

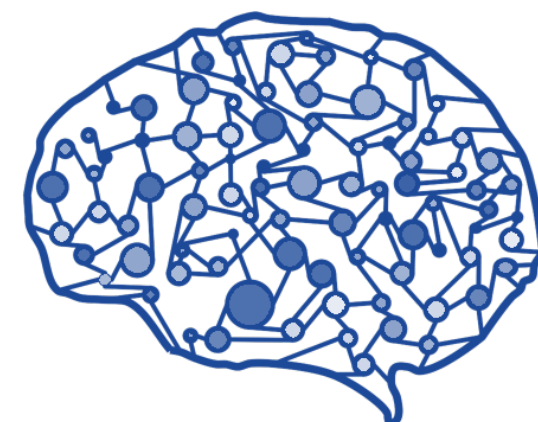


# COC800

# Introdução à Ciência de Dados

# Aprendizado Supervisionado

- Vizinhos mais próximos
- Modelagem bayesiana
- Modelos lineares
- Árvores e regras de decisão
- Redes neurais artificiais

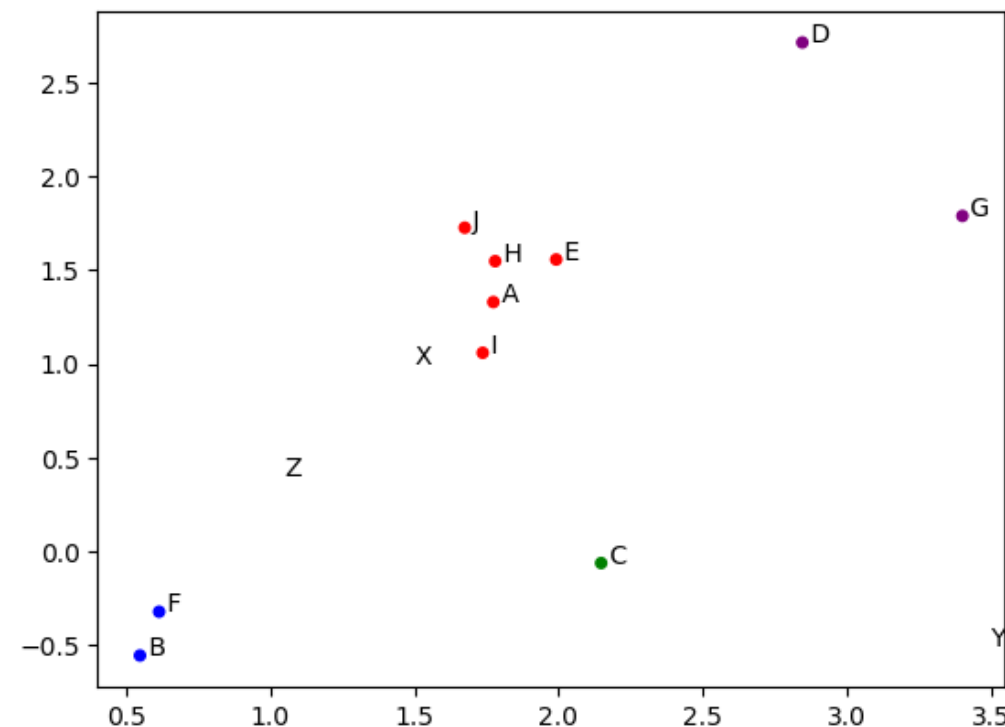


## ***K-nearest neighbors (knn)***

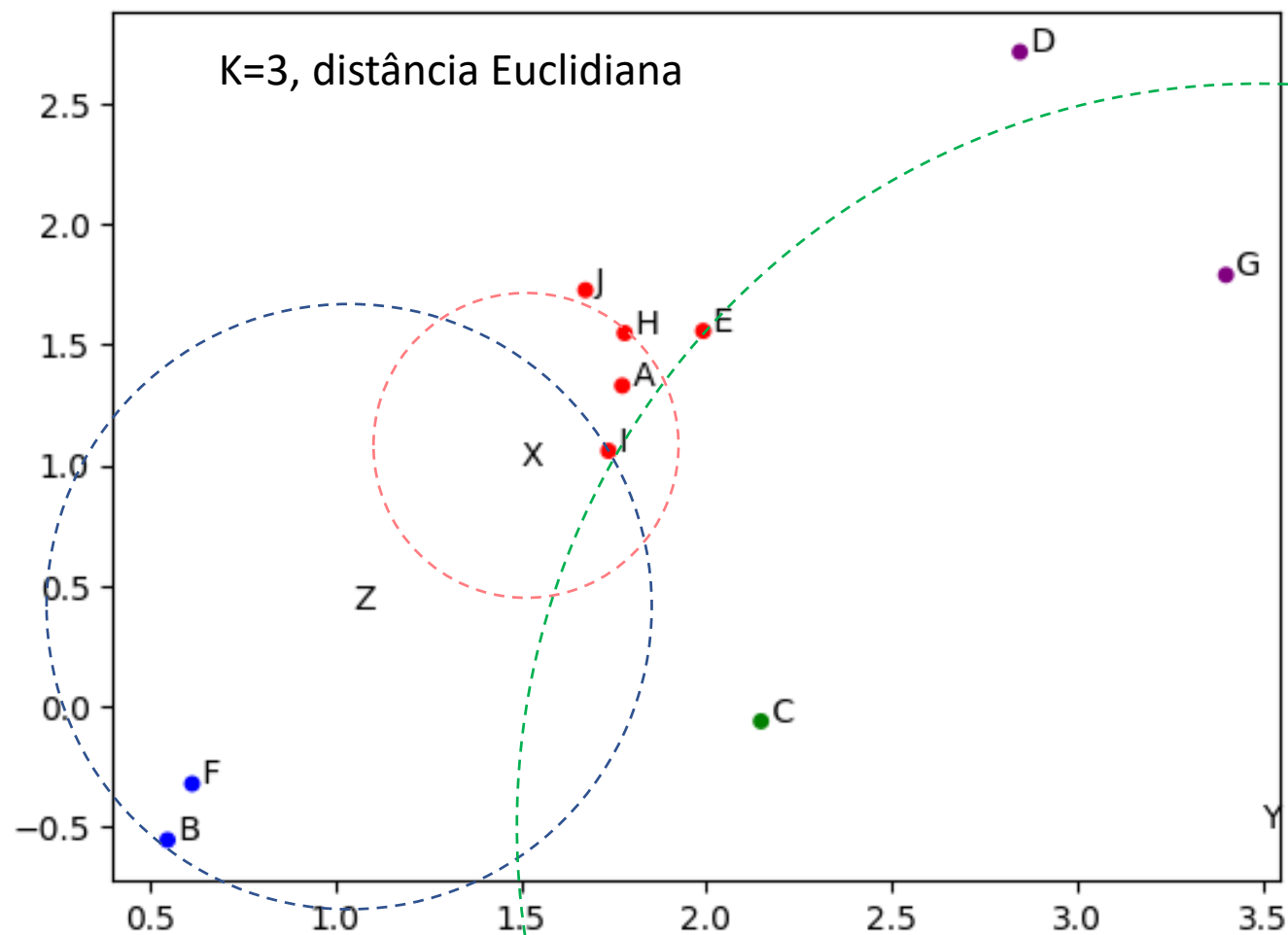
Proposto em 1951, é um dos mais utilizados algoritmos preditivos sobretudo com o surgimento dos sistemas de recomendação. Ele não produz um modelo de representação dos dados e a generalização a partir dos dados de “treinamento” é realizada no momento da tomada de decisão. Algoritmos com esse comportamento são chamados de “preguiçosos” (*lazy*).

Algoritmos “preguiçosos” são interessantes em problemas cujas bases de dados mudam intensamente e são consultadas frequentemente.

O knn é um algoritmo que produz a resposta para um determinado registro a partir dos seus  $k$  registros mais próximos, realizando uma aproximação local da variável resposta. Em problemas de classificação, a saída é a classe mais frequente dentre os vizinhos. Na regressão, a saída é uma combinação dos valores da variável resposta, em geral a média aritmética deles.



## K-nearest neighbors (knn)



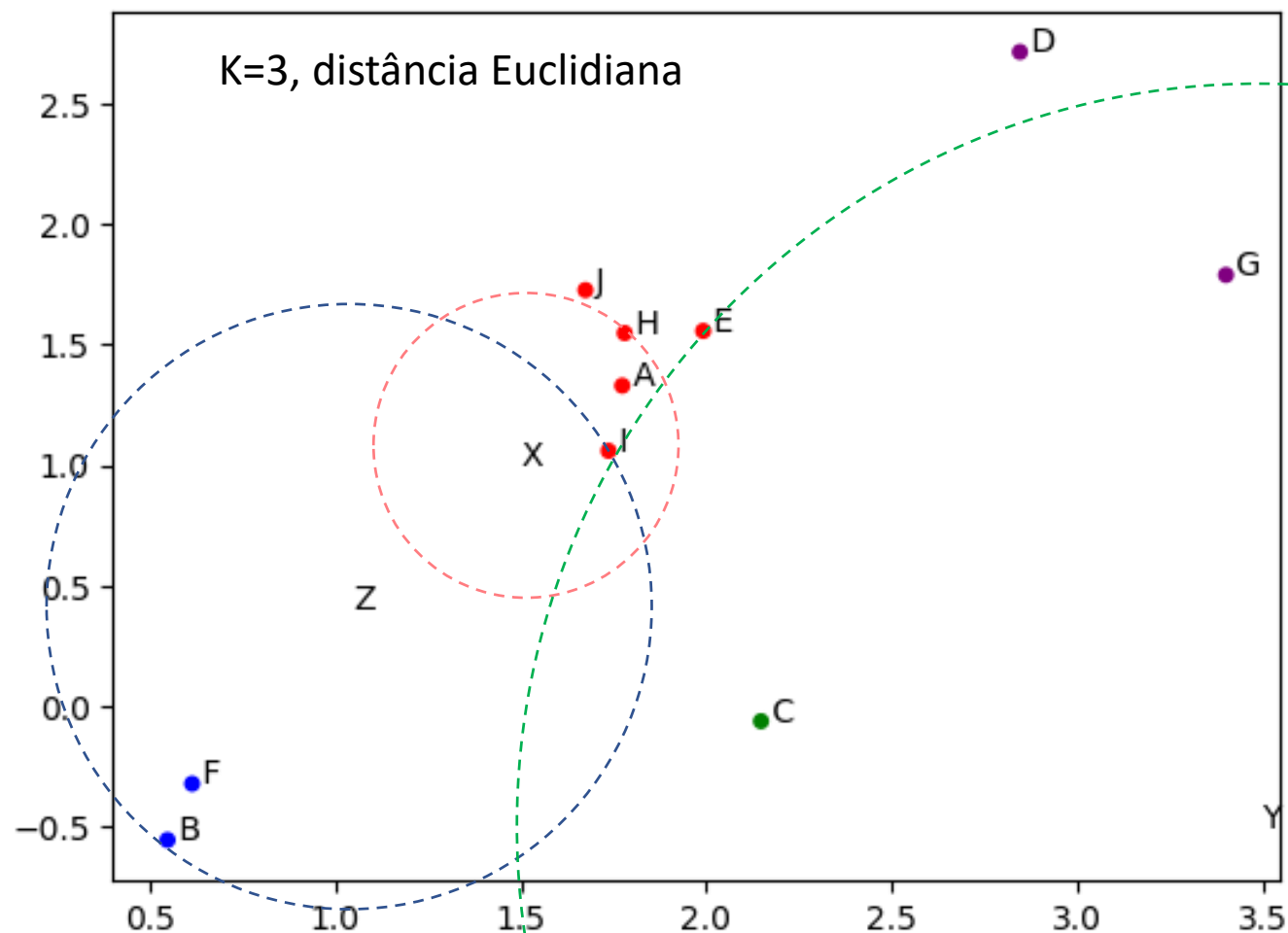
Em geral, ponderar as contribuições dos vizinhos, segundo suas proximidades ao que se quer a resposta, costuma produzir melhores resultados. Assim, registros mais próximos contribuem mais com a resposta.

O “treinamento” apenas armazena os dados em memória.

A escolha do  $k$  depende dos dados. Valores maiores implicam em robustez a ruídos, porém piora a identificação das fronteiras de decisão.

Atributos irrelevantes devem ser excluídos para não interferirem nos cálculos das distâncias.

## K-nearest neighbors (knn)



A escolha da função distância a ser adotada também é importante e pode levar a uma melhor separação entre as classes.

Pode-se aplicar a análise de componentes principais para reduzir o número de atributos, tornando o cálculo de distância mais eficiente.

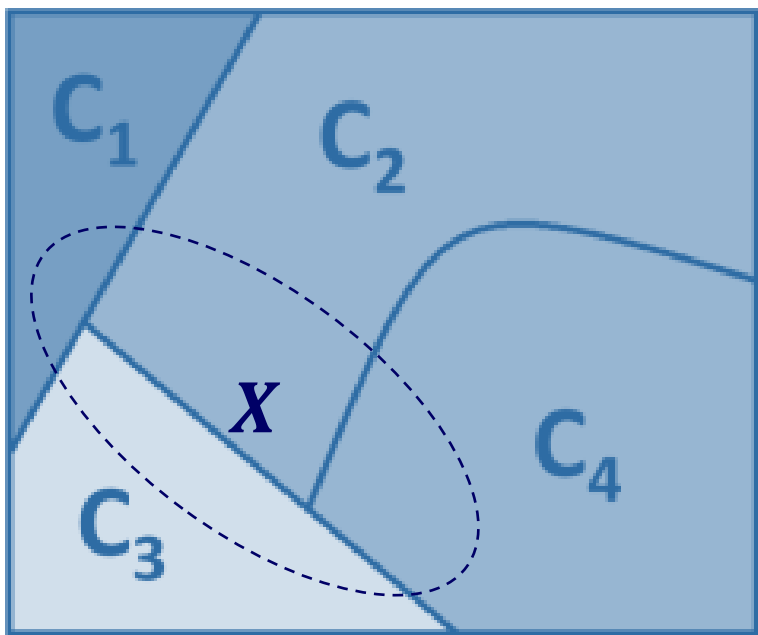
Podem ser arbitrados limites para os raios de observação de vizinhos evitando-se que a computação exaustiva decorrente do cálculo de todas as distâncias seja realizada.

Ele é facilmente paralelizado o que promove uma grande melhoria de desempenho.

## Classificador bayesiano simples (*naïve Bayes classifier*)

São algoritmos de classificação simples, rápidos e de bom desempenho, que usam o Teorema de Bayes para a tomada de decisão supondo a independência entre os atributos.

Eles consideram que cada atributo contribui de modo independente para a probabilidade de um objeto pertencer a uma classe, o que motiva a sua qualificação como “ingênuo”.



Teorema de Bayes

$$P(C_k | \mathbf{x}) = \frac{P(C_k \cap \mathbf{x})}{P(\mathbf{x})} = \frac{P(\mathbf{x} | C_k) \cdot P(C_k)}{P(\mathbf{x})}$$

em que

$$P(\mathbf{x}) = \sum_{k=1}^K P(\mathbf{x} | C_k) \cdot P(C_k)$$

Informações conhecidas

Informações supostamente conhecidas

A classe com maior probabilidade  $P(C_k | \mathbf{x}_0)$  será atribuída a um novo registro  $\mathbf{x}_0$ .

## Classificador bayesiano simples (*naïve Bayes classifier*)

Como  $P(\mathbf{x})$  é igual para toda classe, basta calcular  $P(C_k|\mathbf{x}) = P(\mathbf{x}|C_k) \cdot P(C_k)$  durante a classificação.

Pela suposição simplista de independência de atributos, esses são representados por eventos independentes e, portanto:

$$P(\mathbf{x}|C_k) \approx \prod_{i=1}^n P(x_i|C_k) = P(x_1|C_k) \cdot \dots \cdot P(x_n|C_k), \text{ } i \text{ atributo de } \mathbf{x}$$

Se a variável  $x_i$  for discreta,  $P(x_i|C_k)$  é a frequência relativa de  $x_i$  em  $C_k$ .

Se a variável  $x_i$  for contínua,  $P(x_i|C_k)$  é normalmente estimada pela função densidade de probabilidade de uma distribuição gaussiana:

$$P(x_i|C_k) = \frac{1}{\sqrt{2\pi}\sigma_{C_k}} e^{-\frac{(x_i - \mu_{C_k})^2}{2\sigma_{C_k}^2}}$$

As diferentes abordagens ao classificador bayesiano simples referem-se, principalmente, a forma de estimar a distribuição de  $P(x_i|C_k)$ .

## Classificador bayesiano simples (*naïve Bayes classifier*)

Exemplo:

Registro	Sexo	Curso	Idiomas	Trabalha	Classe
1	M	Eng.Civil	1	N	A
2	M	Eng.Elétrica	1	N	A
3	M	Eng.Elétrica	2	N	B
4	F	Eng.Elétrica	2	N	A
5	F	Eng.Civil	2	N	A
6	M	Eng.Elétrica	2	S	B
7	M	Eng.Elétrica	1	S	B
8	M	Eng.Elétrica	1	S	C
9	M	Eng.Civil	1	S	C
10	M	Eng.Civil	2	S	C
11	F	Eng.Civil	1	S	C
12	F	Eng.Civil	1	N	C
13	F	Eng.Civil	1	S	B
14	F	Eng.Civil	3	S	A
15	M	Eng.Civil	1	S	A

$$P(C_A)=0,40$$

$$P(C_B)=0,27$$

$$P(C_C)=0,33$$

$$P(M/C_A)=0,50$$

$$P(F/C_A)=0,50$$

$$P(M/C_B)=0,75$$

$$P(F/C_B)=0,25$$

$$P(M/C_C)=0,60$$

$$P(F/C_C)=0,40$$

$$P(\text{Eng.Elét}/C_A)=0,33$$

$$P(\text{Eng.Civil}/C_A)=0,67$$

$$P(\text{Eng.Elét}/C_B)=0,75$$

$$P(\text{Eng.Civil}/C_B)=0,25$$

$$P(\text{Eng.Elét}/C_C)=0,20$$

$$P(\text{Eng.Civil}/C_C)=0,80$$

$$P(1/C_A)=0,50$$

$$P(2/C_A)=0,33$$

$$P(3/C_A)=0,17$$

$$P(1/C_B)=0,50$$

$$P(2/C_B)=0,50$$

$$P(3/C_B)=0,00$$

$$P(1/C_C)=0,80$$

$$P(2/C_C)=0,20$$

$$P(3/C_C)=0,00$$

$$P(N/C_A)=0,67$$

$$P(S/C_A)=0,33$$

$$P(N/C_B)=0,25$$

$$P(S/C_B)=0,75$$

$$P(N/C_C)=0,20$$

$$P(S/C_C)=0,80$$

modelo produzido



## Classificador bayesiano simples (*naïve Bayes classifier*)

Exemplo:

Registro	Sexo	Curso	Idiomas	Trabalha	Classe
X	F	Eng.Civil	2	S	?

$$P(C_k|\mathbf{x}) = \frac{P(C_k \cap \mathbf{x})}{P(\mathbf{x})} = \frac{P(\mathbf{x}|C_k) \cdot P(C_k)}{P(\mathbf{x})}$$

$$P(\mathbf{x}|C_A) = P(F|C_A) \cdot P(E. Civil|C_A) \cdot P(2|C_A) \cdot P(S|C_A) = 0,50 \cdot 0,67 \cdot 0,33 \cdot 0,33 = 0,036$$

$$P(\mathbf{x}|C_B) = P(F|C_B) \cdot P(E. Civil|C_B) \cdot P(2|C_B) \cdot P(S|C_B) = 0,25 \cdot 0,25 \cdot 0,50 \cdot 0,75 = 0,023$$

$$P(\mathbf{x}|C_C) = P(F|C_C) \cdot P(E. Civil|C_C) \cdot P(2|C_C) \cdot P(S|C_C) = 0,40 \cdot 0,80 \cdot 0,20 \cdot 0,80 = 0,051$$

$$P(C_A|\mathbf{x}) = P(\mathbf{x}|C_A) \cdot P(C_A) = 0,036 \cdot 0,40 = 0,014$$

$$P(C_B|\mathbf{x}) = P(\mathbf{x}|C_B) \cdot P(C_B) = 0,023 \cdot 0,27 = 0,006$$

$$P(C_C|\mathbf{x}) = P(\mathbf{x}|C_C) \cdot P(C_C) = 0,051 \cdot 0,33 = 0,017$$

**classe C**

$$P(C_A)=0,40$$

$$P(C_B)=0,27$$

$$P(C_C)=0,33$$

$$P(M/C_A)=0,50$$

$$P(F/C_A)=0,50$$

$$P(M/C_B)=0,75$$

$$P(F/C_B)=0,25$$

$$P(M/C_C)=0,60$$

$$P(F/C_C)=0,40$$

$$P(\text{Eng.Elét}/C_A)=0,33$$

$$P(\text{Eng.Civil}/C_A)=0,67$$

$$P(\text{Eng.Elét}/C_B)=0,75$$

$$P(\text{Eng.Civil}/C_B)=0,25$$

$$P(\text{Eng.Elét}/C_C)=0,20$$

$$P(\text{Eng.Civil}/C_C)=0,80$$

$$P(N/C_A)=0,67$$

$$P(S/C_A)=0,33$$

$$P(N/C_B)=0,25$$

$$P(S/C_B)=0,75$$

$$P(N/C_C)=0,20$$

$$P(S/C_C)=0,80$$

$$P(1/C_A)=0,50$$

$$P(2/C_A)=0,33$$

$$P(3/C_A)=0,17$$

$$P(1/C_B)=0,50$$

$$P(2/C_B)=0,50$$

$$P(3/C_B)=0,00$$

$$P(1/C_C)=0,80$$

$$P(2/C_C)=0,20$$

$$P(3/C_C)=0,00$$

## Classificador bayesiano simples (*naïve Bayes classifier*)

Em problemas com dados reais, pode-se evitar as probabilidades nulas aplicando-se a correção de Laplace no atributo: a cada possível valor do atributo, uma unidade é adicionada a sua contagem e ao total de registros.

No exemplo anterior, se um novo registro a ser classificado tivesse o valor de um outro curso de graduação não presente nos dados originais, como Matemática, por exemplo, as contagens de Engenharia Civil e Engenharia Elétrica subiriam para 10 e 7, respectivamente, e o total seria 18.

Assim, as contagens e probabilidades afetadas, anteriores e novas, seriam as seguintes, respectivamente:

Classe	Eng. Civil	Eng. Elétrica	Matemática	Total	Classe	Eng. Civil	Eng. Elétrica	Matemática	Total
A	4	2	0	6	A	4+1	2+1	0+1	6+1+1+1
B	1	3	0	4	B	1+1	3+1	0+1	4+1+1+1
C	4	1	0	5	C	4+1	1+1	0+1	5+1+1+1
Classe	Eng. Civil	Eng. Elétrica	Matemática	Total	Classe	Eng. Civil	Eng. Elétrica	Matemática	Total
A	0,67	0,33	0	1	A	0,56	0,33	0,11	1
B	0,25	0,75	0	1	B	0,29	0,57	0,14	1
C	0,80	0,20	0	1	C	0,625	0,25	0,125	1

## Classificador bayesiano simples (*naïve Bayes classifier*)

A adição de uma unidade apresentada no exemplo anterior não é obrigatória, mas é a mais comum. Pode-se considerar essa quantidade como um hiperparâmetro a ser determinado por validação cruzada. Outra possibilidade é realizar a adição de modo diferenciado por variável, conforme as características de cada uma.

Outro problema surge quando os produtórios se degeneram e, portanto, aplica-se a função logaritmo às probabilidades:

$$P(C_k|\mathbf{x}) = P(\mathbf{x}|C_k) \cdot P(C_k)$$

$$\ln(P(C_k|\mathbf{x})) = \ln(P(\mathbf{x}|C_k) \cdot P(C_k)) = \ln P(\mathbf{x}|C_k) + \ln P(C_k) =$$

linearidade

$$\ln \prod_{i=1}^n P(x_i|C_k) + \ln P(C_k) = \ln P(x_1|C_k) + \cdots + \ln P(x_n|C_k) + \ln P(C_k)$$

## Classificador bayesiano simples (*naïve Bayes classifier*)

Uma outra maneira de categorizar as modelagens preditivas é qualificá-las como discriminativas e generativas.

Considerando a tarefa de classificação, modelos discriminativos aprendem a determinar as fronteiras de decisão que separam as diferentes classes nos dados de entrada. Na tarefa de regressão, eles estimam diretamente a relação entre as variáveis de entrada e a variável de saída.

Já os generativos modelam a distribuição de probabilidade conjunta dos dados de entrada e de saída, estimando os processos geradores desses dados, o que permite produzir novos dados semelhantes aos que foram utilizado durante o aprendizado.

Assim, classificadores bayesianos são da categoria generativa e podem ser utilizados para gerar novas amostras de dados a calculados a partir de  $P(C_k)$  e  $P(\mathbf{x}_i|C_k)$  originais dos dados e de  $P(\mathbf{x}|C_k)$  e  $P(C_k|\mathbf{x})$  calculados. Considerando o exemplo apresentado, podemos usar o modelo para gerar a classe do registro abaixo fazendo um sorteio da classe a ser atribuída baseando-se nas  $P(C_k|\mathbf{x})$  calculadas e normalizadas:

Registro	Sexo	Curso	Idiomas	Trabalha	Classe
X	F	Eng.Civil	2	S	?

$$\begin{array}{l}
 P(C_A|\mathbf{x}) = 0,014 \\
 P(C_B|\mathbf{x}) = 0,006 \\
 P(C_C|\mathbf{x}) = 0,017
 \end{array}
 \xrightarrow[\text{(\div } P(\mathbf{x}))]{\text{normalização}}
 \begin{array}{l}
 P(C_A|\mathbf{x}) = 38\% \\
 P(C_B|\mathbf{x}) = 16\% \\
 P(C_C|\mathbf{x}) = 46\%
 \end{array}$$

## Regressão linear

Modelagem que procura estabelecer uma relação linear entre as variáveis de entrada, independentes, e a variável de saída, dependente. Ela é **simples** quando há uma única independente, ou **múltipla**, quando há várias.

Problema real:

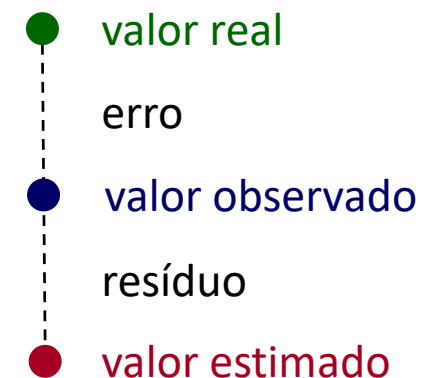
$$y = f(x_1, \dots, x_n) + \varepsilon$$

- $y$  é a variável dependente (de decisão, de resposta, explicada, endógena);
- $x_1, \dots, x_n$  são as variáveis independentes (preditoras, regressoras, explicativas);
- $f$  é desconhecida;
- $\varepsilon$  é um termo de erro, assumido como sendo independente de  $\mathbf{x} = [x_1, \dots, x_n]$ .

Modelo linear aproximado:

$$y \approx \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$

- modelo linear populacional: a melhor aproximação linear desconhecida do modelo real original;
- $\beta_0, \dots, \beta_n$  são os parâmetros desconhecidos do modelo;
- $\varepsilon$  é um termo de erro que abrange tudo o que  $\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$  não conseguiu capturar.

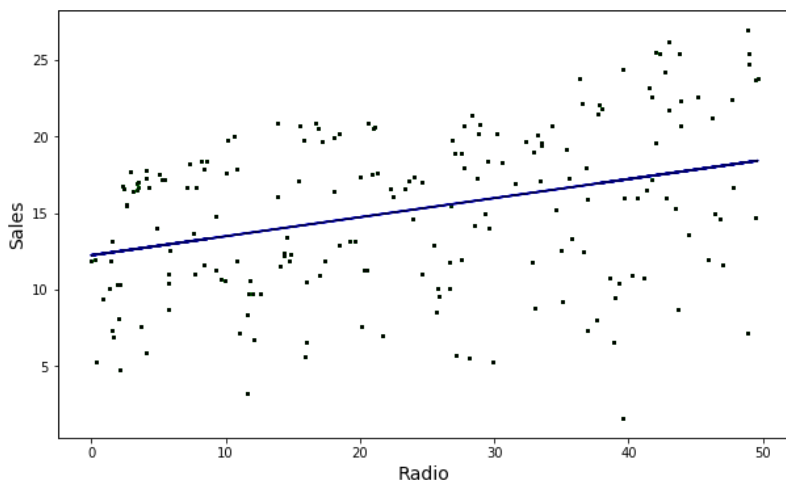


## Regressão linear

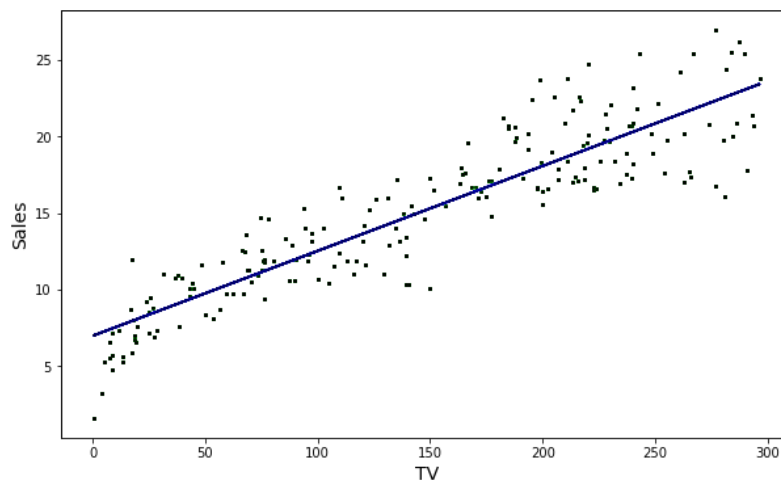
Modelo linear estimado:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_n x_n$

- modelo estimado da aproximação  $y \approx \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$  a partir dos dados disponíveis para gerá-lo;
- $\hat{\beta}_0, \dots, \hat{\beta}_n$  são os **parâmetros** ou **coeficientes** que estimam  $\beta_0, \dots, \beta_n$  do modelo aproximado;
- $\hat{\beta}_0$  é chamado de **intercepto**, ou seja, é o valor da estimativa da função quando os efeitos de todas as variáveis são nulos;
- Com  $\hat{\beta}_0$  nulo, assume-se que a saída é proporcional às entradas.

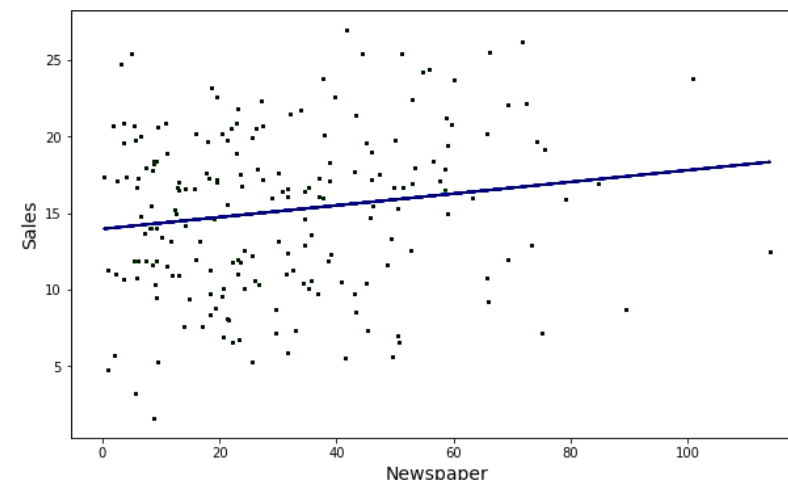
Exemplo: *Advertising dataset*, 200 registros, 4 atributos



$$Sales = 12,236 + 0,124 \text{ Radio}$$



$$Sales = 6,975 + 0,056 \text{ TV}$$



$$Sales = 13,960 + 0,038 \text{ Newspaper}$$

## Regressão linear

Existe relação entre as variáveis dependentes e independentes?

Se os dados apontam que não há relação, descarta-se a variável.

Se houver relação, ela é linear?

Se a relação tem comportamento aproximadamente “linear”, então a regressão linear é interessante. Se não, talvez seja melhor transformar algum preditor ou mesmo a variável de resposta para que a regressão linear seja melhor aplicada.

Se houver relação, qual a intensidade?

Se a intensidade é forte, o uso da relação favorece acertos do modelo. Se é fraca, o modelo pouco difere de uma decisão ao acaso.

Qual o efeito de cada variável preditora na variável de decisão?

A partir da quantificação da contribuição de cada variável pode-se estudar seu efeito e o quão correta ela é.

O quão correta é a resposta do modelo?

O uso de medidas e testes estatísticos adequados permitem saber o quanto alguma resposta do modelo tem correção e a sua confiabilidade.

Há sinergia entre variáveis?

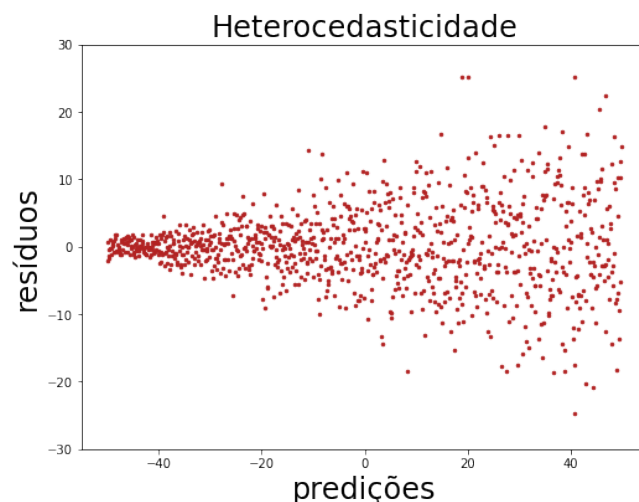
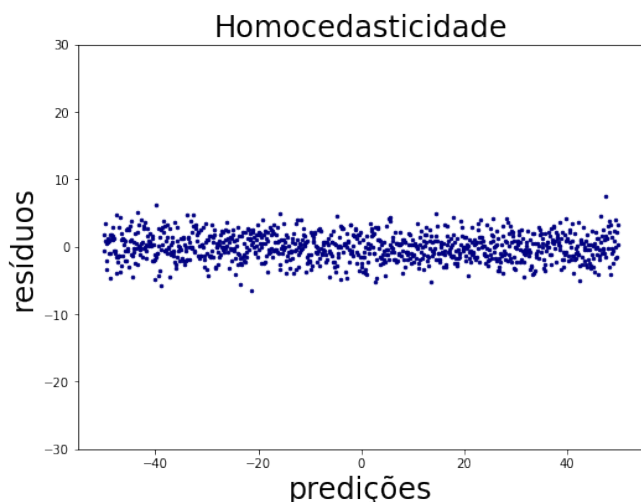
É possível que interações entre variáveis produzam efeitos mais fortes na resposta de um modelo.



## Regressão linear

Para a determinação do hiperplano que melhor descreva a relação entre as variáveis independentes e a dependente, diversos algoritmos de estimação de parâmetros podem ser utilizados e é comum que eles façam certas suposições sobre os comportamentos das variáveis:

- ausência de multicolinearidade: as variáveis preditoras não possuem fortes correlações entre si, simples ou coletivas;
- exogeneidade dos regressores: as variáveis preditoras não são correlacionadas aos erros de observação;
- linearidade de parâmetros: a variável explicada é combinação linear dos parâmetros;
- independência e normalidade de erros: os resíduos não são correlacionados e devem ser normalmente distribuídos em torno de zero;



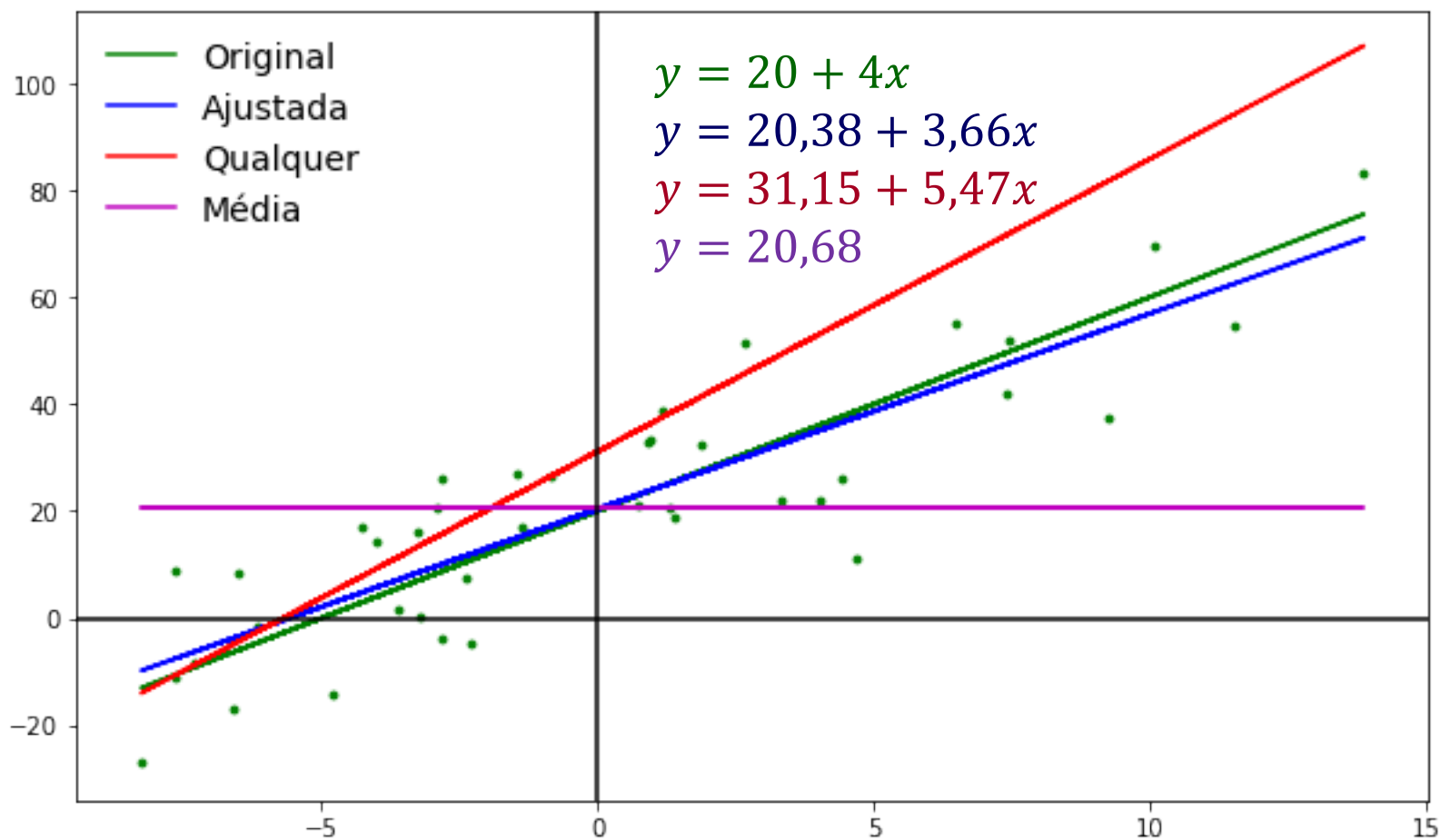
- homocedasticidade: a variância dos resíduos não depende dos valores das variáveis.

Muitas vezes, as suposições são relaxadas ou ignoradas para a aplicação de um algoritmo, promovendo a realização de operações que o tornam dependente de mais dados devido a sua maior complexidade.



## Regressão linear

Exemplo: dados sintéticos, 40 registros, 2 atributos

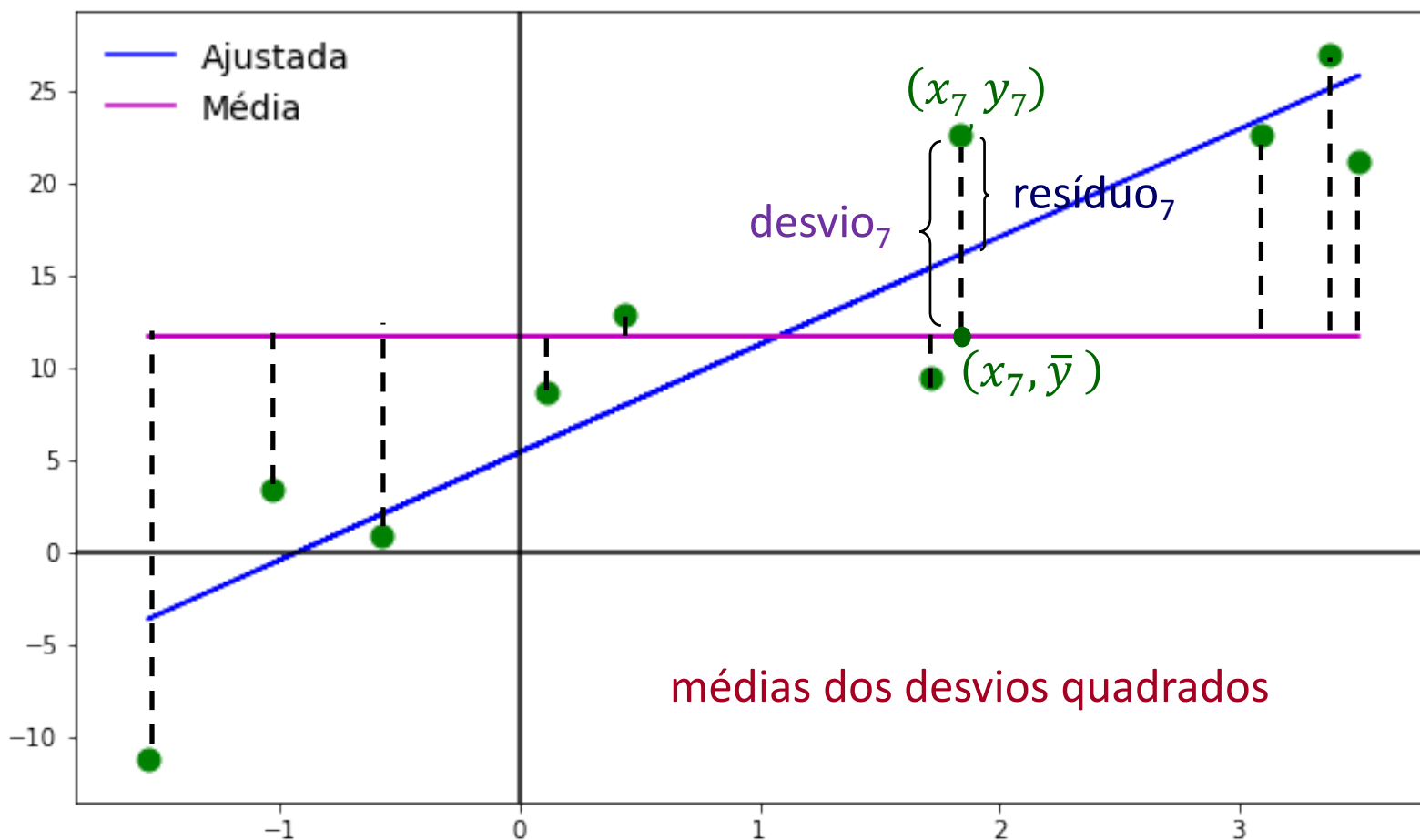


Os passos da modelagem linear podem ser assim resumidos:

- ajustar um modelo linear aos dados;
- avaliar o desempenho do modelo;
- avaliar a significância do desempenho.

## Regressão linear

Exemplo: dados sintéticos, 10 registros, 2 atributos



$$\text{Var}(y) = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n}$$

$$\text{EMQ}(\hat{y}) = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$\text{Var}(y) = 128,23$$

$$\text{EMQ}(\hat{y}) = 20,96$$

Qualquer reta diferente da ajustada terá EMQ maior que o dela.

## Regressão linear

Como o erro médio quadrático é uma função quadrática dos coeficientes de regressão, estes podem ser calculados usando as derivadas parciais da função de erro em relação a cada um e resolvidos de forma analítica. Portanto, os coeficientes ótimos são descobertos pela minimização dos resíduos quadrados:

$$\min_{\widehat{\beta}_0, \widehat{\beta}_1} EMQ(\widehat{\beta}_0, \widehat{\beta}_1) = \min_{\widehat{\beta}_0, \widehat{\beta}_1} \sum_{i=1}^N (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i)^2$$

$$\frac{\partial EMQ}{\partial \text{parâmetros}} = 0 \Rightarrow \frac{\partial EMQ}{\partial \widehat{\beta}_0} = 0 \quad \text{e} \quad \frac{\partial EMQ}{\partial \widehat{\beta}_1} = 0$$

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$$

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{cov(x, y)}{\sigma_x^2}$$



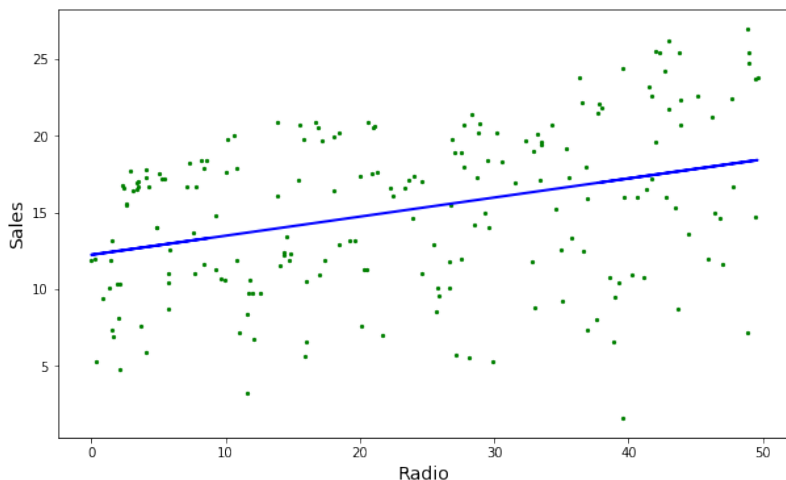
$$EMQ(\widehat{\beta}) = \frac{1}{N} (\mathbf{Y} - \widehat{\mathbf{X}}\widehat{\beta}^T)^T (\mathbf{Y} - \widehat{\mathbf{X}}\widehat{\beta}^T)$$

Forma vetorial:

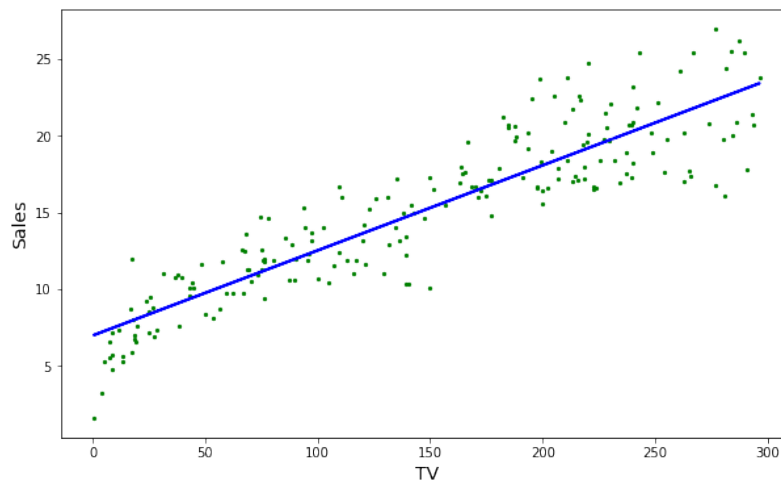
$$\widehat{\beta} = (\widehat{\mathbf{X}}^T \widehat{\mathbf{X}})^{-1} \widehat{\mathbf{X}}^T \mathbf{Y}$$

## Regressão linear

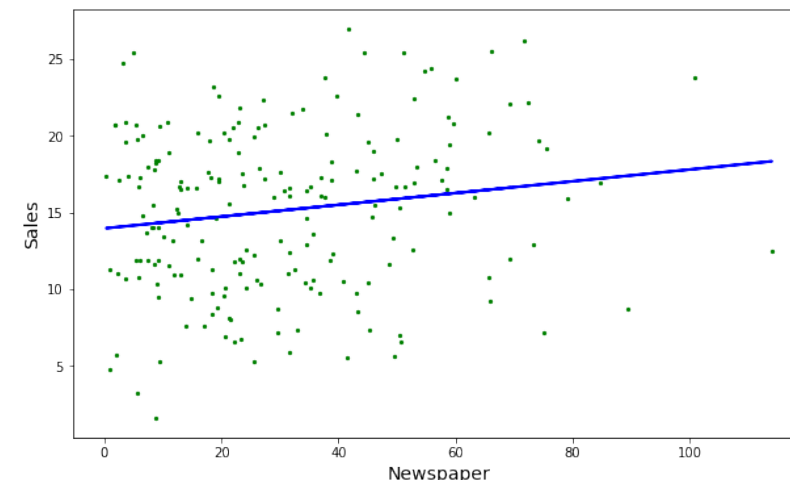
Exemplo: *Advertising dataset*, 200 registros, 4 atributos



$$Sales = 12,236 + 0,124 \text{ Radio}$$



$$Sales = 6,975 + 0,056 \text{ TV}$$



$$Sales = 13,960 + 0,038 \text{ Newspaper}$$

$$\overline{sales} = 15,13$$

$$\overline{tv} = 147,04,$$

$$cov(tv, sales) = 408,83,$$

$$\sigma_{tv}^2 = 7370,95$$

$$\overline{radio} = 23,264,$$

$$cov(radio, sales) = 27,43,$$

$$\sigma_{radio}^2 = 220,43$$

$$\overline{newspaper} = 30,554,$$

$$cov(newspaper, sales) = 18,18,$$

$$\sigma_{newspaper}^2 = 474,31$$

## Regressão linear

Pode-se confiar na correção dos parâmetros encontrados?

A média aritmética é um estimador não enviesado de um conjunto de valores. Assim, as médias dos coeficientes obtidos a partir de uma grande quantidade de regressões produzirão valores muito próximos aos dos coeficientes reais/populacionais desconhecidos.

Entretanto, não serão realizadas inúmeras regressões e pode-se estimar o **erro padrão** de cada parâmetro, que possui as seguintes utilidades:

- informar o quanto o valor do parâmetro está distante do valor real desconhecido;
- encontrar o intervalo de confiança do parâmetro;
- ser usado em testes de hipóteses sobre o parâmetro para verificar se há relação entre a variável dependente e independente.

$$SE(\widehat{\beta}_0) = \sqrt{Var(\varepsilon) \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]} \quad SE(\widehat{\beta}_1) = \sqrt{\frac{Var(\varepsilon)}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad Var(\varepsilon) \cong \frac{\sum_{i=1}^n (y_i - \widehat{y}_i)^2}{n - 2}$$



# Regressão linear

$$\left[ \widehat{\beta}_0 - 2 \cdot SE(\widehat{\beta}_0), \widehat{\beta}_0 + 2 \cdot SE(\widehat{\beta}_0) \right] \quad \text{e} \quad \left[ \widehat{\beta}_1 - 2 \cdot SE(\widehat{\beta}_1), \widehat{\beta}_1 + 2 \cdot SE(\widehat{\beta}_1) \right]$$

Seria necessário calcular a estatística  $t = (\hat{\beta}_i - \beta_i) / SE(\hat{\beta}_i)$  e consultar a distribuição de Student com  $n - 2$  graus de liberdade para se obter o valor crítico tal que:

$$P\left(\widehat{\beta}_i - t_{crítico}(n-2) \cdot SE(\widehat{\beta}_i) \leq \beta_i \leq \widehat{\beta}_i + t_{crítico}(n-2) \cdot SE(\widehat{\beta}_i)\right) = 1 - \text{nível de significância}$$

Considerando o valor típico de 95% de confiança (5% de significância), para amostras acima de 30 elementos o valor crítico é aproximadamente 2 (de 2,048, com 28 graus de liberdade, a 1,960, com infinitos graus de liberdade). Assim, adota-se 2 como uma boa aproximação.

Distribuição T-Student: valores t para os quais $P(-t \leq t \leq t) = 1 - p$																	
p =	90%	80%	70%	60%	50%	40%	30%	20%	10%	8%	6%	5%	4%	3%	2%	1%	0,5%
0	1,28	1,65	2,02	2,37	2,63	2,88	3,12	3,36	3,58	3,79	4,00	4,21	4,41	4,61	4,81	5,01	5,21
1	0,44	0,89	0,45	0,61	0,55	1,06	1,06	1,06	1,06	1,06	1,06	1,06	1,06	1,06	1,06	1,06	1,06
3	0,17	0,27	0,44	0,58	0,78	1,05	1,28	1,43	1,55	1,65	1,75	1,85	1,94	2,03	2,12	2,21	2,30
4	0,14	0,27	0,44	0,58	0,78	1,05	1,29	1,50	1,62	1,73	1,83	1,92	2,00	2,09	2,18	2,27	2,36
5	0,13	0,27	0,44	0,58	0,78	1,05	1,30	1,53	1,74	1,94	2,13	2,31	2,48	2,64	2,80	2,96	3,12
6	0,13	0,25	0,44	0,55	0,78	1,06	1,34	1,60	1,90	2,23	2,61	3,02	3,47	3,96	4,50	5,10	5,75
7	0,13	0,24	0,42	0,54	0,71	1,06	1,39	1,74	2,15	2,64	3,24	3,95	4,79	5,88	7,26	9,00	11,00
8	0,13	0,24	0,42	0,54	0,71	1,06	1,41	1,80	2,28	2,91	3,76	4,88	6,40	8,59	11,71	15,89	21,48
9	0,129	0,241	0,424	0,54	0,703	1,083	1,43	1,87	2,40	3,17	4,25	5,71	7,75	10,91	15,25	21,48	29,74
10	0,129	0,240	0,397	0,542	0,700	1,083	1,437	1,892	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
11	0,129	0,240	0,397	0,542	0,700	1,083	1,437	1,892	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
12	0,128	0,239	0,395	0,539	0,695	1,083	1,436	1,891	2,42	3,19	4,29	5,79	7,89	11,13	15,57	22,03	30,80
13	0,128	0,239	0,394	0,538	0,694	1,079	1,436	1,890	2,42	3,19	4,29	5,79	7,89	11,13	15,57	22,03	30,80
14	0,128	0,238	0,393	0,537	0,693	1,079	1,436	1,890	2,42	3,19	4,29	5,79	7,89	11,13	15,57	22,03	30,80
15	0,128	0,238	0,392	0,535	0,690	1,075	1,437	1,890	2,42	3,19	4,29	5,79	7,89	11,13	15,57	22,03	30,80
17	0,128	0,237	0,392	0,534	0,689	1,063	1,439	1,892	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
18	0,127	0,237	0,391	0,533	0,688	1,061	1,438	1,891	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
19	0,127	0,237	0,391	0,533	0,688	1,061	1,438	1,891	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
20	0,127	0,237	0,391	0,533	0,688	1,061	1,438	1,891	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
21	0,127	0,237	0,391	0,533	0,688	1,061	1,438	1,891	2,43	3,20	4,30	5,80	7,90	11,14	15,58	22,04	30,81
22	0,127	0,236	0,390	0,532	0,688	1,058	1,440	1,893	2,44	3,21	4,31	5,81	7,91	11,15	15,59	22,05	30,82
23	0,127	0,236	0,390	0,532	0,688	1,058	1,440	1,893	2,44	3,21	4,31	5,81	7,91	11,15	15,59	22,05	30,82
24	0,127	0,236	0,390	0,532	0,688	1,058	1,440	1,893	2,44	3,21	4,31	5,81	7,91	11,15	15,59	22,05	30,82
25	0,127	0,236	0,390	0,532	0,688	1,058	1,440	1,893	2,44	3,21	4,31	5,81	7,91	11,15	15,59	22,05	30,82
26	0,127	0,236	0,390	0,532	0,688	1,058	1,440	1,893	2,44	3,21	4,31	5,81	7,91	11,15	15,59	22,05	30,82
27	0,127	0,236	0,389	0,531	0,688	1,055	1,442	1,895	2,45	3,22	4,32	5,82	7,92	11,16	15,60	22,06	30,83
28	0,127	0,236	0,389	0,531	0,688	1,055	1,442	1,895	2,45	3,22	4,32	5,82	7,92	11,16	15,60	22,06	30,83
29	0,127	0,236	0,389	0,531	0,688	1,055	1,442	1,895	2,45	3,22	4,32	5,82	7,92	11,16	15,60	22,06	30,83
30	0,127	0,236	0,389	0,531	0,688	1,055	1,442	1,895	2,45	3,22	4,32	5,82	7,92	11,16	15,60	22,06	30,83
31	0,127	0,236	0,389	0,531	0,688	1,055	1,442	1,895	2,45	3,22	4,32	5,82	7,92	11,16	15,60	22,06	30,83
32	0,127	0,235	0,389	0,530	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
33	0,127	0,235	0,389	0,530	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
34	0,127	0,235	0,389	0,530	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
35	0,127	0,235	0,389	0,530	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
36	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
37	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
38	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
39	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
40	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
41	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
42	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
43	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
44	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
45	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
46	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
47	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
48	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
49	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
50	0,127	0,235	0,388	0,529	0,688	1,052	1,444	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
51	0,126	0,235	0,387	0,527	0,687	1,049	1,446	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
52	0,126	0,235	0,387	0,527	0,687	1,049	1,446	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
53	0,126	0,235	0,387	0,527	0,687	1,049	1,446	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
54	0,126	0,235	0,387	0,527	0,687	1,049	1,446	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
55	0,126	0,235	0,387	0,527	0,687	1,049	1,446	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
56	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
57	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
58	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
59	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
60	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
61	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
62	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
63	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
64	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
65	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
66	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
67	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
68	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
69	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
70	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84
71	0,126	0,234	0,387	0,526	0,687	1,046	1,448	1,897	2,46	3,23	4,33	5,83	7,93	11,17	15,61	22,07	30,84

## Regressão linear

Faz-se necessário também verificar se a variável independente associada ao parâmetro estimado é realmente relevante na explicação da variável dependente por meio de um teste de hipótese:

$H_0$ : Não há relação entre  $x$  e  $y$  ( $\beta_1 = 0$ ) e  $H_1$ : Há relação entre  $x$  e  $y$  ( $\beta_1 \neq 0$ )

É preciso verificar se  $\widehat{\beta}_1$  é suficientemente distante de zero para se ter confiança que  $\beta_1$  não seja zero. Supondo que não exista relação, ou seja,  $\beta_1 = 0$ , então  $y = \beta_0 + \varepsilon$

Calcula-se a estatística  $t = (\widehat{\beta}_1 - 0) / SE(\widehat{\beta}_1)$  para saber a quantos desvios-padrão  $\widehat{\beta}_1$  está distante de zero.

Com o nível de significância do teste e  $n - 2$  graus de liberdade, consultar a distribuição de Student para se obter o valor crítico. Se  $t$  não estiver na região crítica do teste, rejeita-se a hipótese nula. Alternativamente, o valor  $p$  da distribuição  $t$  de Student, com  $n - 2$  graus de liberdade deve ser menor que o nível de significância para a hipótese nula ser rejeitada.

Obviamente, se  $\widehat{\beta}_1$  for distante de zero e se  $SE(\widehat{\beta}_1)$  for bem pequeno, é praticamente certo de que a hipótese nula será rejeitada.



## Regressão linear

Exemplo: *Advertising dataset*, 200 registros, 4 atributos

A partir da distribuição  $t$  com 99.9% de confiança e 198 graus de liberdade, pode-se observar que o valor  $p$  do coeficiente da variável  $TV$  é menor que o nível de significância do teste (0,001). Ou seja, de maneira rigorosa, a hipótese de que não há relação entre anúncios de TV e vendas seria rejeitada.

OLS Regression Results

Dep. Variable:	Sales	R-squared:	0.812
Model:	OLS	Adj. R-squared:	0.811
Method:	Least Squares	F-statistic:	856.2
Date:	Thu, 26 May 2022	Prob (F-statistic):	7.93e-74
Time:	21:47:47	Log-Likelihood:	-448.99
No. Observations:	200	AIC:	902.0
Df Residuals:	198	BIC:	908.6
Df Model:	1		
Covariance Type:	nonrobust		

$$sales = 6,9748 + 0,0555 \cdot TV$$

	coef	std err	t	P> t	[0.025	0.975]
const	6.9748	0.323	21.624	0.000	6.339	7.611
TV	0.0555	0.002	29.260	0.000	0.052	0.059

Alternativamente, a região de aceitação de  $\beta_1$  é  $[-3,291; 3,291]$ . Portanto, como  $t$  de  $\beta_1$  está fora dela, pode-se rejeitar  $H_0$  e afirmar que não há evidências de não haver relação entre as variáveis, ou seja, pode-se sugerir fortemente que há relação entre anúncios nas estações de TV e as vendas.



## Regressão linear

Exemplo: *Advertising dataset*, 200 registros, 4 atributos

A partir da distribuição  $t$  com 99.9% de confiança e 198 graus de liberdade, pode-se observar que o valor  $p$  do coeficiente da variável *Radio* é menor que o nível de significância do teste (0,001). Ou seja, de maneira rigorosa, a hipótese de que não há relação entre anúncios de rádio e vendas seria rejeitada.

OLS Regression Results

Dep. Variable:	Sales	R-squared:	0.122
Model:	OLS	Adj. R-squared:	0.118
Method:	Least Squares	F-statistic:	27.57
Date:	Thu, 26 May 2022	Prob (F-statistic):	3.88e-07
Time:	22:05:27	Log-Likelihood:	-603.18
No. Observations:	200	AIC:	1210.
Df Residuals:	198	BIC:	1217.
Df Model:	1		
Covariance Type:	nonrobust		

$$sales = 12,2357 + 0,1244 \cdot Radio$$

	coef	std err	t	P> t	[0.025	0.975]
const	12.2357	0.653	18.724	0.000	10.947	13.524
Radio	0.1244	0.024	5.251	0.000	0.078	0.171

Alternativamente, a região de aceitação de  $\beta_1$  é  $[-3,291; 3,291]$ . Portanto, como  $t$  de  $\beta_1$  está fora dela, pode-se rejeitar  $H_0$  e afirmar que não há evidências de não haver relação entre as variáveis, ou seja, pode-se sugerir fortemente que há relação entre anúncios nas estações de rádio e as vendas.

## Regressão linear

Exemplo: *Advertising dataset*, 200 registros, 4 atributos

A partir da distribuição  $t$  com 99.9% de confiança e 198 graus de liberdade, pode-se observar que o valor  $p$  do coeficiente da variável *Newspaper* é maior que o nível de significância do teste (0,001). Ou seja, de maneira rigorosa, a hipótese de que não há relação entre anúncios de rádio e vendas seria admitida.

OLS Regression Results						
Dep. Variable:	Sales	R-squared:	0.025			
Model:	OLS	Adj. R-squared:	0.020			
Method:	Least Squares	F-statistic:	5.067			
Date:	Thu, 26 May 2022	Prob (F-statistic):	0.0255			
Time:	22:10:49	Log-Likelihood:	-613.69			
No. Observations:	200	AIC:	1231.			
Df Residuals:	198	BIC:	1238.			
Df Model:	1					
Covariance Type:	nonrobust	$sales = 13,9595 + 0,0383 \cdot Newspaper$				
	coef	std err	t	P> t	[0.025	0.975]
const	13.9595	0.638	21.870	0.000	12.701	15.218
Newspaper	0.0383	0.017	2.251	0.025	0.005	0.072

Alternativamente, a região de aceitação de  $\beta_1$  é  $[-3,291; 3,291]$ . Portanto, como  $t$  de  $\beta_1$  está dentro dela, não se pode rejeitar  $H_0$ , restando afirmar que há evidências de não haver relação entre anúncios em jornais impressos e as vendas. Entretanto, se for utilizada outro critério menos rigoroso, 95% de confiança por exemplo, sendo  $[-1,972; 1,972]$  a região de aceitação, pode-se rejeitar  $H_0$ .

## Regressão linear

O quanto um modelo se ajusta aos dados?

Avaliação da correção do modelo

O **erro padrão residual** (RSE) é uma estimativa do desvio padrão dos erros de um modelo de regressão. Mesmo que os parâmetros populacionais fossem conhecidos, ainda existiria variabilidade nos valores de  $Y$  devido a fatores não modelados. Ele representa, em média, o quanto as observações reais tendem a se afastar da reta de regressão verdadeira. Por definição, ele é o desvio-padrão dos resíduos:

$$RSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - p - 1}} \cong \sqrt{Var(\varepsilon)}$$

em que  $p$  é a quantidade de variáveis (ou parâmetros sem o intercepto) e  $n$  a quantidade de registros.

Quanto menor o seu valor, melhor o ajuste do modelo aos dados pois os valores observados estarão próximos das previsões feitas pelo modelo. Portanto, ele é medida de “desajuste” do modelo: quanto maior o seu valor, maior esse desajuste. Ele não é uma medida padronizada, mas absoluta, o que facilita a interpretação direta, mas dificulta a comparação do seu valor entre modelos ajustados para diferentes problemas.

## Regressão linear

O quanto um modelo se ajusta aos dados?

Avaliação da correção do modelo

Para medir diretamente o ajuste de um modelo aos dados, o coeficiente de determinação ( $R^2$ ) mede a proporção de variabilidade da variável de decisão que é “capturada” pela variabilidade das estimativas, ou, em outras palavras, o quanto da variabilidade da saída é explicada/removida quando se usam as variáveis preditoras.

variabilidade não explicada após a regressão

$$R^2 = \frac{Var(y) - EMQ}{Var(y)} = \frac{Var(média) - Var(modelo)}{Var(média)} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

proporção de variabilidade explicada ou removida ao se utilizar as variáveis preditoras

variabilidade antes da regressão

$$R_a^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - p - 1}$$

$N$  registros,  $p$  variáveis

## Regressão linear

O quanto um modelo se ajusta aos dados?

Avaliação da correção do modelo

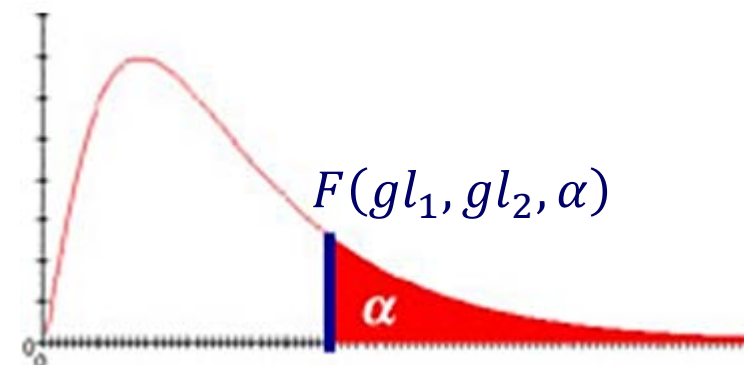
Para saber se os valores de  $R^2$  são estatisticamente significantes, calcula-se a estatística  $F$  (Fisher), de comparação de variâncias, para estimar a quantos resíduos quadrados  $R^2$  está distante de zero:

$$F = \frac{Var(média) - Var(modelo)}{Var(modelo)} \times \frac{n - k - 1}{k} = \frac{Var(y) - EMQ}{EMQ} \times \frac{n - k - 1}{k}$$

em que  $n$  é o tamanho da amostra e  $k$  é a quantidade de regressores excluindo-se o intercepto.

Dado um nível de significância  $\alpha$  e dos graus de liberdade  $gl_1 = k$  e  $gl_2 = n - k - 1$ , obter o valor crítico da distribuição  $F$ .

Se  $F > F_{crítico}$ , então o  $R^2$  é significativo.



## Regressão linear

Exemplo: *Advertising dataset*, 200 registros, 4 atributos

$$sales = 6,975 + 0,055 TV$$

$$Var(sales) = 27,92, \quad RMQ = 5,22, \quad R^2 = 0,81$$

$$F = 861,03$$

$$F_{CRÍTICO}(1; 198; 0,01) = 6,635$$

$$sales = 12,236 + 0,124 radio$$

$$Var(sales) = 27,92, \quad RMQ = 24,38, \quad R^2 = 0,12$$

$$F = 28,75$$

$$F_{CRÍTICO}(1; 198; 0,01) = 6,635$$

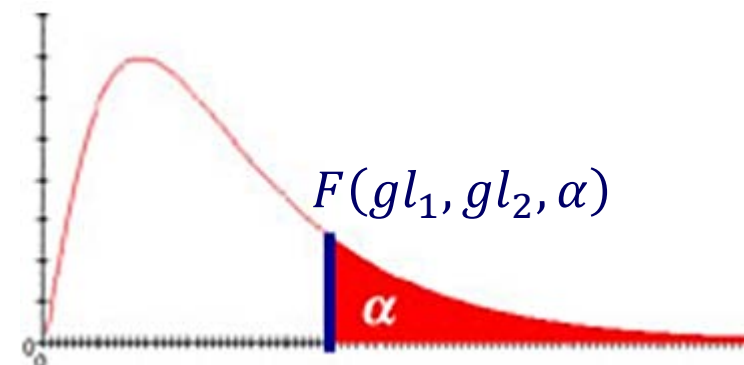
$$sales = 13,960 + 0,038 newspaper$$

$$Var(sales) = 27,92, \quad RMQ = 27,09, \quad R^2 = 0,02$$

$$F = 6,07$$

$$F_{CRÍTICO}(1; 198; 0,01) = 6,635$$

Se  $F > F_{crítico}$ , então o  $R^2$  é significativo.



## Regressão linear

O quanto um modelo se ajusta aos dados?

Avaliação da correção do modelo

O  $R^2$  ajustado penaliza modelos com mais variáveis, o que auxilia a seleção de modelos menos complexos e, conseqüentemente, menos propensos ao superajuste. Quando dois modelos apresentam valores de  $R^2$  próximos, é importante realizar testes de hipóteses adicionais para avaliar se a diferença entre eles é estatisticamente significativa.

Como alternativa mais objetiva, pode-se utilizar os critérios AIC (*Akaike Information Criterion*) e BIC (*Bayesian Information Criterion*), que, diferentemente do  $R^2$ , comparam modelos com base no equilíbrio entre qualidade de ajuste e complexidade. Nesses critérios, modelos com menores valores de AIC ou BIC são considerados preferíveis, dispensando a necessidade de testes formais adicionais.

Ambos seguem a lógica:

$$\text{Critério} = \text{Erro do modelo} + \text{Penalidade pela complexidade}$$

O erro do modelo é medido a partir da verossimilhança ( $P(\text{Dados Observados}|\text{Parâmetros})$ ), em que quanto maior o valor, melhor o ajuste). A complexidade é medida pela quantidade de parâmetros usados.



## Regressão linear

O quanto um modelo se ajusta aos dados?

Avaliação da correção do modelo

*Akaike Information Criterion:*  $AIC(M) = -2 \cdot \ln \mathcal{L}(M) + 2 \cdot p$

modelo  $M$ ,  $p$  parâmetros,  $n$  registros

*Bayesian Information Criterion:*  $BIC(M) = -2 \cdot \ln \mathcal{L}(M) + p \cdot \ln n$

BIC é mais rigoroso que o AIC

$\mathcal{L}(M)$  é a função de máxima verossimilhança dos parâmetros do modelo, geralmente gaussiana. Dada a complexidade de cálculos e a rápida degeneração da função de verossimilhança provocada pelo produto de probabilidades, as versões logarítmicas são as adotadas:

$$\mathcal{L}(M) = \prod_{i=1}^n P(y_i | \mathbf{x}_i; \hat{\boldsymbol{\beta}}, s^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi \cdot s^2}} e^{-\frac{(\text{resíduo}_i)^2}{2s^2}}$$

AIC tende a ser menos severo que o BIC quando um modelo usa mais parâmetros, ou seja, pode resultar em um valor menor (melhor) para modelos mais complexos desde que a melhoria no ajuste compense a penalização. No BIC, a penalidade cresce muito rapidamente, ultrapassando o valor de AIC a partir de 7 registros!



## Regressão linear múltipla

Modelo linear estimado:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_n x_n$

- modelo estimado da aproximação  $y \approx \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$  a partir dos dados disponíveis para gerá-lo;
- $\hat{\beta}_j$  pode ser visto como o efeito médio da variável  $x_j$  em  $\hat{y}$  quando todos os demais parâmetros estão fixos;
- o uso de mais variáveis permite que o modelo se ajuste mais aos dados do que o modelo linear simples, independentemente da relevância estatística da variável.

Exemplo:

*Advertising dataset,*  
200 registros,  
4 atributos

$\alpha = 0,05$   
 $t_{crítico} = 1,653, gl = 196$   
 $F_{crítico} = 2,651$   
 $gl1 = 3, gl2 = 196$

$\alpha = 0,05$   
 $t_{crítico} = 1,653, gl = 197$   
 $F_{crítico} = 3,042 gl1 = 2, gl2 = 197$

Dep. Variable:	Sales	R-squared:	0.903
Model:	OLS	Adj. R-squared:	0.901
Method:	Least Squares	F-statistic:	605.4
Date:	Thu, 26 May 2022	Prob (F-statistic):	8.13e-99
Time:	21:30:29	Log-Likelihood:	-383.34
No. Observations:	200	AIC:	774.7
Df Residuals:	196	BIC:	787.9
Df Model:	3		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	4.6251	0.308	15.041	0.000	4.019	5.232
TV	0.0544	0.001	39.592	0.000	0.052	0.057
Radio	0.1070	0.008	12.604	0.000	0.090	0.124
Newspaper	0.0003	0.006	0.058	0.954	-0.011	0.012

Dep. Variable:	Sales	R-squared:	0.903
Model:	OLS	Adj. R-squared:	0.902
Method:	Least Squares	F-statistic:	912.7
Date:	Thu, 26 May 2022	Prob (F-statistic):	2.39e-100
Time:	21:33:56	Log-Likelihood:	-383.34
No. Observations:	200	AIC:	772.7
Df Residuals:	197	BIC:	782.6
Df Model:	2		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	4.6309	0.290	15.952	0.000	4.058	5.203
TV	0.0544	0.001	39.726	0.000	0.052	0.057
Radio	0.1072	0.008	13.522	0.000	0.092	0.123

## Regressão linear múltipla

Exemplo: Base de dados sobre os bairros do Rio de Janeiro, 152 registros, 24 atributos

OLS Regression Results

```
=====
Dep. Variable:    Indice desenvolvimento social    R-squared:                0.664
Model:            OLS                            Adj. R-squared:           0.661
Method:           Least Squares                  F-statistic:              296.0
Date:            Thu, 26 May 2022                 Prob (F-statistic):       2.56e-37
Time:            17:17:04                         Log-Likelihood:           124.05
No. Observations: 152                            AIC:                      -244.1
Df Residuals:     150                            BIC:                      -238.1
Df Model:         1
Covariance Type:  nonrobust
=====
```

```
=====
               coef    std err          t      P>|t|      [0.025    0.975]
-----
const                0.8916     0.028    31.422     0.000     0.836     0.948
Densidade domiciliar -0.7714     0.045   -17.204     0.000    -0.860    -0.683
=====
```

$\alpha = 0,05$

$t_{crítico} = 1,655, gl = 150$

$F_{crítico} = 3,904, gl1 = 1, gl2 = 150$

OLS Regression Results

```
=====
Dep. Variable:    Indice desenvolvimento social    R-squared:                0.684
Model:            OLS                            Adj. R-squared:           0.679
Method:           Least Squares                  F-statistic:              161.0
Date:            Thu, 26 May 2022                 Prob (F-statistic):       5.78e-38
Time:            19:43:27                         Log-Likelihood:           128.70
No. Observations: 152                            AIC:                      -251.4
Df Residuals:     149                            BIC:                      -242.3
Df Model:         2
Covariance Type:  nonrobust
=====
```

```
=====
               coef    std err          t      P>|t|      [0.025    0.975]
-----
const                0.7674     0.049    15.663     0.000     0.671     0.864
Densidade domiciliar -0.6328     0.063   -10.073     0.000    -0.757    -0.509
Grau de verticalizacao 0.2325     0.076     3.068     0.003     0.083     0.382
=====
```

$\alpha = 0,05$

$t_{crítico} = 1,655, gl = 150$

$F_{crítico} = 3,057, gl1 = 2, gl2 = 149$

## Regressão linear múltipla

Exemplo: Base de dados sobre os bairros do Rio de Janeiro, 152 registros, 24 atributos

OLS Regression Results						
=====						
Dep. Variable:	Indice desenvolvimento social		R-squared:	0.707		
Model:	OLS		Adj. R-squared:	0.701		
Method:	Least Squares		F-statistic:	118.8		
Date:	Thu, 26 May 2022		Prob (F-statistic):	3.27e-39		
Time:	19:53:32		Log-Likelihood:	134.41		
No. Observations:	152		AIC:	-260.8		
Df Residuals:	148		BIC:	-248.7		
Df Model:	3					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
const	0.8745	0.057	15.374	0.000	0.762	0.987
Densidade domiciliar	-0.7065	0.064	-10.958	0.000	-0.834	-0.579
Grau de verticalizacao	0.2322	0.073	3.171	0.002	0.088	0.377
Percentual de imoveis alugados	-0.1645	0.048	-3.397	0.001	-0.260	-0.069

$$\alpha = 0,05$$

$$t_{\text{crítico}} = 1,655, gl = 150$$

$$F_{\text{crítico}} = 2,666, gl1 = 3, gl2 = 148$$

$$\alpha = 0,05$$

$$t_{\text{crítico}} = 1,655, gl = 150$$

$$F_{\text{crítico}} = 2,433, gl1 = 4, gl2 = 147$$

OLS Regression Results						
Dep. Variable:	Indice desenvolvimento social	R-squared:	0.708			
Model:	OLS	Adj. R-squared:	0.700			
Method:	Least Squares	F-statistic:	89.04			
Date:	Thu, 26 May 2022	Prob (F-statistic):	2.80e-38			
Time:	21:15:05	Log-Likelihood:	134.75			
No. Observations:	152	AIC:	-259.5			
Df Residuals:	147	BIC:	-244.4			
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	0.8668	0.058	15.010	0.000	0.753	0.981
Densidade domiciliar	-0.7182	0.066	-10.859	0.000	-0.849	-0.587
Grau de verticalizacao	0.2241	0.074	3.028	0.003	0.078	0.370
Percentual de imoveis alugados	-0.1752	0.050	-3.487	0.001	-0.274	-0.076
Percentual de area urbanizada	0.0270	0.033	0.811	0.419	-0.039	0.093

## Regressão linear múltipla

Como lidar com interações entre variáveis?

Supondo que uma variável  $x_k$  tenha influência em outra variável  $x_l$ , pode-se relaxar a suposição aditiva dos modelos lineares:

$$\hat{y} = \widehat{\beta}_0 + \cdots + \widehat{\beta}_k x_k + \widehat{\beta}_l x_l + \boxed{\widehat{\beta}_{n+1} x_k x_l} + \cdots + \widehat{\beta}_n x_n$$

termo de interação

Reescrevendo o modelo:

$$\hat{y} = \widehat{\beta}_0 + \cdots + \widehat{\beta}_k x_k + (\widehat{\beta}_l + \widehat{\beta}_{n+1} x_k) x_l + \cdots + \widehat{\beta}_n x_n$$

Quando há interação, os coeficientes isolados das variáveis envolvidas deixam de ter interpretação direta.

Após o cálculo dos parâmetros, o teste de hipótese a respeito da relevância de  $\widehat{\beta}_{n+1}$  pode indicar que a interação é estatisticamente insignificante e que o termo pode ser excluído da modelagem, ou seja, o modelo resultante é simplificado. Atenção: nunca excluir os termos originais ao manter o termo de interação.

## Regressão linear múltipla

Como lidar com interações entre variáveis?

Interações como a exemplificada propiciam a ocorrência de multicolineariedade entre as variáveis envolvidas. Centralizar as variáveis originais ajuda a reduzir o problema sem afetar o desempenho da modelagem. O impacto será apenas nos valores dos parâmetros ajustados e nos erros padrão, que ficarão mais estáveis.

Outras não linearidades podem ser incorporadas, caracterizando não linearidades entre preditoras e preditas, e o modelo continua linear nos parâmetros. Por exemplo,  $\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x_1^3 + \widehat{\beta}_2 x_2 + (\widehat{\beta}_3 + \widehat{\beta}_4 x_2) x_3$  é um modelo linear chamado de polinômio de regressão com termo de interação.

Neste caso, pode-se fazer uma mudança de variável para melhor observar a estrutura linear:  $\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 z_1 + \widehat{\beta}_2 z_2 + \widehat{\beta}_5 z_3$ . Entretanto,  $\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x_1^{\beta_0}$ , por exemplo, não é linear, ou seja, não é possível reescrevê-la em uma estrutura linear. Neste caso, a regressão linear não seria a modelagem a ser adotada.

Assim, pode-se usar a modelagem linear para um ajuste não-linear aos dados desde que permaneça linear nos coeficientes. Não linearidades entre variáveis não constituem um problema.

## Regressão linear múltipla

Quais variáveis são realmente importantes para o modelo?

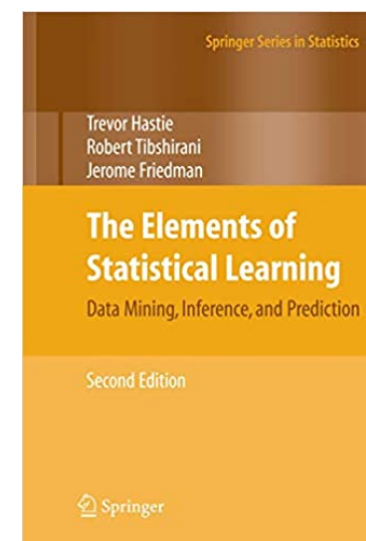
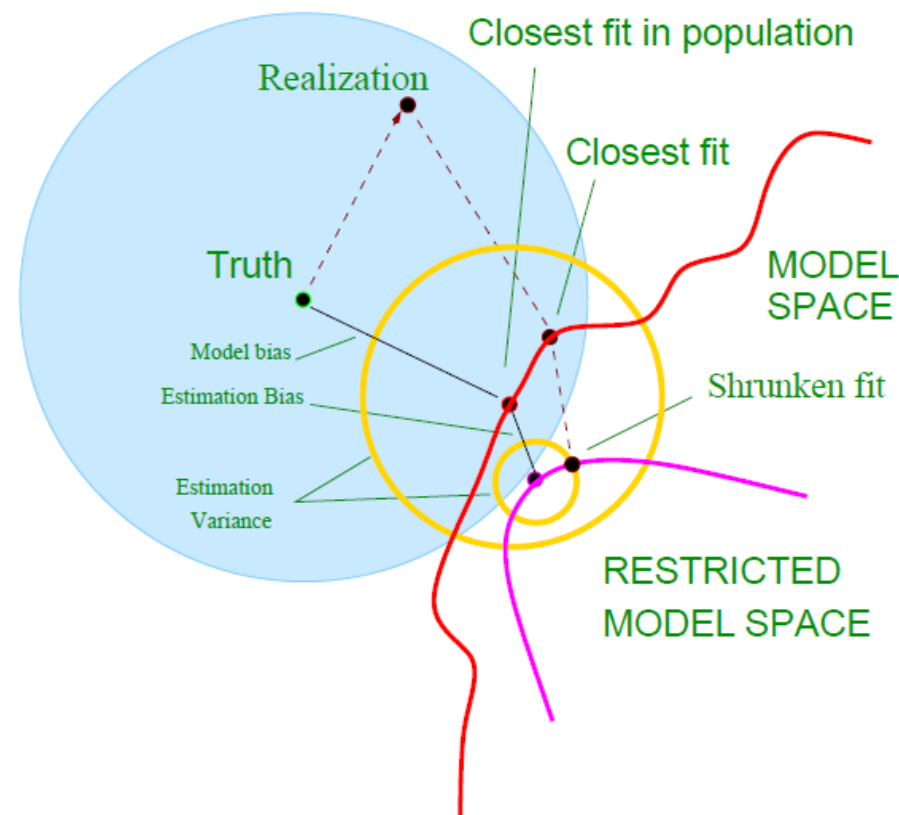
Pode-se adotar uma estratégia de seleção de variáveis:

- *Progressiva/Forward*: a partir de um modelo básico, apenas com o intercepto, com a lista de regressores vazia, executam-se  $p$  modelos lineares simples e adiciona-se ao modelo básico a variável com menor EMQ; e, assim, mais variáveis são adicionadas ao último modelo obtido até que algum critério de parada seja atingido;
- *Regressiva/Backward*: a partir de um modelo múltiplo completo, exclui-se a variável com maior valor  $p$  do teste  $t$ ; assim, variáveis são excluídas do último modelo produzido até que um critério de parada seja atingido;
- *Híbrida*: começa como a progressiva e age como a regressiva caso alguma adição tenha causado aumento de valor  $p$  acima do tolerado.

## Regressão linear múltipla com regularização

Como estratégia alternativa à busca pelo dimensionamento adequado da modelagem, pode-se utilizar todas as variáveis regressoras para a geração de um modelo e adotar uma técnica que restringe ou regulariza os coeficientes do modelo.

Chamada de regularização, essa técnica torna menor, menos complexo, o espaço de busca de parâmetros de modelos e promove o encontro de modelos com melhores capacidades de generalização e com variâncias bem menores do que os que seriam obtidos sem ela.



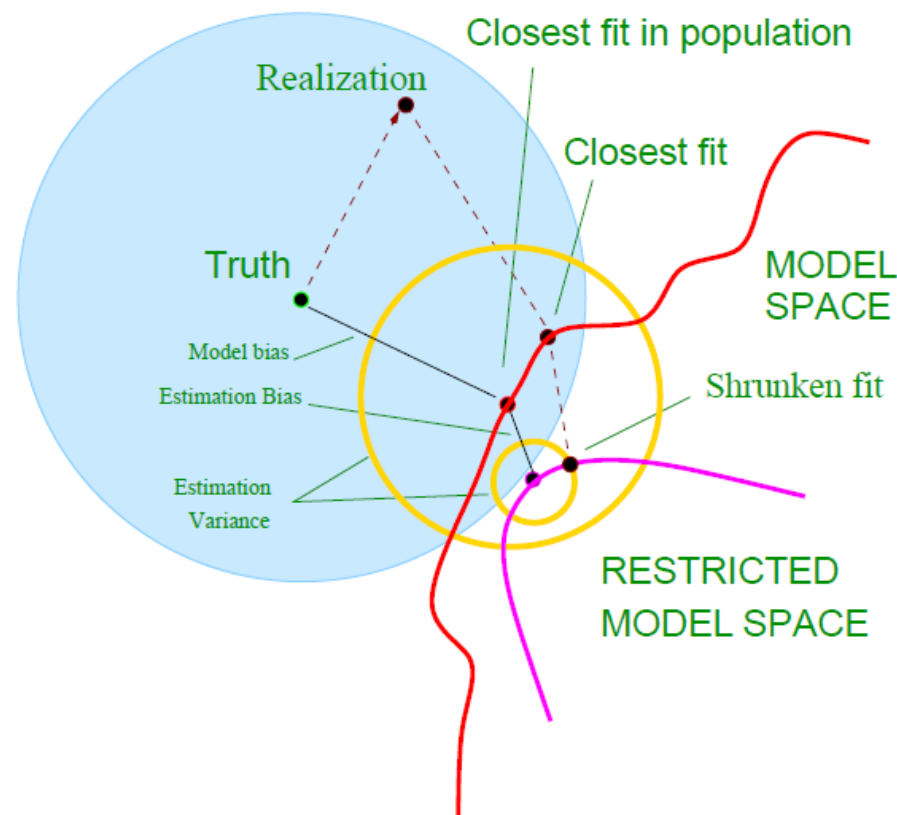


## Regressão linear múltipla com regularização

Como não se conhece a real função geradora dos dados, o modelo produzido possui um erro de predição e minimizá-lo constitui um problema **mal posto** ou **mal condicionado**. Assim, resumidamente, a regularização procura tornar o problema da minimização do erro de predição em um problema menos mal posto.

Para ser **bem posto** ou **bem condicionado**, é preciso atender aos seguintes critérios:

- Existir uma função geradora e o modelo ser capaz de aproximá-la;
- A solução da minimização ser única;
- A solução deve ser estável, isto é, deve haver continuidade de comportamento diante de diferentes condições iniciais e de contorno.

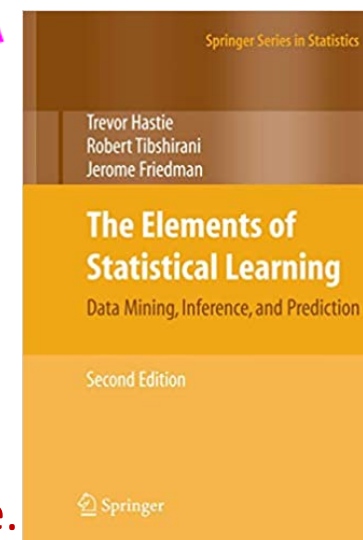
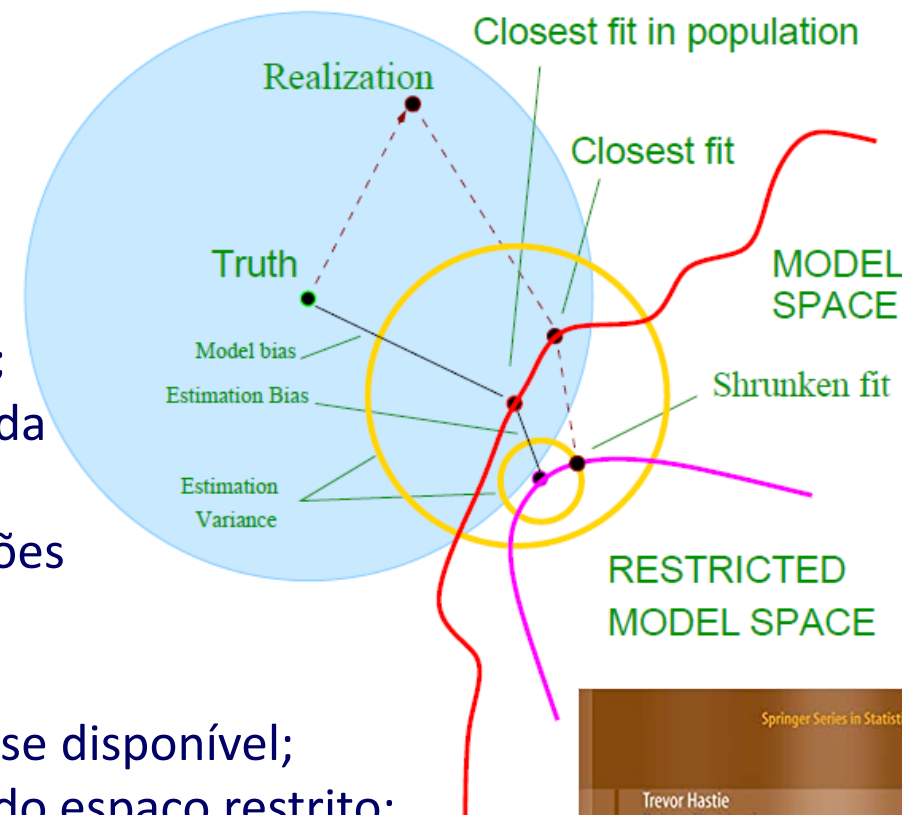




## Regularização

Na figura, pode-se fazer a seguinte interpretação:

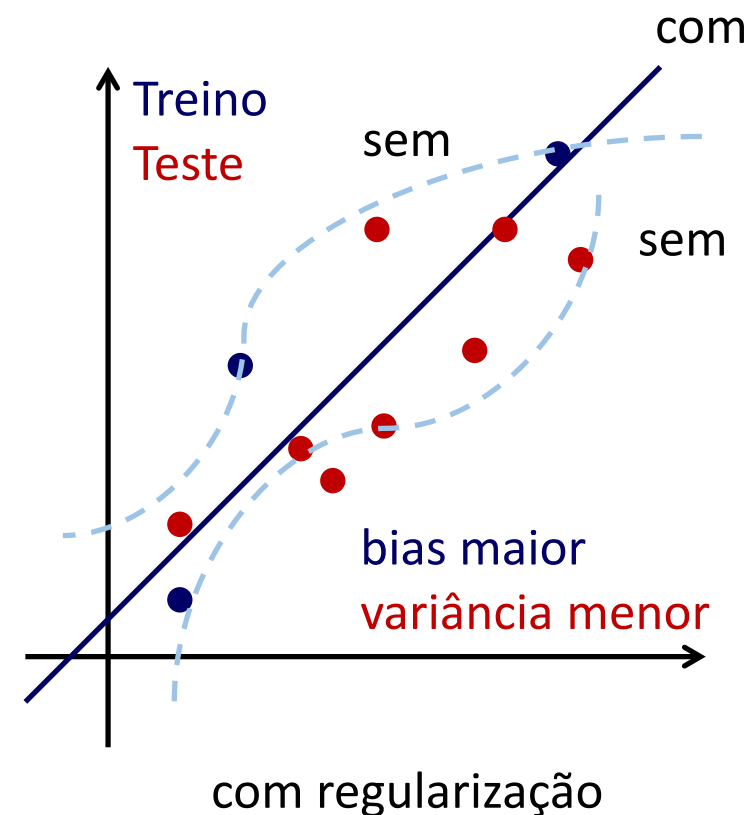
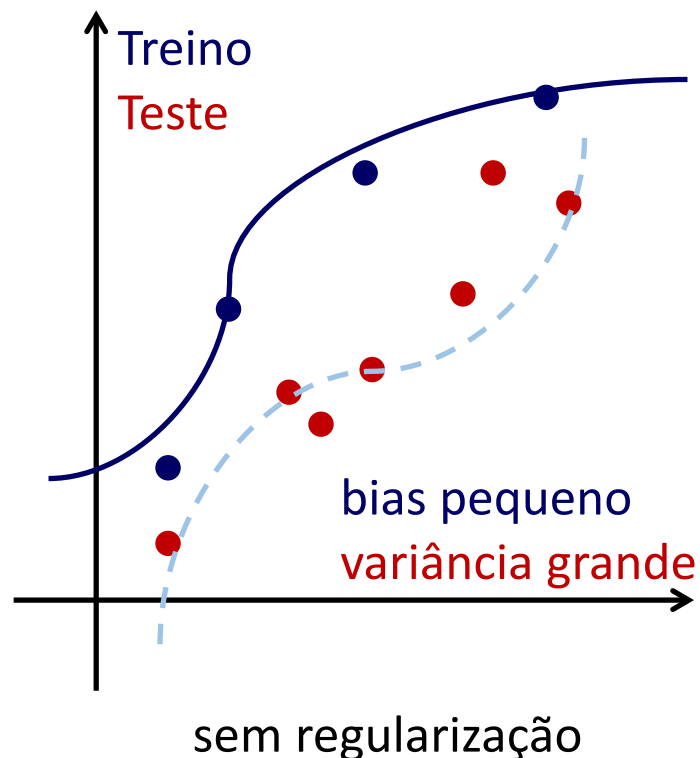
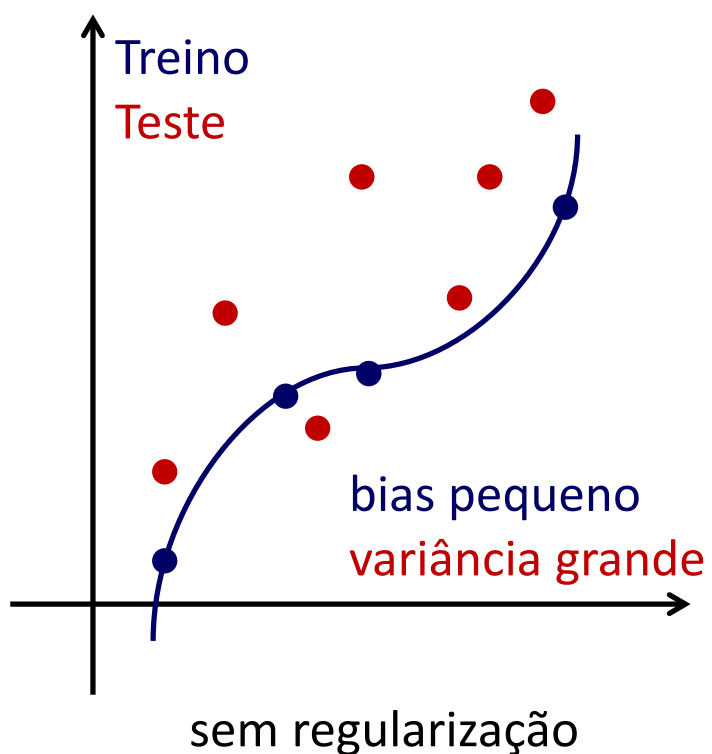
- *Truth* é a verdade dos dados, nem sempre observada;
- *Realization* é verdade observada, possivelmente com erros na coleta;
- *Model Space* é o conjunto de todas as modelagens possíveis a partir da verdade observada;
- *Restricted Model Space* é o subconjunto das modelagens com restrições (menos variáveis, domínios reduzidos, uso de regularização etc.);
- *Closest fit* é o melhor modelo ajustado aos dados observados;
- *Closest fit in population* é o modelo ideal se toda a população estivesse disponível;
- *Shrunk fit* é o modelo final ajustado aos dados observados dentro do espaço restrito;
- *Estimation Variance* refere-se à variabilidade de modelos com outros dados observados (quanto maior a complexidade, maior a variabilidade);
- *Model Bias* é o erro entre a predição do melhor modelo teórico e a verdade real;
- *Estimation Bias* é o erro adicional que surge por restringir o espaço de modelos.



Modelagens complexas podem se ajustar demais aos dados observados, tendo alta variabilidade.

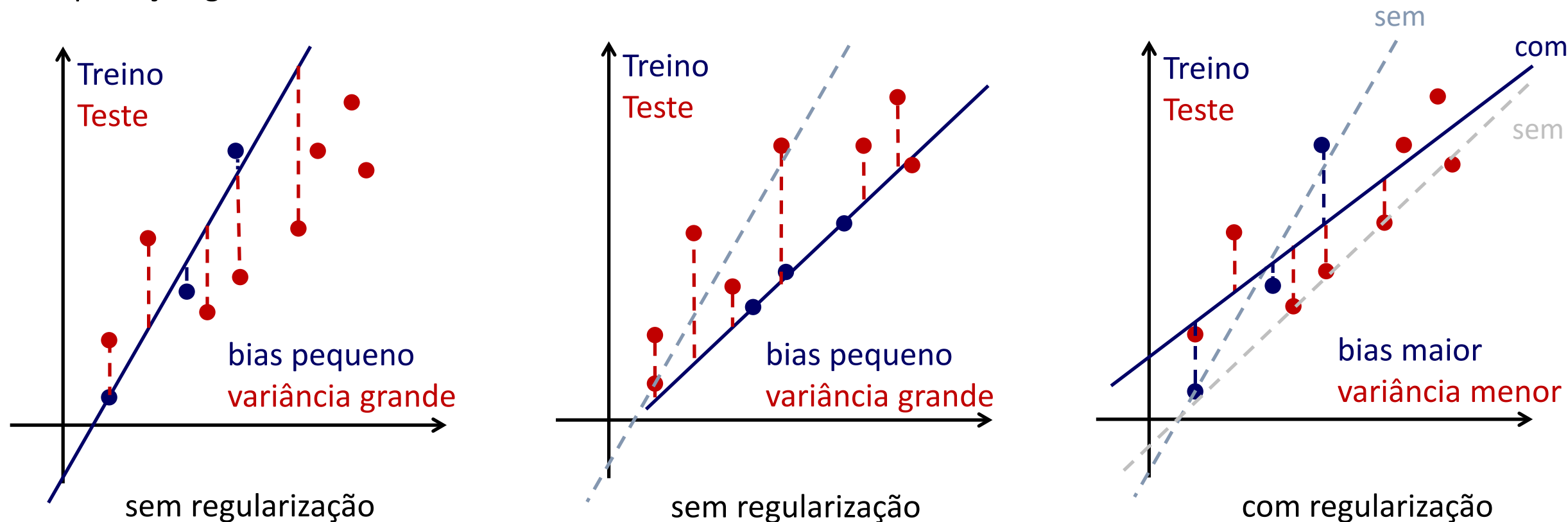
## Regularização

Interpretação geométrica de um superajuste



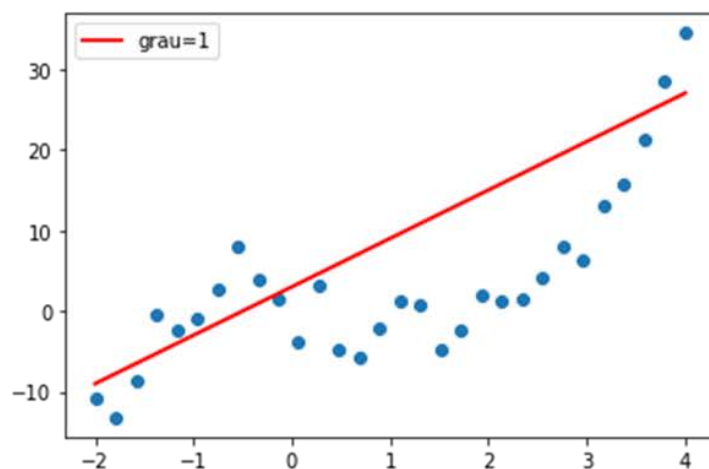
## Regressão linear múltipla com regularização

Interpretação geométrica

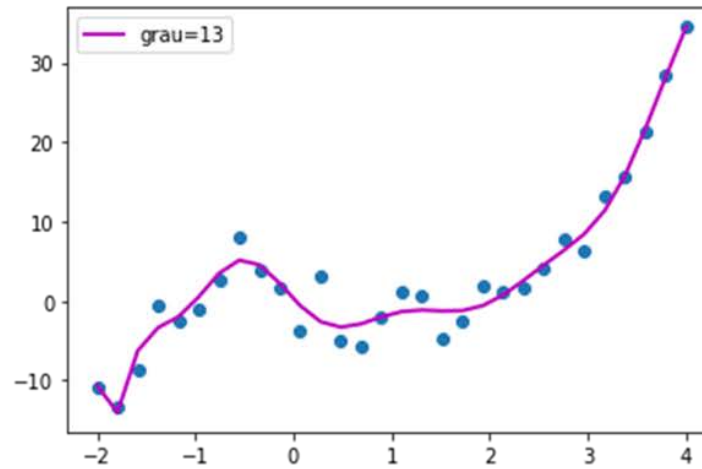


## Regressão linear múltipla com regularização

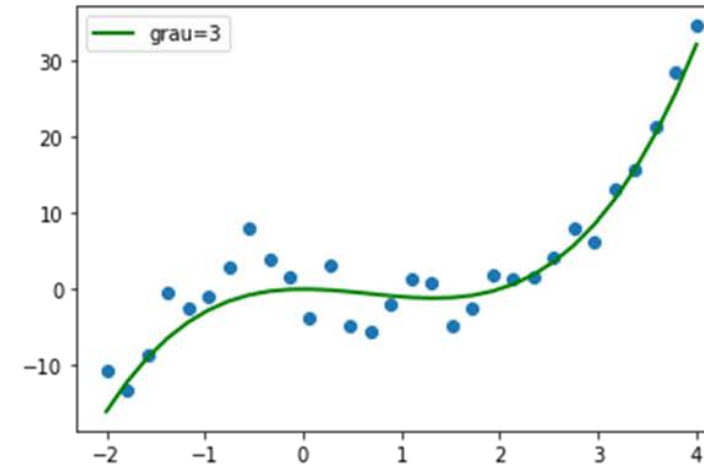
Uma observação empírica: um dimensionamento inadequado tende a provocar superajuste ou subajuste.



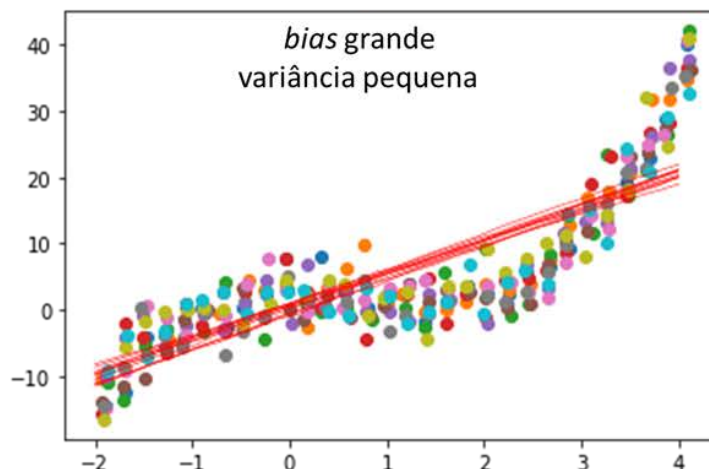
Subajustado



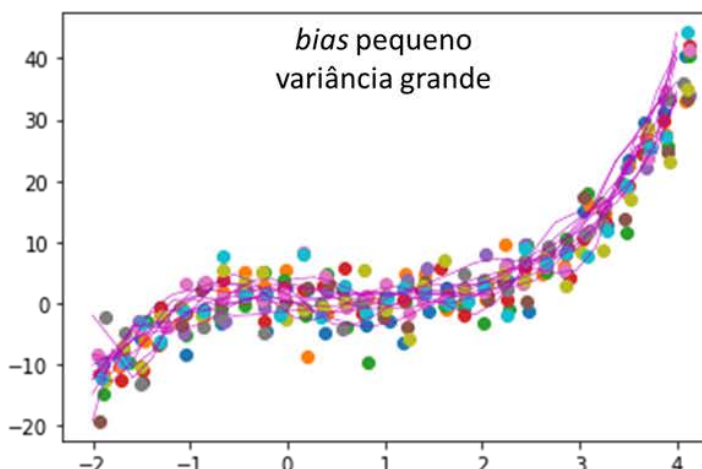
Superajustado



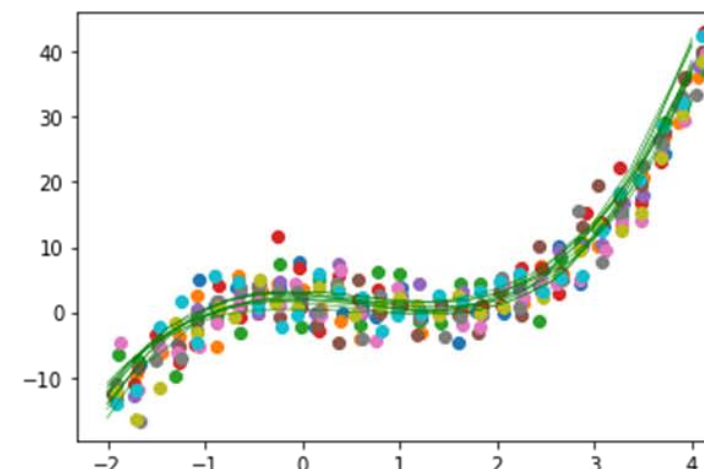
Bem ajustado



*bias grande*  
*variância pequena*

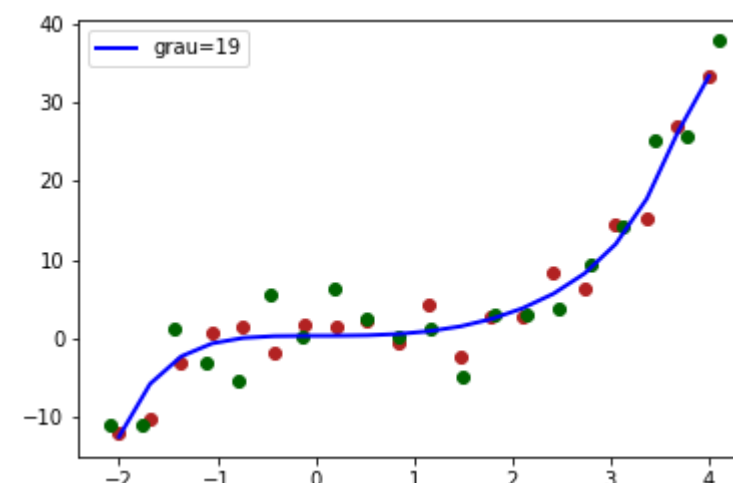
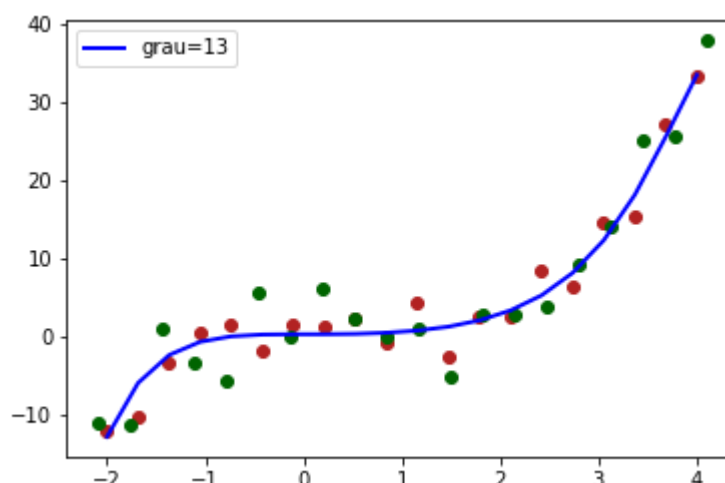
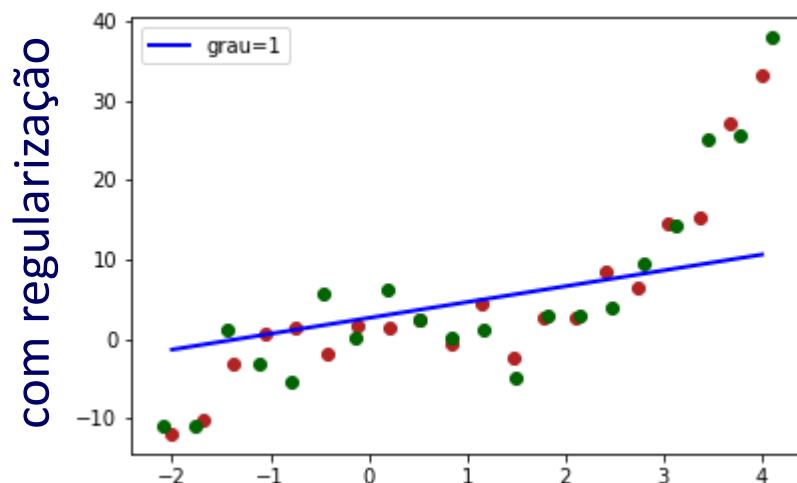
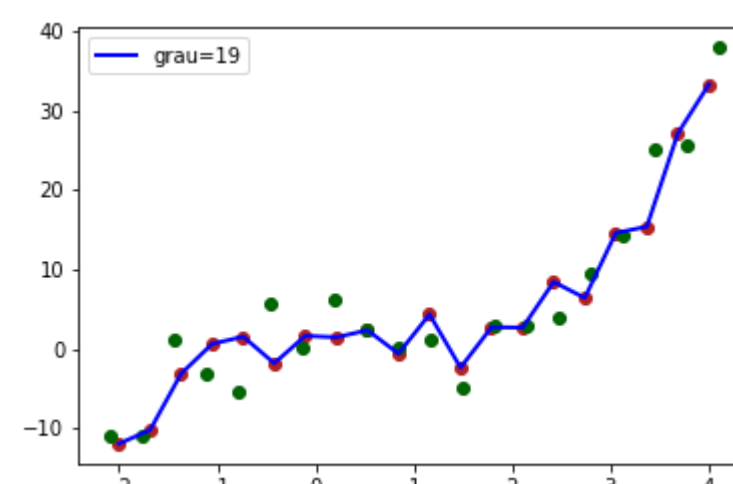
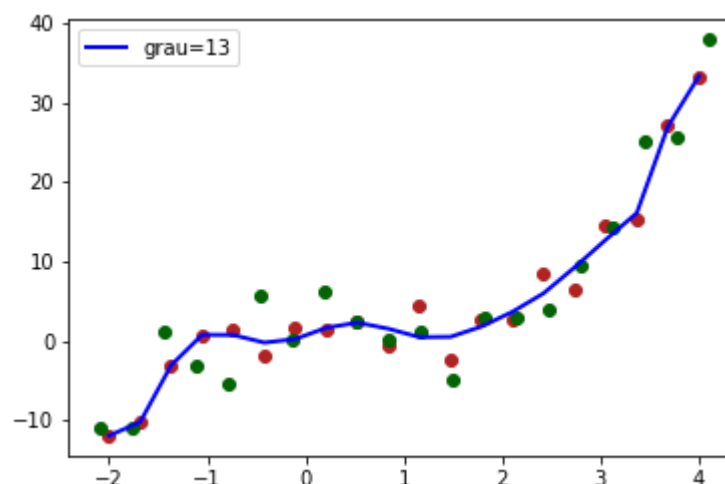
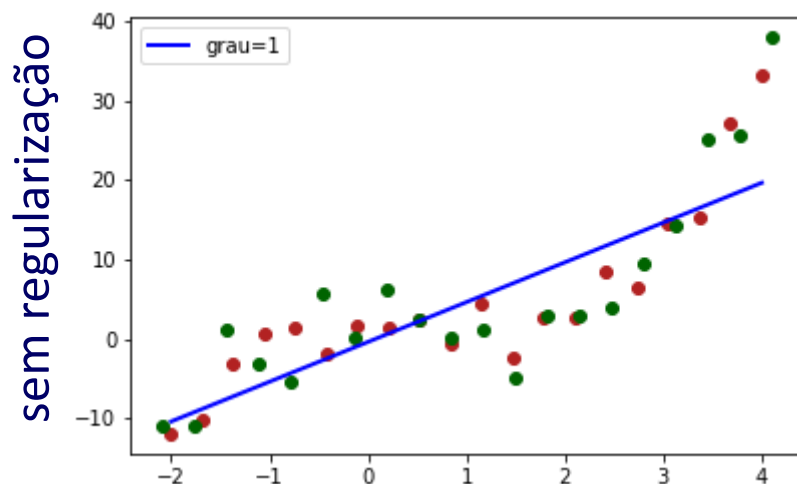


*bias pequeno*  
*variância grande*



## Regressão linear múltipla com regularização

Uma observação empírica: com regularização, um superdimensionamento pode ser corrigido.

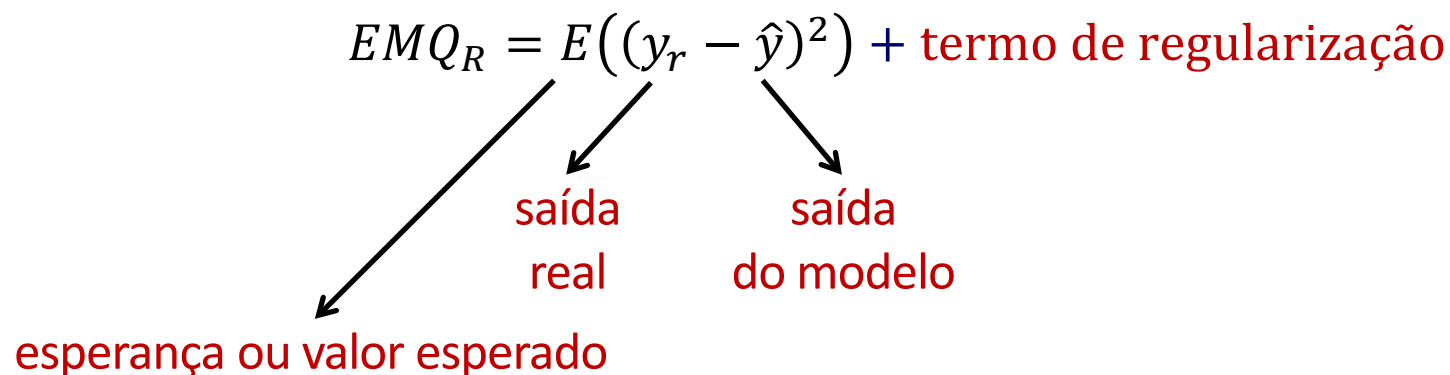


## Regressão linear múltipla com regularização

Mais observações empíricas:

- O superajuste está frequentemente associado a valores mais altos de coeficientes e é difícil controlá-los durante o aprendizado sem intervenção;
- Em muitos problemas, a quantidade de dados disponíveis para o aprendizado é insuficiente para a adequada geração de modelos.

Assim, as técnicas de regularização mais comuns procuram evitar o superajuste mantendo os coeficientes com valores pequenos por meio de uma penalização proporcional à soma dos valores **absolutos** ou **quadrados** de todos os coeficientes de um modelo de regressão. Muitas vezes, ambas as formas de penalização são adotadas.

$$EMQ_R = E((y_r - \hat{y})^2) + \text{termo de regularização}$$


saída real

saída do modelo

esperança ou valor esperado

Ao se reduzir os domínios dos coeficientes de um modelo, alguns eventualmente poderão assumir valores nulos, o que reduz a dimensão do modelo, e promove-se uma menor variância de coeficientes em diferentes amostragens, o que resulta em menor variância das respostas de diferentes modelos.

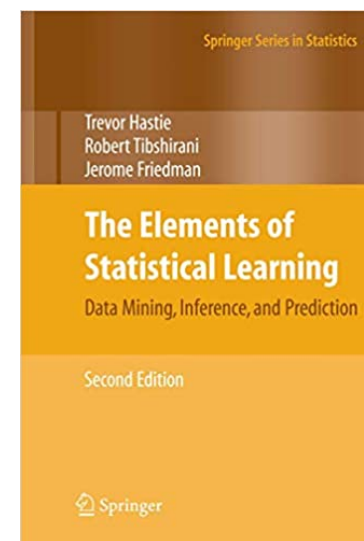
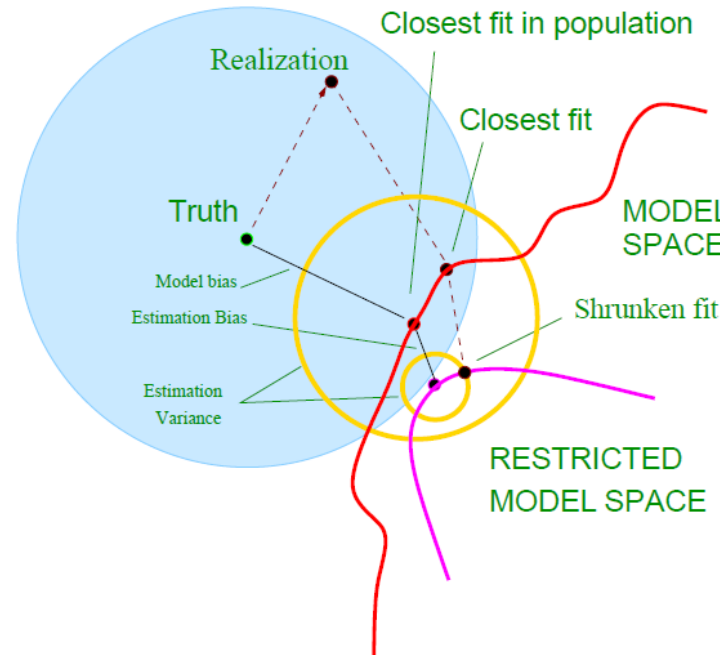
$$EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot \Sigma_R \rightarrow \text{termo de regularização}$$

$\Sigma_R$  é o regularizador

$\lambda \in [0, +\infty)$  é arbitrado e pode ser escolhido por validação cruzada.

O efeito do termo de regularização é o de redução de possibilidades de representação das funções geradoras dos dados, tornando-as mais “similares”.

Assim, o aprendizado se torna mais estável diante de diferentes dados de treinamento apresentados.





## Regressão linear múltipla com regularização

A regularização abrange todas as técnicas que adotam termos de penalização à função de erro/custo/perda de predição para evitar o superajuste.

Ao escolher o parâmetro de regularização de um modelo, é preciso encontrar um bom equilíbrio entre simplicidade e ajuste:

- se o valor for alto demais, o modelo tenderá a ser muito simples e correrá o risco de subajuste;
- se o valor for pequeno demais, o modelo tenderá a ser mais complexo do que o necessário, ocorrendo o risco de superajuste.

**L2/Tikhonov:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum \text{parâmetros}^2)$  (ridge, de cume, de crista)

**L1/LASSO:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum |\text{parâmetros}|)$  (Least Absolute Shrinkage and Selection Operator)

**Elastic net:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda_1 \cdot (\sum \text{parâmetros}^2) + \lambda_2 \cdot (\sum |\text{parâmetros}|)$  (rede elástica)

## Regressão linear múltipla com regularização

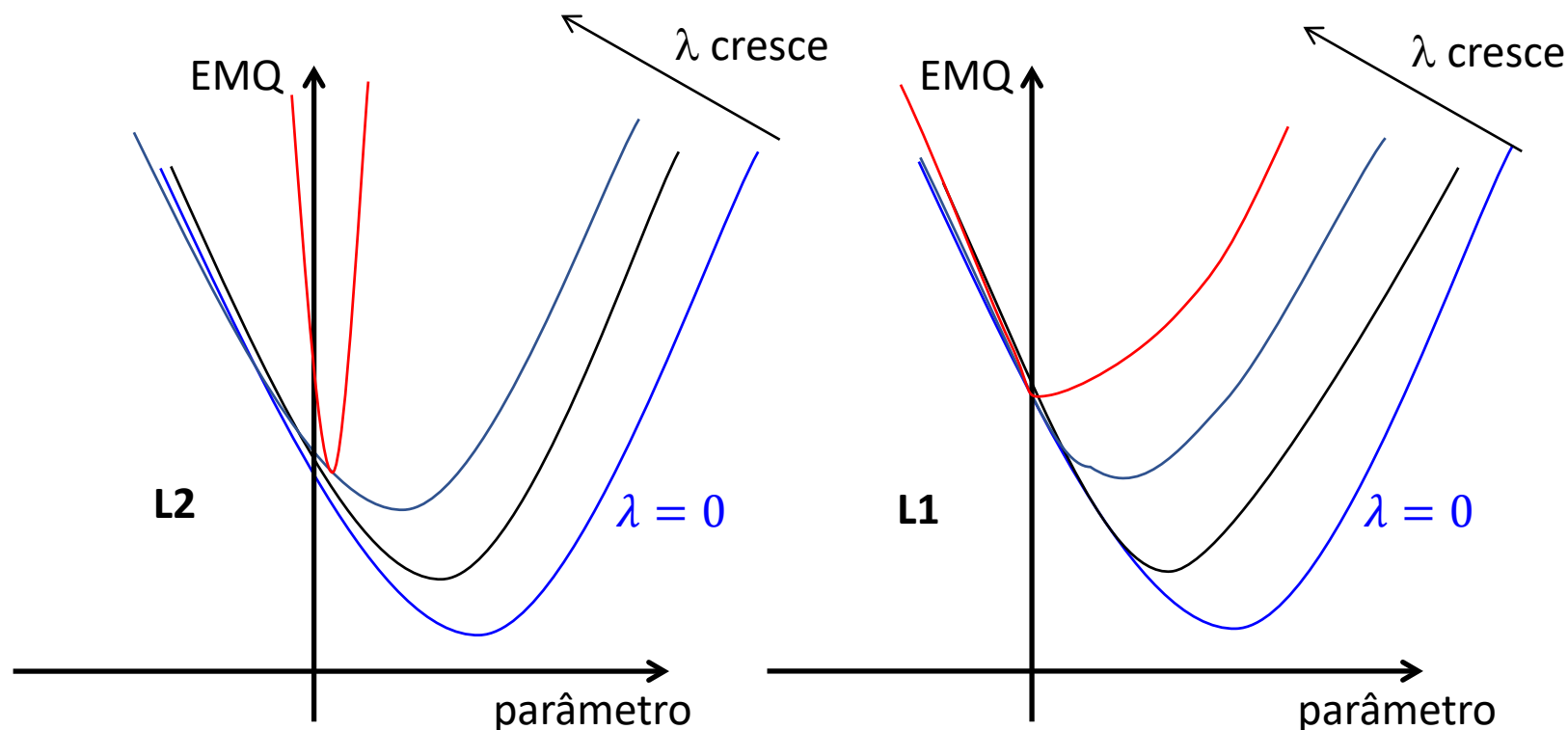
**L2/Tikhonov:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum \text{parâmetros}^2)$  (encolhimento, redução)

**L1/LASSO:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum |\text{parâmetros}|)$  (encolhimento e seleção de atributos)

**Elastic net:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda_1 \cdot (\sum \text{parâmetros}^2) + \lambda_2 \cdot (\sum |\text{parâmetros}|)$  (ambos)

Na *ridge*, os parâmetros podem se aproximar de zero e na LASSO eles podem alcançá-lo, significando a remoção das variáveis a que os coeficientes nulos estão associados.

Assim, se há certeza que todas as variáveis são relevantes, *ridge* é a escolha a ser feita, por não realizar qualquer seleção.



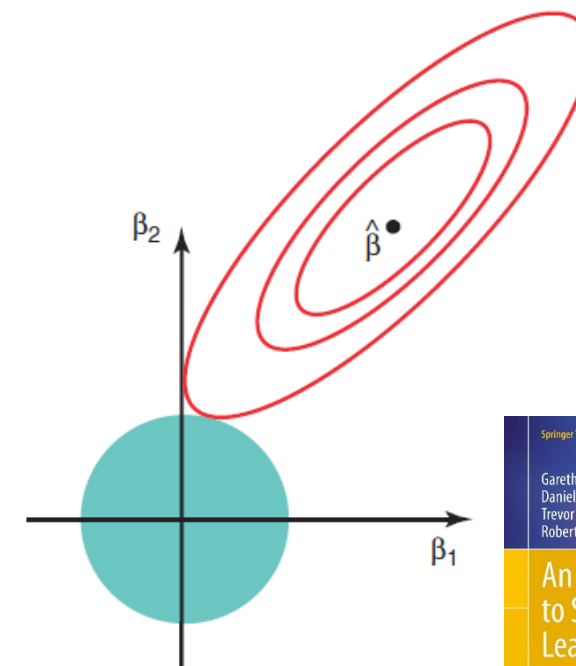
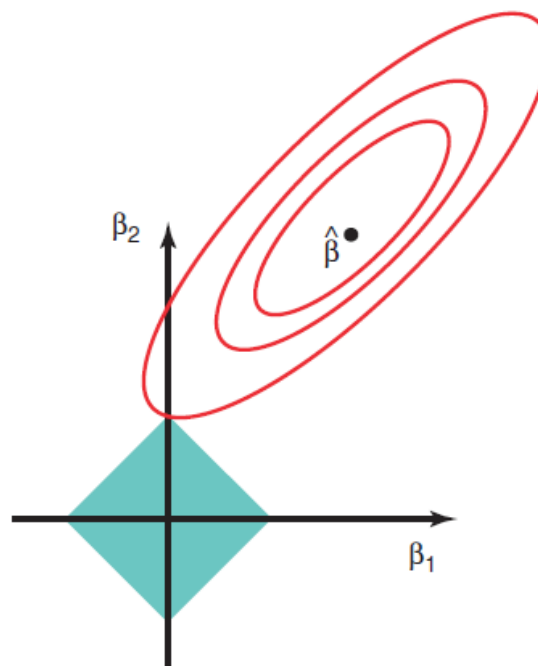
## Regressão linear múltipla com regularização

Interpretação geométrica

É possível mostrar que os coeficientes estimados pelas regularizações *ridge* e LASSO resolvem os seguintes problemas de otimização, respectivamente:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j \cdot x_{ij} \right)^2 \text{ sujeito a } \sum_{j=1}^p \beta_j^2 \leq s$$

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j \cdot x_{ij} \right)^2 \text{ sujeito a } \sum_{j=1}^p |\beta_j| \leq s$$



As estimativas dos coeficientes são dadas pelo primeiro ponto em que uma curva de contorno do erro quadrático (ellipse) toca a região de restrições (losango para o LASSO e círculo para o *ridge*).

## Regularização

Empiricamente, observa-se que a regularização LASSO obtém sucesso em muitos problemas e seleciona no máximo  $\min(n, p)$  variáveis,  $n$  é a quantidade de registros e  $p$  a quantidade de variáveis preditoras. Isso é especialmente desvantajoso em problemas em que  $p \gg n$ , quando muitas variáveis relevantes podem ser descartadas, levando à instabilidade no desempenho.

Por outro lado, quando  $n > p$  em problemas com variáveis correlacionadas, frequentemente a *ridge* obtém melhores resultados que a LASSO, distribuindo os valores de coeficientes de forma mais homogênea, promovendo uma generalização mais suave e estável.

A *Elastic net* procura superar as dificuldades da LASSO combinando sua atuação com a da *ridge*:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{f(\lambda_1, \lambda_2, \beta)\}$$

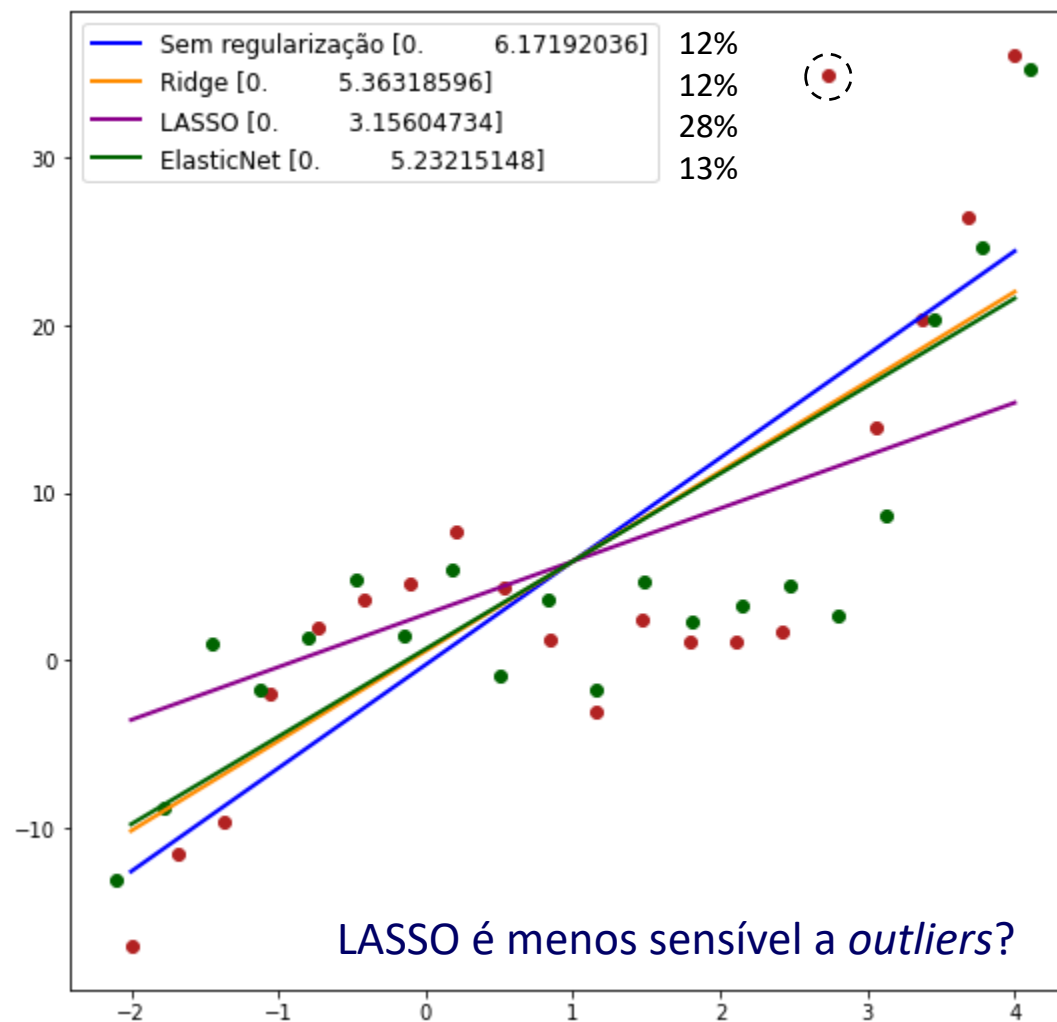
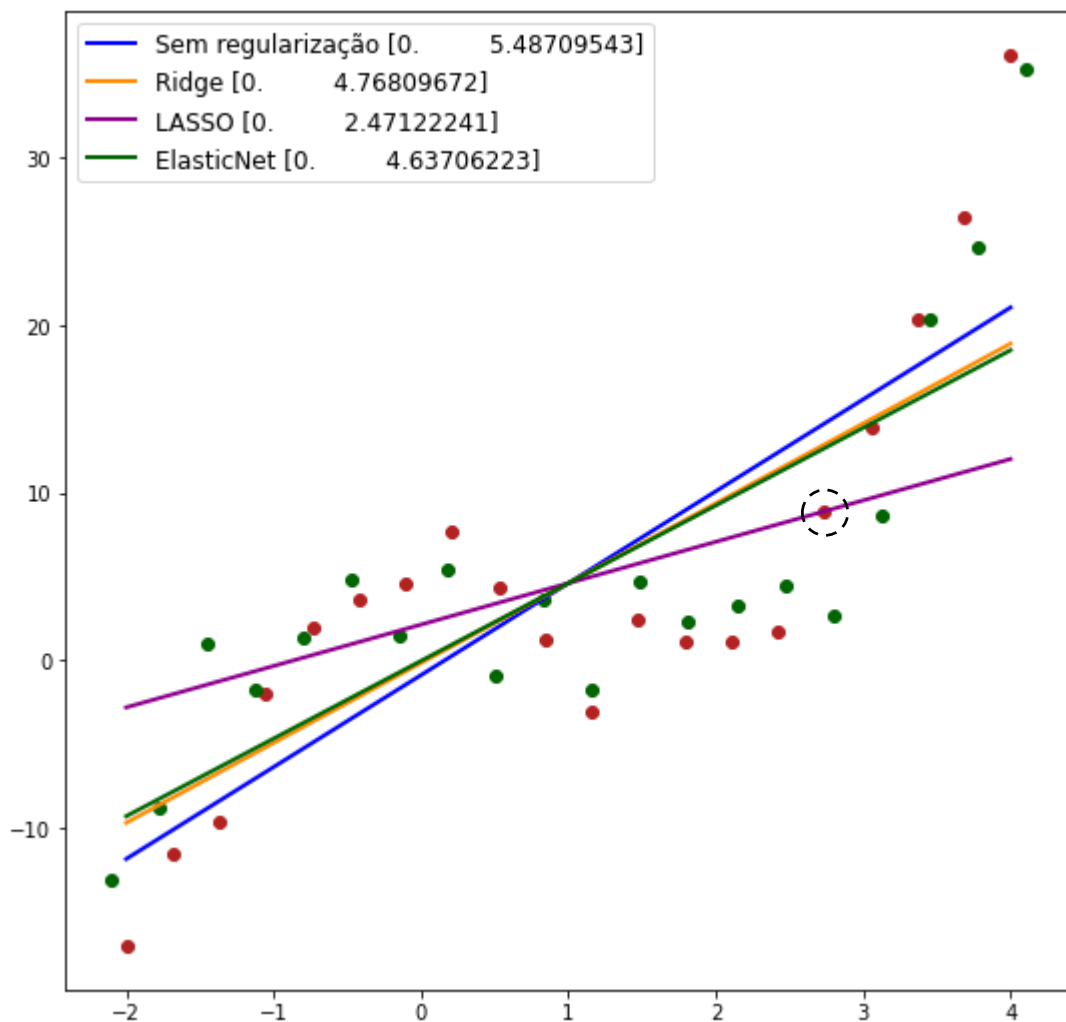
**L2/Tikhonov:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum \text{parâmetros}^2)$  (*ridge, de cume, de crista*)

**L1/LASSO:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda \cdot (\sum |\text{parâmetros}|)$  (*Least Absolute Shrinkage and Selection Operator*)

**Elastic net:**  $EMQ_R = E((y_r - \hat{y})^2) + \lambda_1 \cdot (\sum \text{parâmetros}^2) + \lambda_2 \cdot (\sum |\text{parâmetros}|)$  (rede elástica)

## Regularização

### Interpretação geométrica



## Regressão linear

Como lidar com variáveis discretas?

Se os possíveis valores são ordenados, basta mapeá-los no conjunto dos naturais.

Se as variáveis não são ordenadas e dicotômicas, pode-se:

$$x'_i = \begin{cases} 1, & \text{se } x_i = k \\ 0, & \text{se } x_i \neq k \end{cases}$$

sendo  $k$  um dos possíveis valores da variável original

Considerando a regressão simples:

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x'_i = \begin{cases} \widehat{\beta}_0 + \widehat{\beta}_1, & \text{se } x'_i = 1 \\ \widehat{\beta}_0, & \text{se } x'_i = 0 \end{cases}$$

$\widehat{\beta}_0$  é a média de  $\hat{y}$  quando  $x_i \neq k$ ;

$\widehat{\beta}_0 + \widehat{\beta}_1$  é a média de  $\hat{y}$  quando  $x_i = k$

$\widehat{\beta}_1$  é a média das diferenças entre os dois tipos

$$x'_i = \begin{cases} 1, & \text{se } x_i = k \\ -1, & \text{se } x_i \neq k \end{cases}$$

sendo  $k$  um dos possíveis valores da variável original

Considerando a regressão simples:

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x'_i = \begin{cases} \widehat{\beta}_0 + \widehat{\beta}_1, & \text{se } x'_i = 1 \\ \widehat{\beta}_0 - \widehat{\beta}_1, & \text{se } x'_i = -1 \end{cases}$$

$\widehat{\beta}_0 - \widehat{\beta}_1$  é a média de  $\hat{y}$  quando  $x_i \neq k$ ;

$\widehat{\beta}_0 + \widehat{\beta}_1$  é a média de  $\hat{y}$  quando  $x_i = k$  ; e

$\widehat{\beta}_0$  é a média das diferenças entre os dois tipos

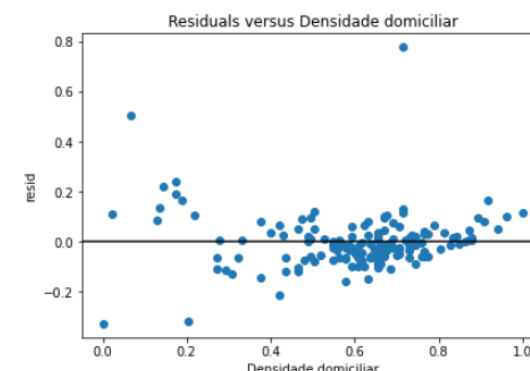
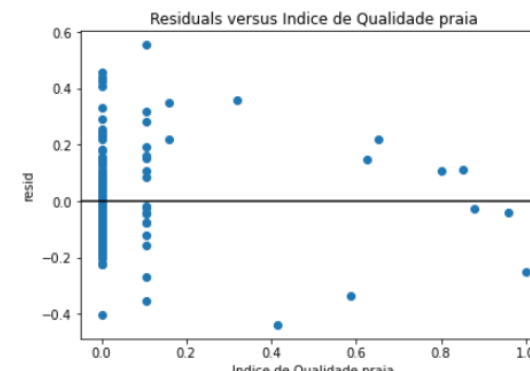
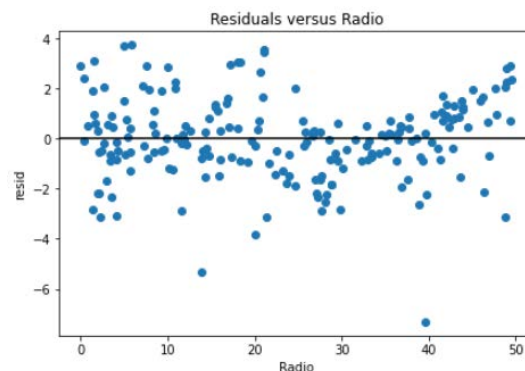
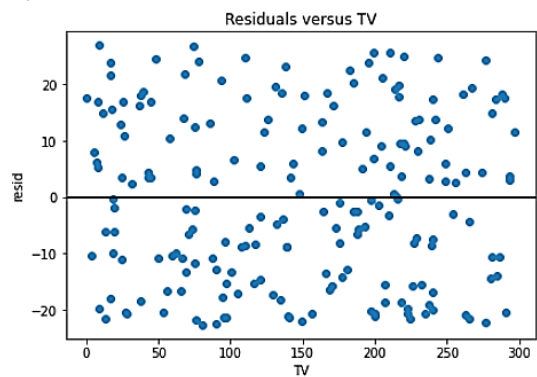
## Regressão linear

Possíveis problemas na modelagem linear:

- não linearidade dos dados, isto é, entre variáveis preditoras e variável de resposta, ou heterocedasticidade dos erros

Um modelo linear dificilmente capturará a forma real da relação entre variáveis preditoras e predita. Como consequências prováveis, os resíduos terão um padrão de comportamento, o que viola a suposição de homocedasticidade, e o ajuste será de má qualidade. O uso do gráfico de dispersão de resíduos (saídas preditas x resíduos) ajuda a identificar o problema: se houver alguma tendência de comportamento dos resíduos, há indícios de heterocedasticidade.

**Possível solução:** transformar os preditores com funções não lineares tais como  $\exp x$ ,  $\log x$ ,  $\sqrt{x}$  e  $x^2$ , usar pesos nos dados de modo a reduzir a influência de quem mais varia, dentre outros.





## Regressão linear

Possíveis problemas na modelagem linear:

- presença de *outliers*

Os *outliers* afetam as medidas de dispersão e de acurácia do modelo.

**Possível solução:** excluí-los, pois provavelmente nada significativo em relação ao ajuste será observado.

- colinearidade de preditores

Além de prejudicar a identificação do efeito de cada preditora envolvida, traz incertezas às estimativas dos coeficientes e pode causar a exclusão de alguma variável indevidamente.

**Possível solução:** excluir variáveis a partir da matriz de correlações.

- multicolinearidade de preditores

É possível que variáveis não correlacionadas apresentem correlações em grupo.

**Possível solução:** calcular o fator de inflação de variância (VIF) para identificar variáveis candidatas à exclusão; centralizar variáveis; usar PCA para combinar as variáveis envolvidas; usar regularização.

## Regressão linear

Possíveis problemas na modelagem linear:

- multicolinearidade de preditores

Se uma variável  $x_j$  está fortemente correlacionada com outras variáveis do modelo, haverá dificuldade em determinar o efeito dela na predição. Como consequência, a variância estimada do parâmetro  $\hat{\beta}_j$  aumenta (infla). Assim, o VIF mede o quanto a variância estimada de um parâmetro  $\hat{\beta}_j$  aumenta devido à colinearidade com os demais preditores do modelo.

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_{j,\sim j}^2}$$

em que  $R_{j,\sim j}^2$  é o coeficiente de determinação da regressão de  $x_j$  a partir de todas as outras variáveis regressoras do problema.

$VIF = 1$ : sem multicolinearidade  
 $5 < VIF \leq 10$ : multicolinearidade alta

$1 \leq VIF \leq 5$ : multicolinearidade moderada  
 $VIF > 10$ : multicolinearidade altíssima

## Regressão linear

Possíveis problemas na modelagem linear:

### correlação entre erros

A suposição da modelagem linear de que os erros (ou resíduos) são independentes entre si é violada quando se tem a autocorrelação dos erros: o erro de uma observação é influenciado pelo erro de outra.

Se os resíduos dos registros são correlacionados, então estão comprometidas as decisões em que foram usados os erros padrão dos parâmetros, que tenderão a subestimar os erros reais e a superestimar a confiança no modelo.

Analisar autocorrelação apenas faz sentido em dados sequenciais, sejam eles temporais, espaciais ou outra sequência lógica.

O uso do gráfico de dispersão de resíduos ajuda a identificar o problema. Pode-se aplicar o teste de Durbin-Watson, para autocorrelações sequenciais diretas, ou o teste de Breusch-Godfrey, mais geral e capaz de detectar autocorrelação com defasagens mais altas.

**Possível solução:** não utilizar o método de mínimos quadrados tradicional para encontrar os parâmetros e escolher outro capaz de lidar com dados sequenciais, tais como as modelagens ARIMA.

## Regressão logística

Modelagem estatística que utiliza uma função logística para estimar a probabilidade de pertinência de pontos em uma classe.

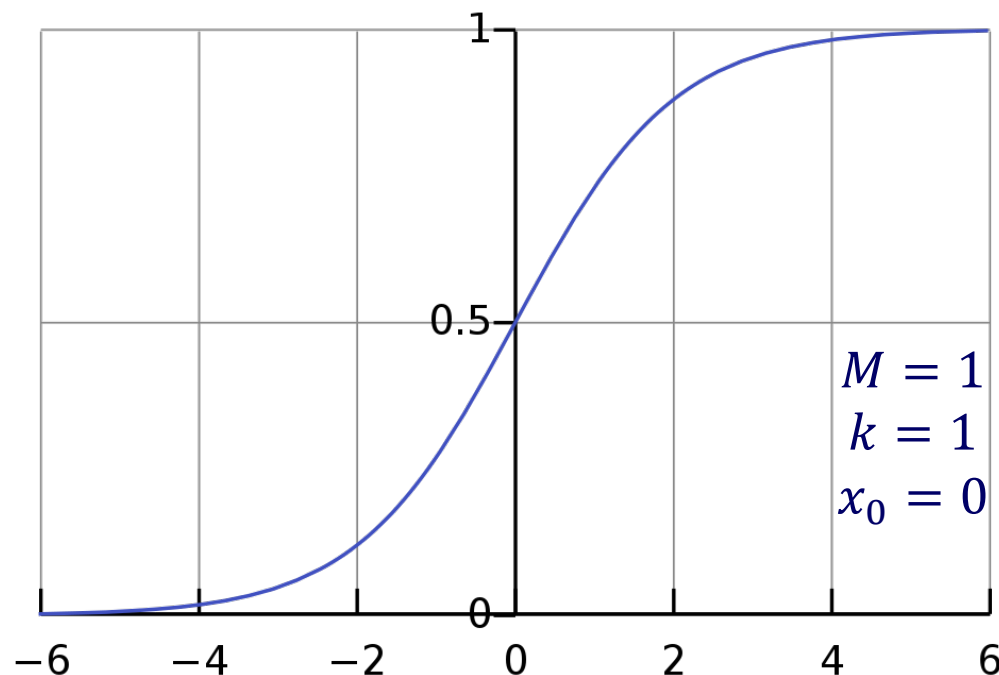
$$f(x) = \frac{M}{1 + e^{-k(x-x_0)}}$$

$M$ : valor máximo

$k$ : inclinação

$x_0$ : centro da sigmoide

$$p(y = 1|x) = f(x)$$



$$f(x) = \frac{1}{1 + e^{-x}}$$

O valor máximo da sigmoide é 1 (código da classe alvo). Os demais parâmetros da função são estimados por máxima verossimilhança a partir dos dados de treinamento.

## Regressão logística

$$f(x) = \frac{1}{1 + e^{-k(x-x_0)}} = \frac{e^{k(x-x_0)}}{e^{k(x-x_0)} + 1} \Rightarrow (e^{k(x-x_0)} + 1) \cdot f(x) = e^{k(x-x_0)}$$

$$\Rightarrow f(x) = e^{k(x-x_0)} - e^{k(x-x_0)} \cdot f(x) \Rightarrow f(x) = e^{k(x-x_0)} \cdot (1 - f(x))$$

$$\Rightarrow \boxed{\frac{f(x)}{1 - f(x)}} = e^{k(x-x_0)} \Rightarrow \boxed{\ln \frac{f(x)}{1 - f(x)}} = \boxed{k(x - x_0)}$$

log odds  
logit

chance/odds  $\in [0, +\infty[$       função linear

$$\ln \frac{f(x)}{1 - f(x)} = k(x - x_0) = \beta_0 + \beta_1 x$$

Assim, estimar os parâmetros  $\beta_0$  e  $\beta_1$  da função linear permitirá encontrar a melhor função logística de acordo com os dados. Os melhores parâmetros são os que maximizam a função de verossimilhança dos dados de treinamento.

## Regressão logística

Os melhores parâmetros são os que maximizam a verossimilhança dos dados de treinamento:

$$\mathcal{L}(\beta_0, \beta_1) = \prod_{x,y=1} p(x) \cdot \prod_{x,y=0} (1 - p(x))$$

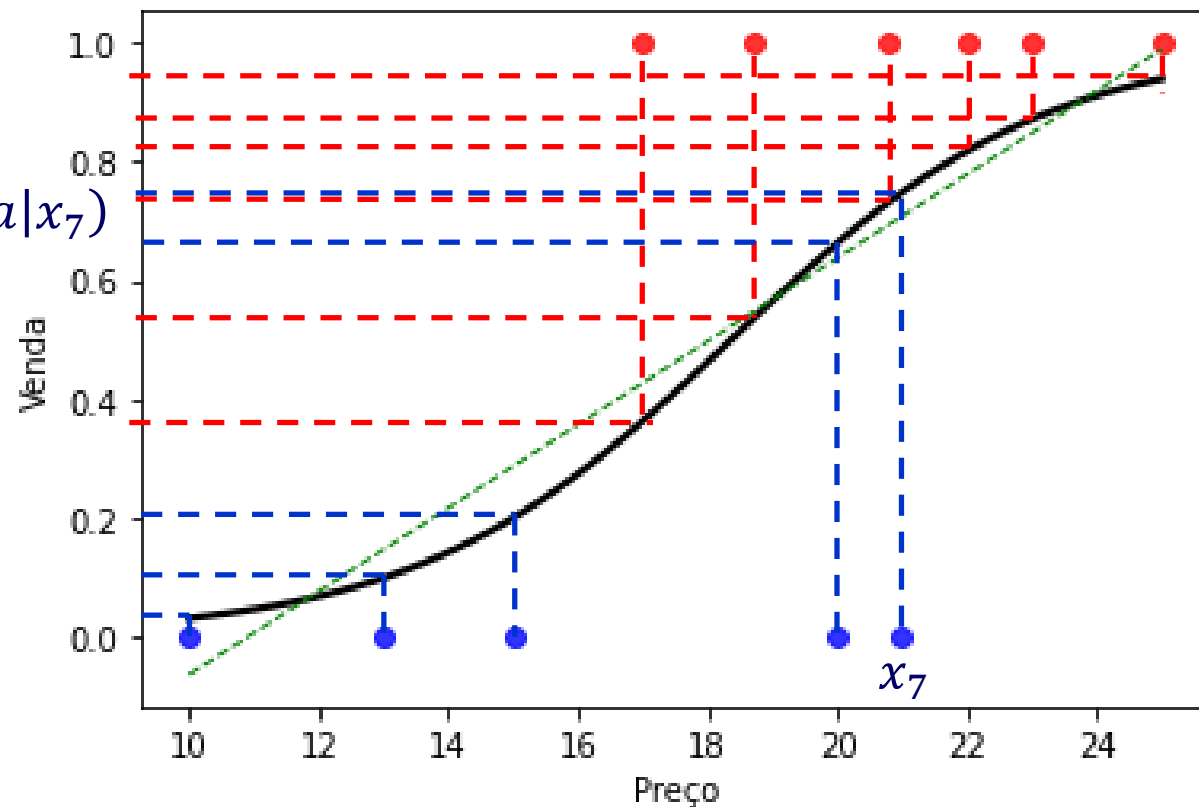
$x$  é registro da classe alvo

$x$  não é registro da classe alvo

Como os produtórios acima resultam em valores muito pequenos, na prática usa-se o logaritmo da função de verossimilhança:

$$\ln \mathcal{L}(\beta_0, \beta_1) = \sum_{x,y=1} \ln p(x) + \sum_{x,y=0} \ln(1 - p(x))$$

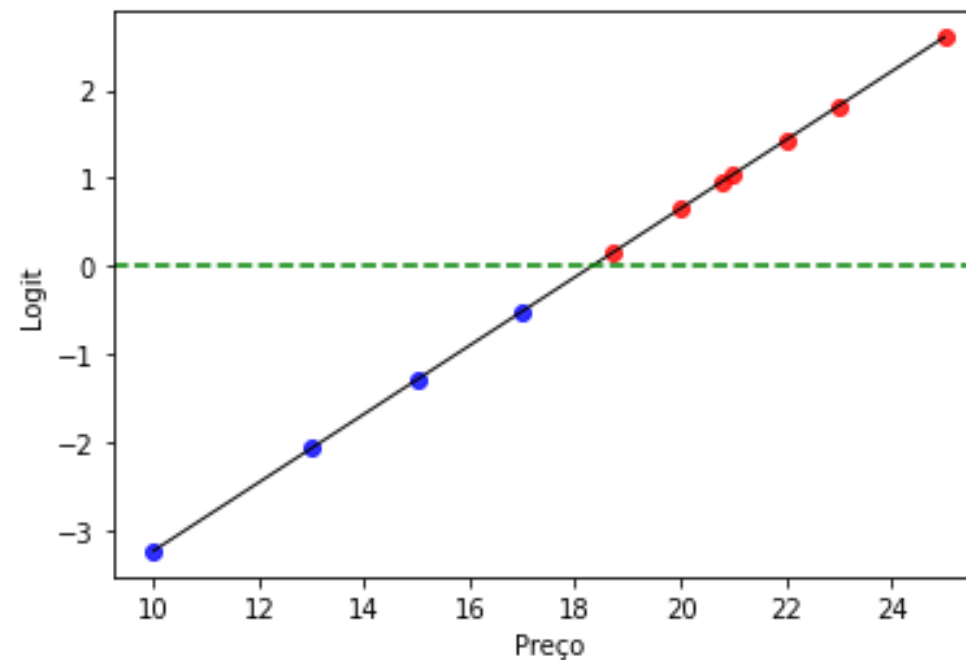
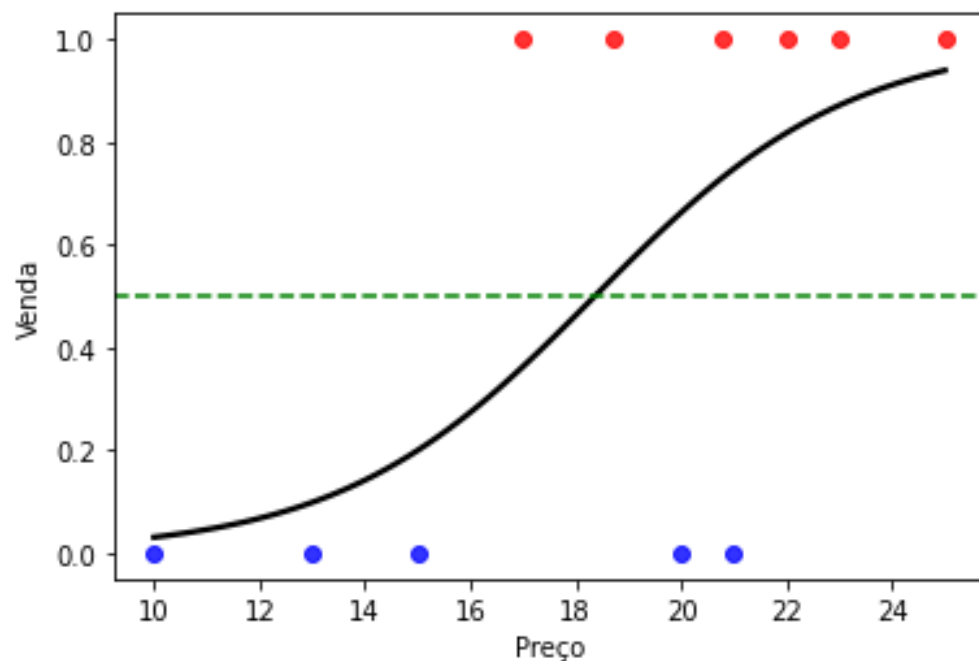
$p(Venda|x_7)$



$$Prob.Vendas(x) = \frac{1}{1 + e^{7,132 - 0,389x}}$$

## Regressão logística

Durante o aprendizado, o problema é levado do espaço original, em que o domínio da variável resposta (probabilidade de pertinência à classe alvo) é  $[0,1]$ , para outro espaço em que a variável resposta (logit) tem domínio  $[-\infty, +\infty]$ .



Assim, os parâmetros da função logística do espaço original serão dados pelo hiperplano produzido no espaço de busca do problema de otimização  $\max \ln \mathcal{L}(\beta_0, \beta_1)$ .



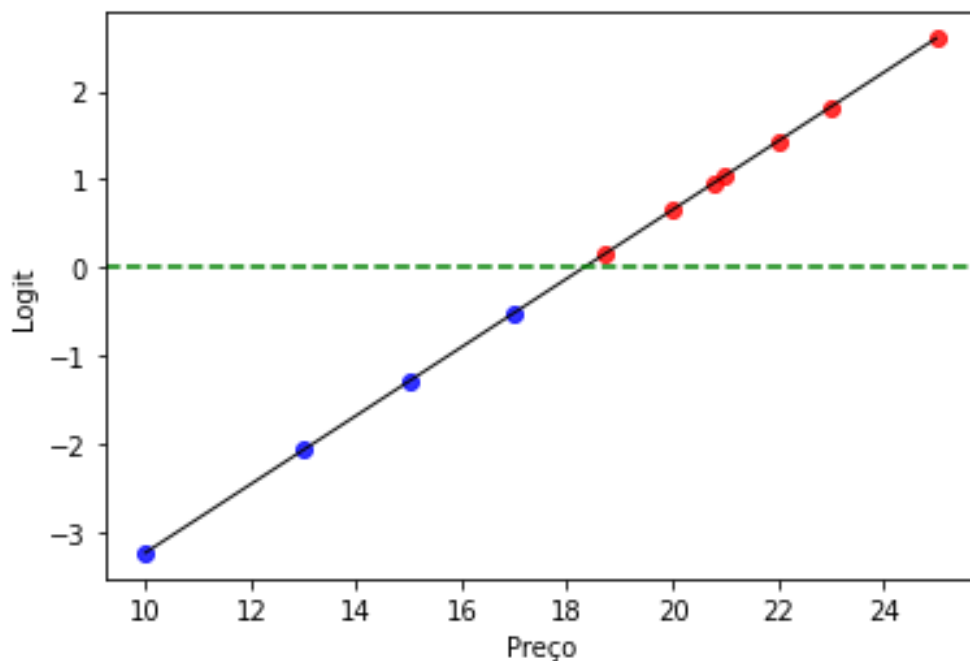
## Regressão logística

Uso do modelo obtido:

$\beta_0 + \beta_1 x \geq 0$ : é da classe alvo

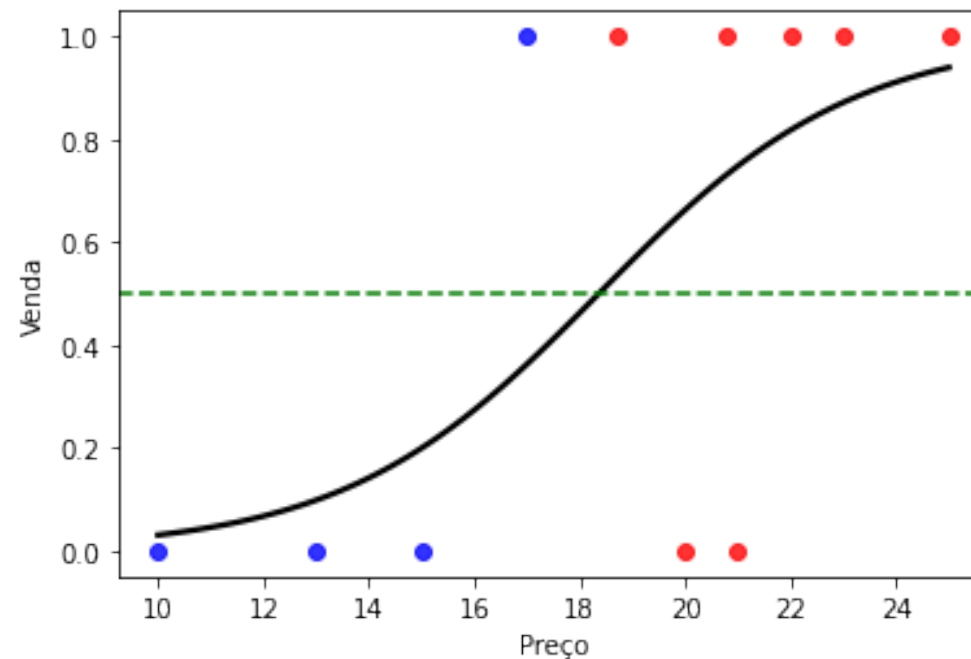
$\beta_0 + \beta_1 x < 0$ : não é da classe alvo

Exemplo:



$-7,132 + 0,389x \geq 0$ : é da classe alvo

$-7,132 + 0,389x < 0$ : não é da classe alvo



$$Prob.Vendas(x) = \frac{1}{1 + e^{7,132 - 0,389x}}$$

## Regressão logística

Passos:

1. dados  $\widehat{\beta}_0, \dots, \widehat{\beta}_n$  candidatos, obtém-se um modelo linear múltiplo  $\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 x_1 + \dots + \widehat{\beta}_n x_n$ ;
2. transformam-se as respostas deste modelo para os dados de treino em suas probabilidades/verossimilhanças:  $p(x_{i1}, \dots, x_{in}) = e^{\widehat{y}_i} / (e^{\widehat{y}_i} + 1)$
3. calcula-se o logaritmo da verossimilhança  $\mathcal{L}(\widehat{\beta}_0, \dots, \widehat{\beta}_n)$  da curva;
4. ajustam-se os parâmetros  $\widehat{\beta}_0, \dots, \widehat{\beta}_n$  segundo alguma técnica de otimização numérica a partir das (estimativas das) derivadas parciais de  $\ln \mathcal{L}(\widehat{\beta}_0, \dots, \widehat{\beta}_n)$  ou de outras técnicas;
5. os passos acima são repetidos até que  $\mathcal{L}(\widehat{\beta}_0, \dots, \widehat{\beta}_n)$  atinja seu máximo valor ou algum critério de parada.

Ao final, tem-se o modelo linear mais ajustado no espaço transformado que corresponderá à melhor função logística no espaço original. Como não é possível calcular o coeficiente de determinação sem resíduos, pode-se calcular um falso  $R^2$  para ser interpretado como um indicador de desempenho relativo para a seleção de modelos, assim como AUC e tantos outros.

$$R^2_{McFadden} = 1 - \frac{\ln \mathcal{L}(\beta_0, \dots, \beta_n)}{\ln \mathcal{L}(\beta_0)}$$