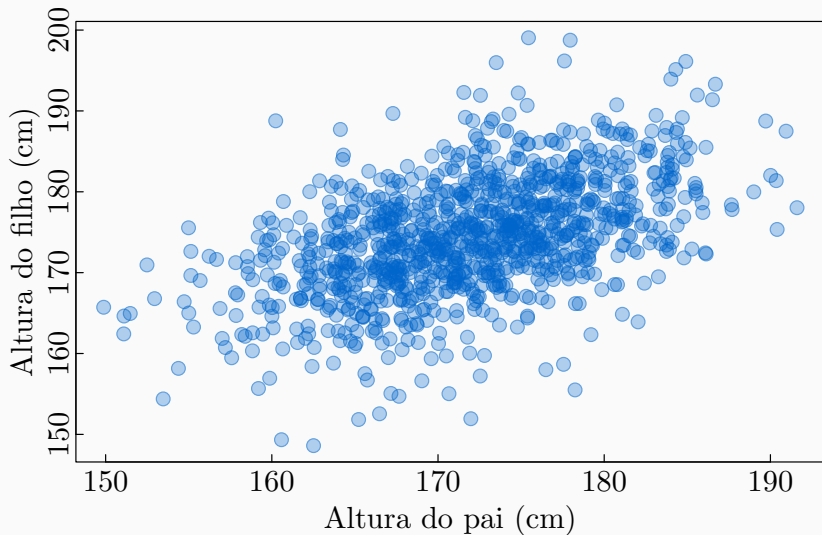
---

# Modelos de regressão linear

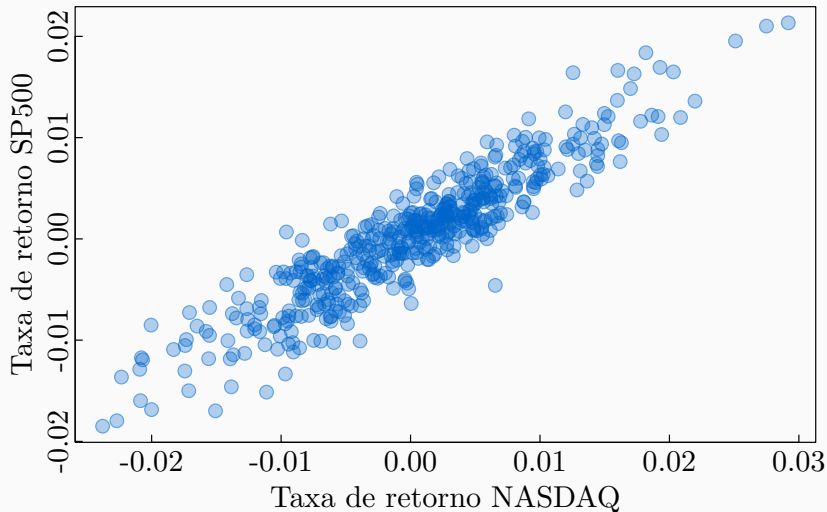
- Vamos começar com alguns exemplos (utilizando dados reais!):



## Modelos de regressão linear



## Modelos de regressão linear



# Modelos de regressão linear

- O modelo de regressão linear estipula que a relação entre duas variáveis aleatórias  $X$  e  $Y$  é descrita pela equação

$$Y = \alpha + \beta X + \varepsilon,$$

onde assume-se que  $\mathbb{E}(\varepsilon | X = x) = 0$ .

- Aqui,  $\alpha$  e  $\beta$  são números reais fixos (isto é, *parâmetros*), enquanto  $X$ ,  $Y$  e  $\varepsilon$  são variáveis aleatórias.
- Observação: nesse caso, pode-se mostrar que

$$\rho(X, Y) = \beta \frac{\sigma_X}{\sigma_Y}.$$

onde  $\sigma_X^2 = \text{Var}(X)$ , etc.

## Modelos de regressão linear: terminologia

- A variável  $Y$  é dita o **regressando**, ou a **variável explicada**, ou a **resposta**.
- A variável  $X$  é dita o **regressor**, ou a **variável explicativa**, ou a **covariável**, ou o **preditor**.
- A variável  $\varepsilon$  é dita o **termo de erro**, ou **ruído**.
- O parâmetro  $\alpha$  é dito o **intercepto** do modelo.
- O parâmetro  $\beta$  é dito o **coeficiente de inclinação** do modelo.
- A variável aleatória  $\varepsilon$  é dita o **termo de erro** do modelo.
- Observação: em alguns contextos, é usual chamar  $Y$  de variável dependente, e  $X$  de variável independente. Essa terminologia pode se sobrepôr ao conceito de variáveis aleatórias independentes e, portanto, deve ser evitada.

## Modelos de regressão linear: comentários

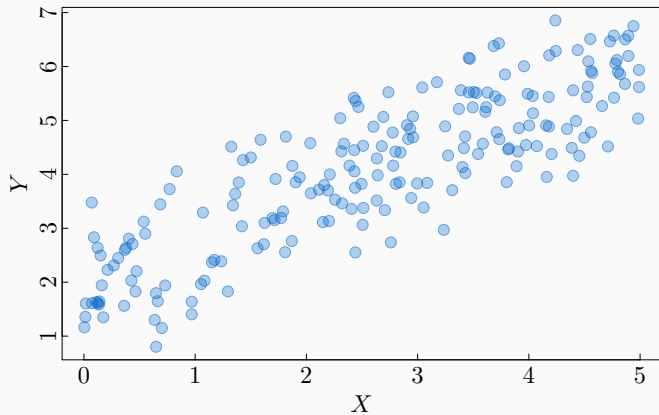
- A ideia do modelo proposto é que somente o preditor  $X$  e a resposta  $Y$  são observáveis, mas o termo de erro  $\varepsilon$  não. Além disso,  $\alpha$  e  $\beta$  são parâmetros desconhecidos;
- **Em aplicações**, coletamos uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$  de uma população a qual supostamente é descrita por um modelo linear  $Y_i = \alpha + \beta X_i + \varepsilon_i$ , com os parâmetros  $\alpha$  e  $\beta$  desconhecidos, e onde os termos de erro  $\varepsilon_i$  não são observados.

## Modelo de regressão linear: exemplo simulado

A figura a seguir mostra o gráfico de dispersão de uma amostra iid  $(X_1, Y_1), \dots, (X_n, Y_n)$ , **simulada** do modelo  $Y = 2 + 0.8 X + \varepsilon$ , onde  $X \sim \text{Uniforme}(0, 5)$  e  $\varepsilon \sim \text{Normal}(\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 0.7^2)$ , com  $X$  e  $\varepsilon$  mutuamente independentes.



# Modelo de regressão linear: exemplo simulado



## Modelo de regressão linear: exemplo simulado

- A sintaxe em R para gerar o gráfico acima é a seguinte:

```
set.seed(2)
```

```
alpha = 2
```

```
beta = .8
```

```
n = 200
```

```
epsilon = rnorm(n, sd = .7)
```

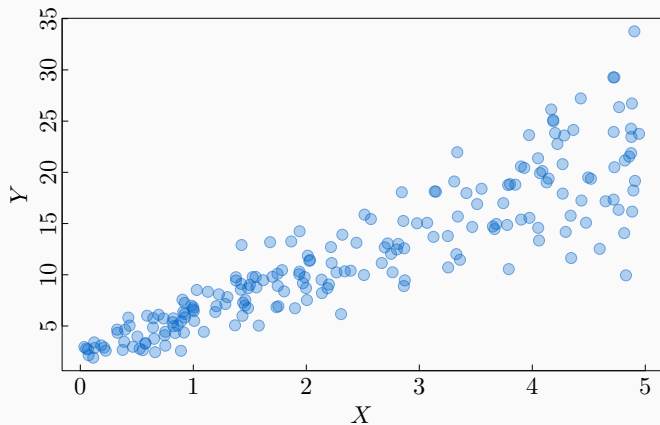
```
X = runif(n, min = 0, max = 5)
```

```
Y = alpha + beta*X + epsilon
```

## Modelo de regressão linear: outro exemplo simulado

A figura a seguir mostra o gráfico de dispersão de uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$ , **simulada** do modelo  $Y = 2 + 4X + \varepsilon$ , onde  $X \sim \text{Uniforme}(0, 5)$  e onde  $\varepsilon$ , condicional em  $X = x$ , tem distribuição Normal com  $\mathbb{E}(\varepsilon|X = x) = 0$  e  $\text{Var}(\varepsilon|X = x) = 0.8 + 0.7X$ .

## Modelo de regressão linear: outro exemplo simulado



## Modelos de regressão linear: comentários

- Algumas observações sobre o modelo linear: em muitas circunstâncias, é possível que se tenha algum tipo de controle sobre os valores da covariável  $X$  ao se executar um experimento ou coletar uma amostra.

## Modelos de regressão linear: comentários

- Isso é particularmente verdadeiro em experimentos laboratoriais, em que o experimentador de fato pode determinar os níveis de  $X$ : por exemplo, se  $Y$  é a temperatura de ebulição da água e  $X$  é o nível de concentração de alguma substância química (digamos, sal), então podemos preestabelecer alguns níveis de concentração  $X_1, X_2, \dots, X_n$  (possivelmente com repetições), e efetuar medições da temperatura de ebulição da água nesses diferentes níveis.
- Em outras situações, tipicamente não se tem esse tipo de controle (por exemplo, dados econômicos, biológicos, meteorológicos). Nos três exemplos acima, a covariável é intrinsecamente aleatória.

# **Modelos de regressão linear: Estimação**

---

## Modelos de regressão linear: Estimação

- Utilizando tão somente os valores amostrais  $(X_1, Y_1), \dots, (X_n, Y_n)$ , vamos propor estimadores  $\hat{\alpha}$  e  $\hat{\beta}$  para os parâmetros  $\alpha$  e  $\beta$ , dados respectivamente por

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2}, \quad \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$



# Modelos de regressão linear: Estimação

- O par  $(\hat{\alpha}, \hat{\beta})$  minimiza a função

$$L(a, b) = \sum_{i=1}^n (Y_i - a - bX_i)^2.$$

- Dito de outra forma:  $(\hat{\alpha}, \hat{\beta})$  é a solução do problema de minimização

$$\min_{a,b} L(a, b);$$

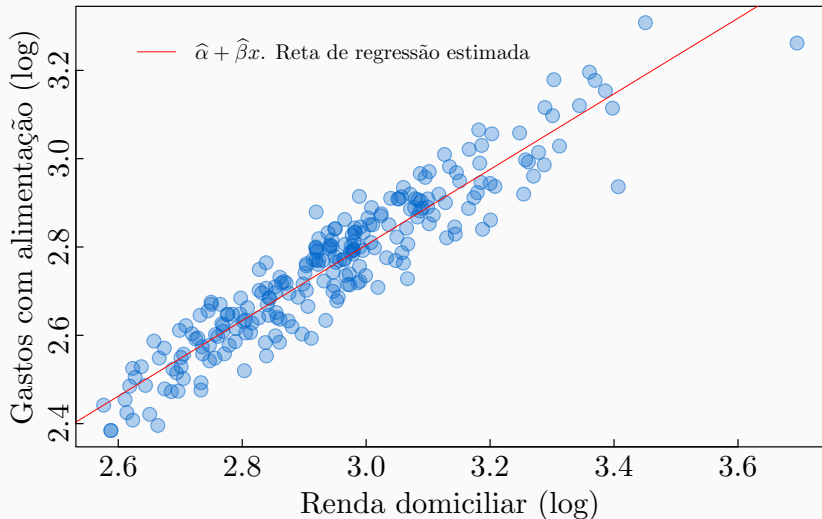
- Por esta razão,  $\hat{\alpha}$  e  $\hat{\beta}$  são chamados de estimadores de mínimos quadrados (ordinários), MQO.

## Modelos de regressão linear: exemplo aplicado

- Para o conjunto de dados sobre renda domiciliar e gastos com alimentação ilustrado anteriormente, vamos denotar por  $Y_i$  o “logaritmo dos gastos com alimentação da  $i$ -ésima unidade amostral” e por  $X_i$  o “logaritmo da renda domiciliar da  $i$ -ésima unidade amostral”.
- Nessa amostra, os valores computados dos estimadores de MQO são

$$\hat{\beta} \approx 0.856, \quad \hat{\alpha} \approx 0.237$$

## Modelos de regressão linear: exemplo aplicado



## Regressão linear: observações e propriedades

- É importante ressaltar (e o exemplo acima ilustra isso) que o termo *linear* refere-se à relação entre os parâmetros do modelo, e não entre as variáveis envolvidas.
- Por exemplo, digamos que estamos interessados em estudar a relação entre duas variáveis aleatórias  $Y$  e  $X$ , onde

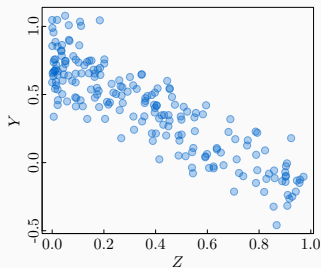
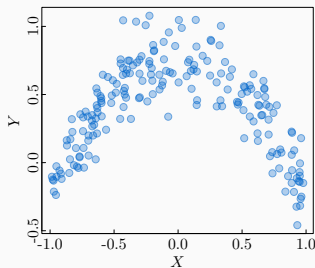
$$Y = \alpha + \beta X^2 + \varepsilon,$$

Em particular, queremos estimar os parâmetros  $\alpha$  e  $\beta$ .

- Isso pode ser feito escrevendo-se  $Z = X^2$ : a relação entre  $Y$  e  $Z$  é linear, logo podemos estimar os parâmetros de interesse através do método de MQO (aplicado aos dados  $(Z_1, Y_1) \dots, (Z_n, Y_n)$ ).

# Regressão linear: observações e propriedades

Gráficos de dispersão de uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$ , **simulada** do modelo  $Y = 3/4 - X^2 + \varepsilon$ , onde  $X \sim \text{Uniforme}(-1, 1)$  e onde  $\varepsilon$ , condicional em  $X = x$ , tem distribuição Normal com  $\mathbb{E}(\varepsilon|X = x) = 0$  e  $\text{Var}(\varepsilon|X = x) = 0.2 - 0.08|X|$ . Aqui,  $Z = X^2$ .



# Regressão linear: ajuste e previsão

---

# Regressão linear: ajuste e previsão

- No modelo de regressão linear

$$Y = \alpha + \beta X + \varepsilon$$

com  $\mathbb{E}(\varepsilon | X = x)$ , vale que a reta  $f(x) = \alpha + \beta x$  corresponde à esperança condicional de  $Y$  dado que  $X = x$ :

$$\begin{aligned}\mathbb{E}(Y | X = x) &= \mathbb{E}(\alpha + \beta X + \varepsilon | X = x) \\ &= \mathbb{E}(\alpha | X = x) + \beta \mathbb{E}(X | X = x) + \mathbb{E}(\varepsilon | X = x) \\ &= \alpha + \beta x\end{aligned}$$

## Regressão linear: ajuste e previsão

- Por essa razão, definimos o **valor ajustado (ou predito) de  $Y$  dado que  $X = x$**  como sendo a quantidade

$$\hat{Y}(x) = \hat{\alpha} + \hat{\beta}x$$

e escrevemos também

$$\hat{Y}_i = \hat{Y}(X_i)$$

quando o valor de  $x$  corresponde a um dos valores na amostra.

- Definimos ainda os **resíduos** do ajuste por

$$\hat{\varepsilon}_i = Y_i - (\hat{\alpha} + \hat{\beta}X_i), \quad i = 1, \dots, n$$



# Regressão linear: ajuste e previsão

- **Importante:** não confundir os resíduos  $\hat{\varepsilon}_i$  com os termos de erro  $\varepsilon_i$ .
- Os *resíduos* são quantidades estimadas a partir da amostra:

$$\hat{\varepsilon}_i = Y_i - (\hat{\alpha} + \hat{\beta}X_i), \quad i = 1, \dots, n$$

- Os *termos de erro* são quantidades não observáveis:

$$\varepsilon_i = Y_i - (\alpha + \beta X_i), \quad i = 1, \dots, n$$

- **Proposição:** valem as seguintes propriedades:

1. A média *amostral* dos resíduos  $\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_n$  é zero:

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i &= \frac{1}{n} \sum_{i=1}^n Y_i - (\hat{\alpha} + \hat{\beta} X_i) \\ &= \frac{1}{n} \sum_{i=1}^n Y_i - \frac{n}{n} \hat{\alpha} - \frac{\hat{\beta}}{n} \sum_{i=1}^n X_i \\ &= (\bar{Y} - \hat{\beta} \bar{X}) - \hat{\alpha} = 0\end{aligned}$$

Em particular,  $\sum_{i=1}^n \hat{\varepsilon}_i = 0$ .

2. Os termos de erro são *ortogonais* aos regressores:

$$\sum_{i=1}^n X_i \hat{\varepsilon}_i = 0.$$

## Avaliação de ajuste *in-sample*

- Uma vez ajustado o modelo de regressão (isto é, uma vez estimados os parâmetros  $\alpha$  e  $\beta$ ), podemos avaliar a qualidade desse ajuste respondendo à seguinte pergunta: *quanto da variabilidade observada na variável resposta pode ser explicada pela reta de regressão obtida?*
- Formalmente, procedemos da seguinte maneira: definimos as quantidades

$$SQ_{Tot} = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad SQ_{Res} = \sum_{i=1}^n \hat{\varepsilon}_i^2, \quad SQ_{Reg} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2,$$

denominados, respectivamente, a **soma dos quadrados totais**, a **soma dos quadrados dos resíduos** e a **soma dos quadrados da regressão**.

## Avaliação de ajuste *in-sample*

- Temos a seguinte identidade:

$$SQ_{Tot} = SQ_{Reg} + SQ_{Res}$$

- De fato, como  $Y_i - \bar{Y} = (Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y}) = \hat{\varepsilon}_i + (\hat{Y}_i - \bar{Y})$ , temos que (elevando ambos os lados ao quadrado e somando em  $i$ )

$$SQ_{Tot} = SQ_{Reg} + SQ_{Res} + 2 \sum_{i=1}^n \hat{\varepsilon}_i (\hat{Y}_i - \bar{Y}).$$

onde o terceiro termo da soma no lado direito da equação acima se anula: de fato

$$\sum_{i=1}^n \hat{\varepsilon}_i (\hat{Y}_i - \bar{Y}) = \hat{\alpha} \sum_{i=1}^n \hat{\varepsilon}_i + \hat{\beta} \sum_{i=1}^n \hat{\varepsilon}_i X_i - \bar{Y} \sum_{i=1}^n \hat{\varepsilon}_i = 0$$

## Avaliação de ajuste *in-sample*

- A expressão

$$SQ_{Tot} = SQ_{Reg} + SQ_{Res}$$

nos diz que a variabilidade dos  $Y_i$ 's em torno da média  $\bar{Y}$  se decompõe em duas componentes: a primeira correspondendo àquela parte dessa variabilidade que pode ser explicada pela reta de regressão e a segunda correspondendo à variabilidade proveniente dos resíduos.

- Sendo assim, a quantidade

$$R^2 \stackrel{\text{def}}{=} \frac{SQ_{Reg}}{SQ_{Tot}} = 1 - \frac{SQ_{Res}}{SQ_{Tot}}$$

captura o quanto da variabilidade amostral da resposta pode ser atribuído a variações no preditor.

## Regressão linear: predição e ajuste *out-of-sample*

- Na prática, nem sempre estaremos satisfeitos em obter um bom ajuste *in-sample*.
- De fato, uma das principais aplicações do modelo de regressão linear é a possibilidade de serem feitas predições *fora da amostra*
- Digamos que, a partir de uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$  retirada da população  $(X, Y)$ , obtivemos as estimativas  $\hat{\alpha}$  e  $\hat{\beta}$  para os parâmetros  $\alpha$  e  $\beta$ .
- Se obtivermos uma nova unidade para nossa amostra da qual só conhecemos o valor  $X_{n+1}$  (mas não o de  $Y_{n+1}$ ), ainda assim podemos oferecer uma *previsão* para o valor de  $Y_{n+1}$ , qual seja

$$\hat{Y}(X_{n+1}) = \hat{\alpha} + \hat{\beta}X_{n+1}.$$

- O valor predito acima é uma estimativa de  $\mathbb{E}(Y_{n+1} | X_{n+1})$ .

# **Propriedades estatísticas dos estimadores de MQO**

---

- Consideremos uma amostra aleatória  $(X_1, Y_1), \dots, (X_n, Y_n)$  da população descrita por uma função de distribuição  $F_{X,Y}$ , onde  $X$  e  $Y$  são relacionados pelo modelo

$$Y = \alpha + \beta X + \varepsilon,$$

com  $\mathbb{E}(\varepsilon | X = x) = 0$  e  $\text{Var}(\varepsilon | X = x) = \sigma_\varepsilon^2(x)$ .



# Propriedades estatísticas dos estimadores de MQO

- **Proposição:** Nas condições acima, valem as seguintes propriedades para o estimador  $\hat{\beta}$ :

- (1)  $\hat{\beta}$  é não-viesado para  $\beta$ :

$$\mathbb{E}(\hat{\beta}) = \beta$$

- (2) Escrevendo  $\tilde{X} = (X_1, \dots, X_n)$ , a variância de  $\hat{\beta}$  condicional em  $\tilde{X}$  é dada por:

$$\text{Var}(\hat{\beta} | \tilde{X}) = \frac{\sum_{i=1}^n \sigma_{\varepsilon}^2(X_i)(X_i - \bar{X})^2}{\left(\sum_{i=1}^n (X_i - \bar{X})^2\right)^2},$$

onde a variância incondicional de  $\hat{\beta}$  satisfaz

$$\text{Var}(\hat{\beta}) \approx \frac{\mathbb{E}\{\sigma_{\varepsilon}^2(X)(X - \mathbb{E}X)\}}{n\sigma_X^4}.$$

Em particular,  $\hat{\beta}$  é consistente para  $\beta$ .

# Propriedades dos estimadores de MQO

- Segue que:

(3) Condicional em  $X_1, \dots, X_n$ , o estimador  $\hat{\beta}$  tem distribuição aproximadamente normal:

$$\hat{\beta}|\tilde{X} \sim N(\beta, \text{Var}(\hat{\beta}|\tilde{X}))$$

(4) Sob a hipótese nula  $H_0 : \beta = \beta_0$ , a estatística de teste

$$T = \frac{\hat{\beta} - \beta_0}{\hat{V}}$$

tem distribuição aproximadamente Normal padrão. Aqui,  $\hat{V}^2$  é um estimador de  $\text{Var}(\hat{\beta}|\tilde{X})$  dado por

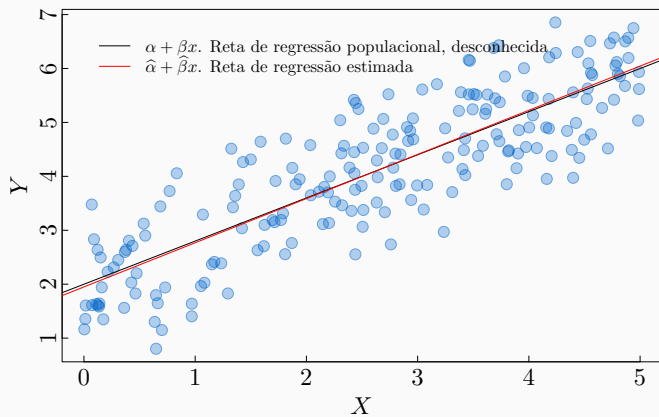
$$\hat{V}^2 = \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2 (X_i - \bar{X})^2}{\left(\sum_{i=1}^n (X_i - \bar{X})^2\right)^2}$$

## Exemplo

- Voltemos ao nosso **exemplo** simulado: temos uma amostra iid  $(X_1, Y_1), \dots, (X_n, Y_n)$  do modelo  $Y = 2 + 0.8X + \varepsilon$ , onde  $X \sim \text{Uniforme}(0, 5)$  e  $\varepsilon \sim \text{Normal}(\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 0.7^2)$ , com  $X$  e  $\varepsilon$  mutuamente independentes. Aqui temos  $n = 200$ .
- Os valores de  $\hat{\alpha}$  e  $\hat{\beta}$  correspondentes a essa amostra são

$$\hat{\alpha} = 1.94934 \quad \text{e} \quad \hat{\beta} = 0.81907$$

# Exemplo



## Exemplo: Testes de hipóteses

- Se quisermos testar a hipótese nula “ $H_0 : \beta = 0$ ” contra a alternativa “ $H_1 : \beta > 0$ ”, devemos obter o valor da estatística de teste

$$T = \frac{\hat{\beta}}{\widehat{V}},$$

a qual segue (sob  $H_0$ ) uma distribuição (aproximadamente) Normal padrão (conforme vimos acima).

- Nessa particular amostra o valor da estatística de teste é 23.767. Ao nível de significância de 1%, portanto, rejeitamos  $H_0$ . (Por quê?)

**Observação:** se quiséssemos ser inteiramente rigorosos na notação, deveríamos escrever  $X_i(\omega)$ ,  $\bar{X}_n(\omega)$ ,  $\hat{\alpha}_n(\omega)$ , etc, que são os valores numéricos das variáveis aleatórias em questão *na particular amostra*  $\omega$  *de tamanho*  $n$ . Essa notação fica um tanto carregada e, por isso, é costumeiro omitir o “ $\omega$ ” e o  $n$ .

# **Uma aplicação: classificação binária**

---

## Uma aplicação: classificação binária

- Pela definição do modelo de regressão linear, temos

$$\mathbb{E}(Y | X = x) = \alpha + \beta x.$$

- Quer dizer, a reta dada pela equação  $\alpha + \beta x$  descreve o comportamento médio da resposta nos diferentes *níveis* do preditor  $X$ .
- Esse modelo tem uma interpretação interessante no caso em que a resposta é uma variável binária (isto é, quando  $Y$  tem distribuição Bernoulli).



## Uma aplicação: classificação binária

- **Lembrando:**  $Y$  tem distribuição Bernoulli com parâmetro  $u$  se

$$\mathbb{P}(Y = 1) = u, \quad \mathbb{P}(Y = 0) = 1 - u.$$

- Nesse caso, temos a identidade

$$\mathbb{E}(Y) = (0 \times \mathbb{P}(Y = 0)) + (1 \times \mathbb{P}(Y = 1)) = \mathbb{P}(Y = 1) = u.$$

- Semelhantemente, condicionando em uma covariável  $X$ ,

$$\mathbb{E}(Y | X = x) = \mathbb{P}(Y = 1 | X = x) =: u(x)$$

## O modelo de probabilidade linear

- Se  $Y$  é uma variável com distribuição Bernoulli e  $X$  é uma covariável, o *modelo de probabilidade linear* estipula que

$$\mathbb{P}(Y = 1 \mid X = x) = \alpha + \beta x$$

- Evidentemente, o lado esquerdo da equação acima deve ser um número entre 0 e 1, para qualquer nível  $x$  da covariável  $X$ :

$$0 \leq \alpha + \beta x \leq 1.$$

Em particular, tal modelo só faz sentido se o preditor for uma variável aleatória limitada.

- Por exemplo, se  $\beta > 0$ , então deve valer

$$\mathbb{P}\left(-\frac{\alpha}{\beta} \leq X \leq \frac{1 - \alpha}{\beta}\right) = 1$$

# O modelo de probabilidade linear

- Suponhamos, então, que  $\mathbb{P}(Y = 1 | X = x) = \alpha + \beta x$ .
- Definindo  $\varepsilon = Y - (\alpha + \beta X)$ , podemos escrever

$$Y = \alpha + \beta X + \varepsilon$$

e é de fácil verificação que  $\mathbb{E}(\varepsilon | X = x) = 0$ .

- Quer dizer, de fato o modelo de probabilidade linear nada mais é do que um modelo de regressão linear no caso em que a resposta tem distribuição Bernoulli.

## O modelo de probabilidade linear

- Esse modelo ilustra como a suposição (usual em textos introdutórios) de que termo de erro e covariável são independentes é facilmente violada.
- De fato, temos  $\text{Var}(\varepsilon|X = x) = \mathbb{E}(\varepsilon^2|X = x)$  já que  $\mathbb{E}(\varepsilon|X = x) = 0$ . Logo, escrevendo  $u(x) = \mathbb{E}(Y|X = x) = \alpha + \beta x$ , vemos que

$$\text{Var}(\varepsilon|X = x) = \mathbb{E}(Y^2 - 2Yu(x) + u(x)^2 | X = x).$$

- Como  $Y$  só pode assumir os valores 0 ou 1, vale que  $Y^2 = Y$ . Assim,

$$\text{Var}(\varepsilon|X = x) = \mathbb{E}(Y|X = x) - 2u(x)\mathbb{E}(Y|X = x) + u(x)^2.$$

- Segue, da definição de  $u(x)$ , que

$$\text{Var}(\varepsilon|X = x) = u(x) \times (1 - u(x)).$$

## O modelo de probabilidade linear: simulação

- A figura a seguir mostra o gráfico de dispersão de uma amostra aleatória de tamanho  $n = 200$  do modelo de probabilidade linear descrito acima, com  $\alpha = 1/8$  e  $\beta = 6/8$ .
- Nessa amostra, o valor das estimativas dos parâmetros de intercepto e inclinação foram, respectivamente  $\hat{\alpha} = 0.1058$  e  $\hat{\beta} = 0.7264$ .

# O modelo de probabilidade linear: simulação

