

Introdução a Ciências de Dados

Aula 6: Classificação

Francisco A. Rodrigues
ICMC/USP
francisco@icmc.usp.br



Aula 6: Classificação

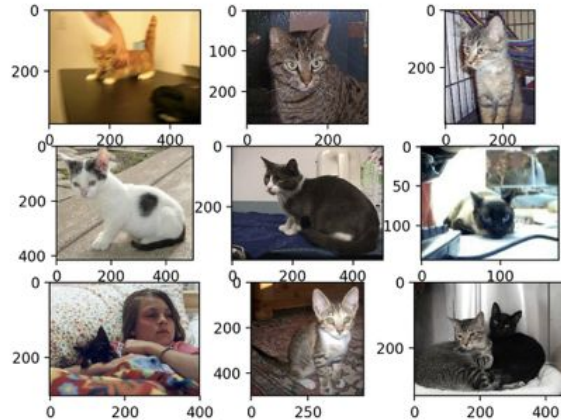
- **k-vizinhos mais próximos**
- **Regressão Logística**
- **Naive Bayes**

k-vizinhos mais próximos

K-vizinhos

- Uma maneira simples de definimos a classificação de objetos é através de distância entre eles:

Classe: gatos

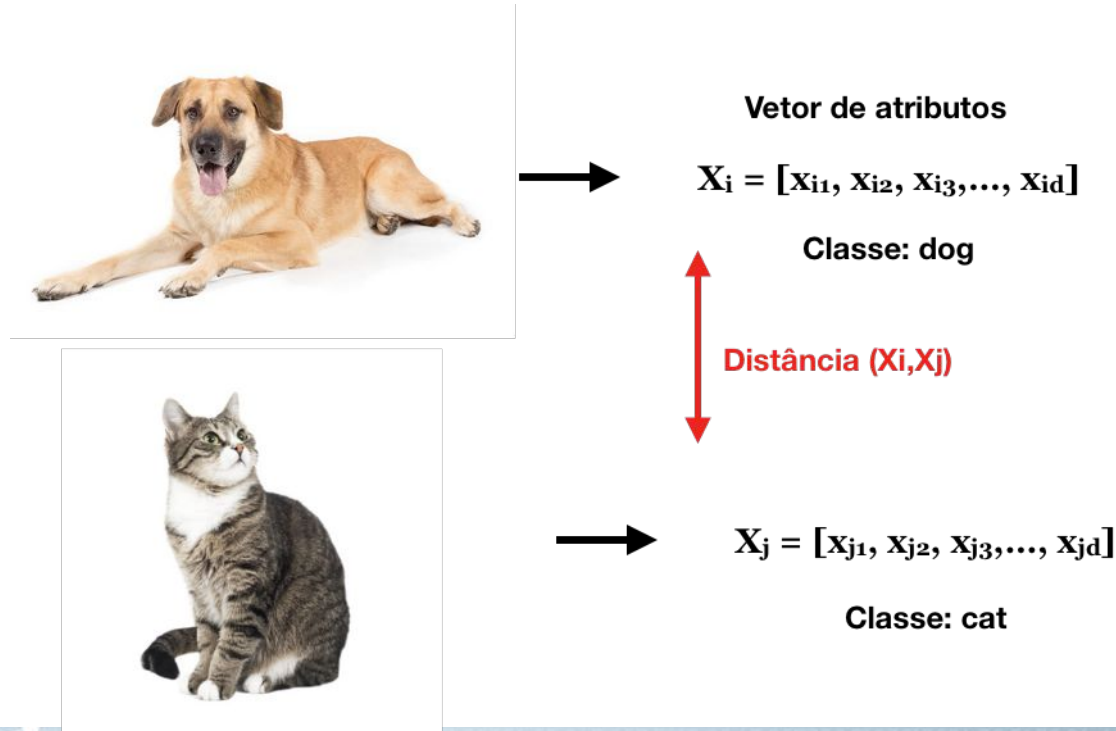


Classe: Cachorros



K-vizinhos

- Como definir a distância entre objetos?



K-vizinhos

- Precisamos definir uma medida de proximidade.
 - **Medida de similaridade:** $d(X_i, X_i)$ é máxima.
 - Exemplo: Número de amigos compartilhados em uma rede social.
 - **Medida de dissimilaridade:** $d(X_i, X_i) = 0$.
 - Exemplo: distância entre cidades (distância Euclidiana).

K-vizinhos

Medida de dissimilaridade:

- $d(p, q) \geq 0$ para todo p e q , e $d(p, q) = 0$ se, e somente se, $p = q$,
- $d(p, q) = d(q, p)$ para todo p e q ,
- $d(p, r) \leq d(p, q) + d(q, r)$ para todo p, q , e r , onde $d(p, q)$ é a distância de dissimilaridade entre os pontos (objetos) p e q .

Medida de similaridade:

- $s(p, q) = 1$ (ou máximo de similaridade) se $p = q$,
- $s(p, q) = s(q, p)$ para todo p e q , onde $s(p, q)$ é a similaridade entre os objetos p e q .

K-vizinhos

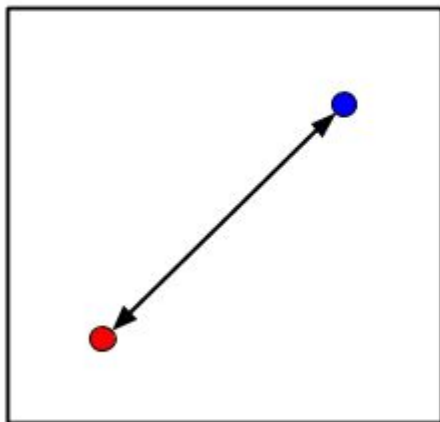
Métricas de distância:

- **Euclidiana** $D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$ $[0, \infty)$
- **Minkowski** $D(X, Y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$ $[0, \infty)$
- **Cosseno** $D(X, Y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n (x_i)^2} \sqrt{\sum_{i=1}^n (y_i)^2}}$ $[0, 1]$
- **Pearson** $D(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$ $[-1, 1]$

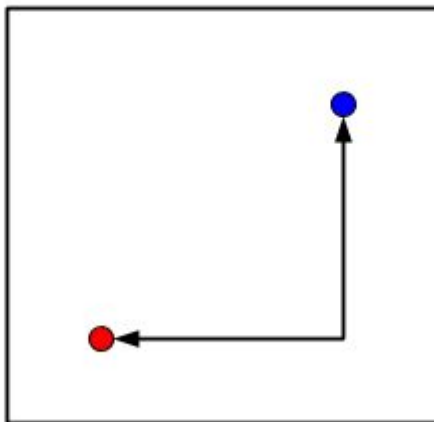
K-vizinhos

Métricas de distância:

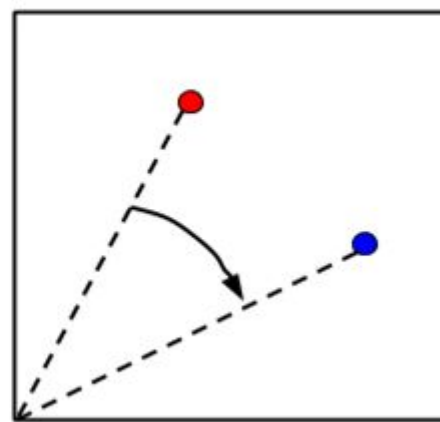
Euclidean



Manhattan



Cosine Similarity



K-vizinhos

Métricas de distância:

- Dados nominais

Similaridade

$$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$$

Dissimilaridade

$$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$$

- Dados ordinais

Similaridade

$$s = 1 - \frac{\|p - q\|}{n - 1}$$

Dissimilaridade

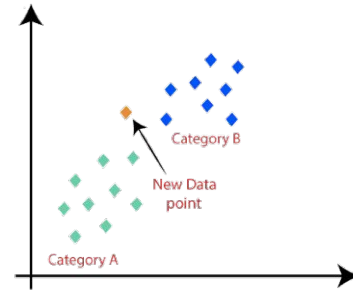
$$d = \frac{\|p - q\|}{n - 1}$$

K-vizinhos

Algoritmo:

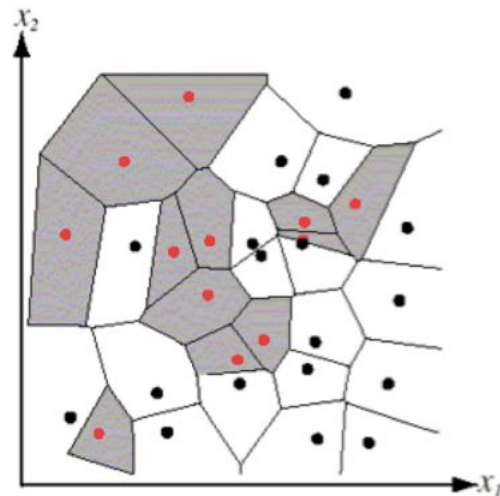
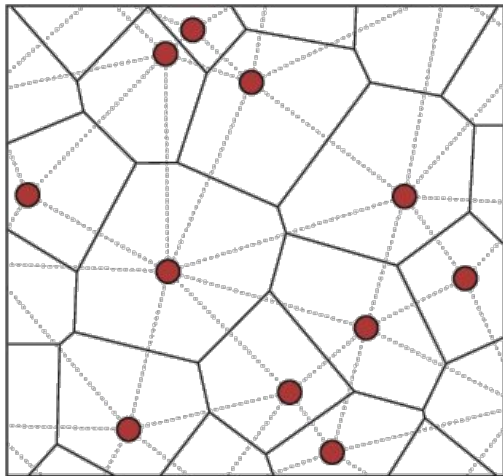
1. Identifique os k-vizinhos mais próximos do vetor de atributos **X** que se quer classificar.
2. Determine o número de vizinhos em cada classe.
3. Classifique **X** com pertencente à classe que resultou em um maior número de vizinhos (a moda entre o número de classes).

$$p(y = j | \mathbf{x}_*) = \frac{1}{k} \sum_{i \in R_*} \mathbb{I}\{y_i = j\}$$



K-vizinhos

Regiões de separação formam telhas de Voronoi.



K-vizinhos

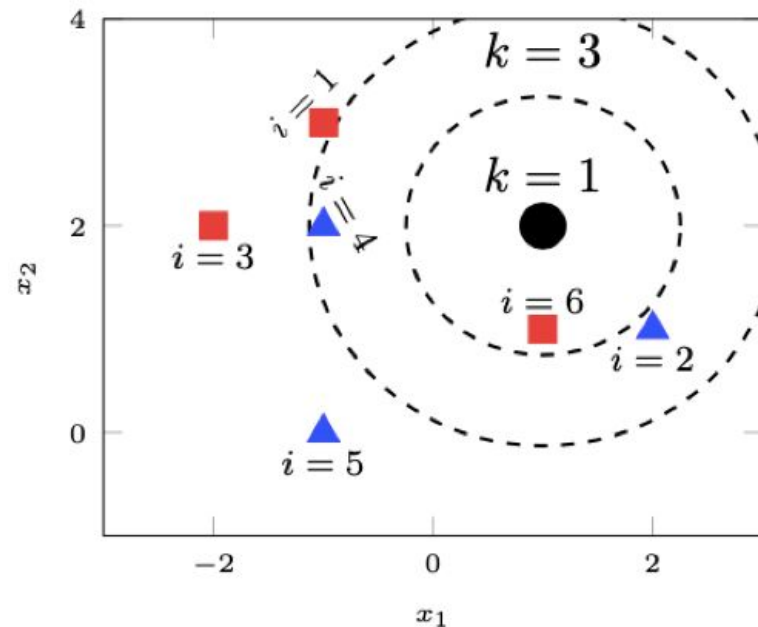
Exemplo:

Observação: $\mathbf{x}_* = [1 \ 2]^\top$

i	x_1	x_2	y
1	-1	3	Red
2	2	1	Blue
3	-2	2	Red
4	-1	2	Blue
5	-1	0	Blue
6	1	1	Red

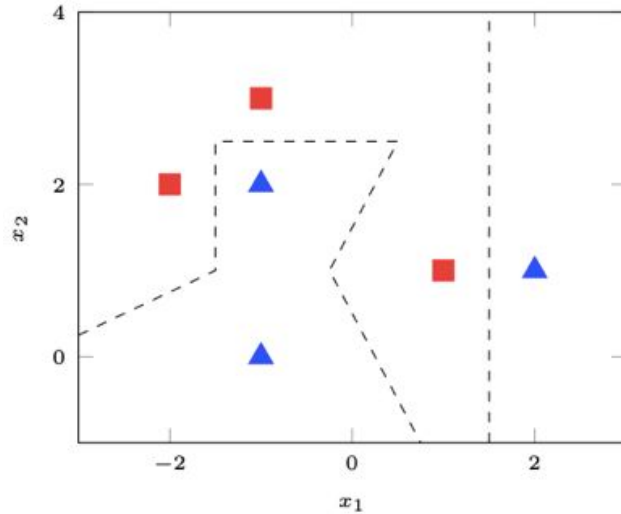
Distâncias

i	$\ \mathbf{x}_i - \mathbf{x}_*\ $	y_i
6	$\sqrt{1}$	Red
2	$\sqrt{2}$	Blue
4	$\sqrt{4}$	Blue
1	$\sqrt{5}$	Red
5	$\sqrt{8}$	Blue
3	$\sqrt{9}$	Red

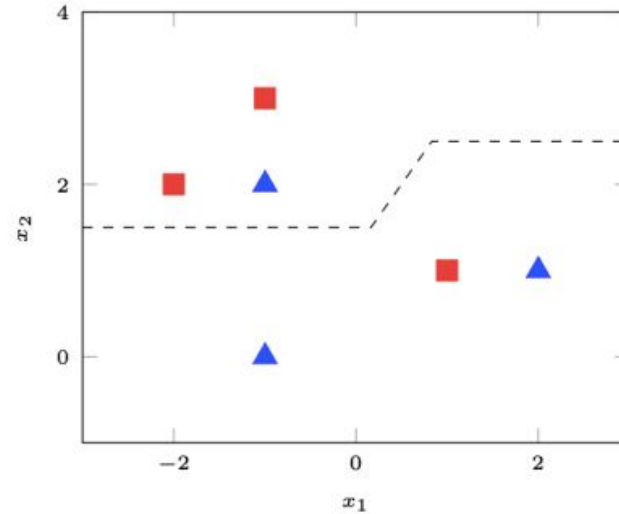


K-vizinhos

A região de decisão pode depender de k .



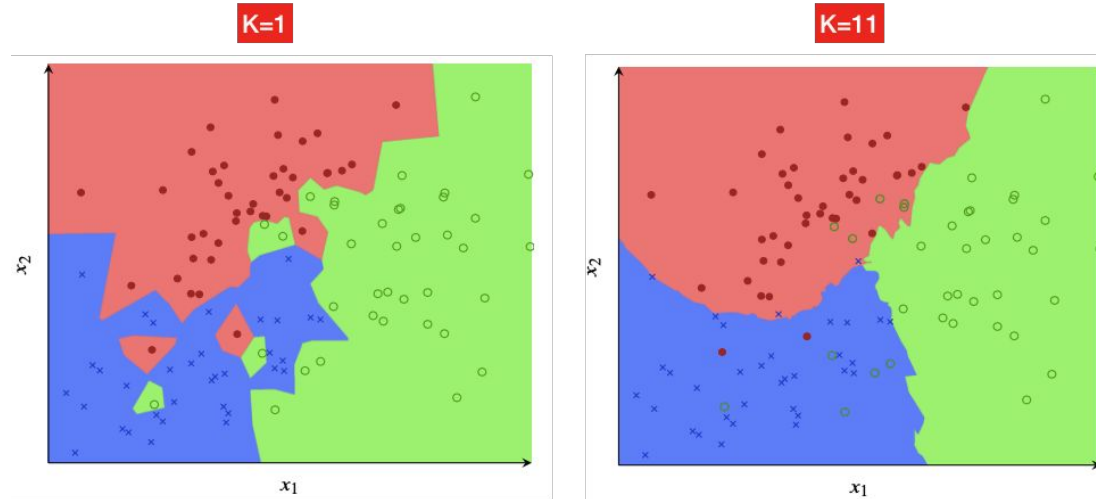
(a) $k = 1$



(b) $k = 3$

K-vizinhos

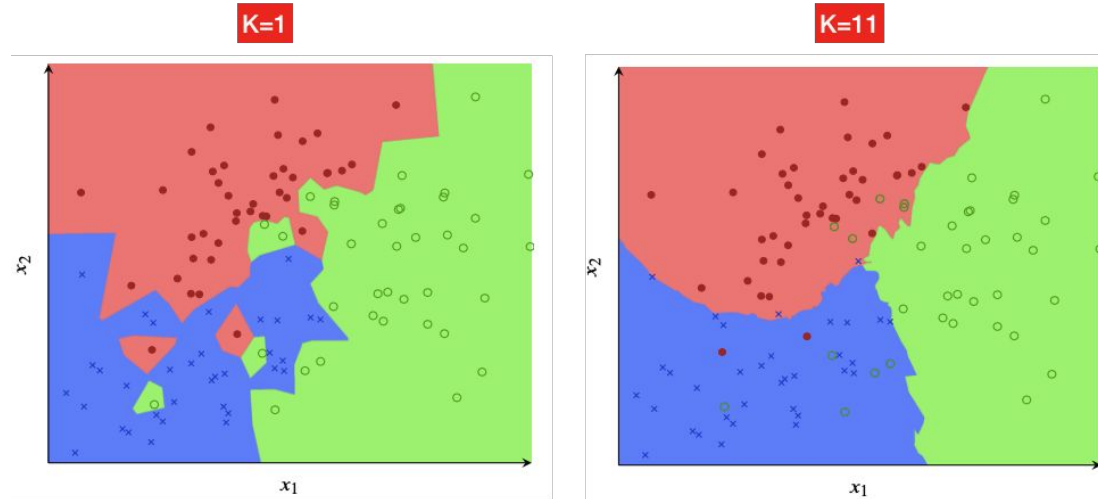
- Para $k=1$, vemos que a região de decisão se “**ajusta muito**” ao dados, ocorrendo overfitting.
- Para $k = 11$, vemos que região é bastante **suave**, o que sugere underfitting.



K-vizinhos

Qual o melhor valor de k?

A melhor maneira de encontrar o melhor valor de k é usar validação cruzada e uma medida para avaliar o resultado da classificação, como a acurácia.



K-vizinhos

Propriedades:

- O algoritmo não “aprende” um modelo, apenas memoriza objetos de treinamento
- Adia computação para a fase de classificação
- O algoritmo pode ser entendido como não paramétrico, dependendo apenas do número de vizinhos k .
- É um classificador não-linear, não sendo restrito a regiões de separação lineares.
- Como geralmente a distância Euclidiana é considerada, é necessário normalizar ou padronizar os dados.
- Dado que o conjunto de treinamento seja relativamente grande, pode-se provar que o erro cometido na classificação é no máximo duas vezes maior do que o classificador Bayesiano, que é ótimo.

Regressão Logística

Regressão Logística

- Vimos que modelos de regressão linear são dados por:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d = \beta^T \mathbf{x}$$

onde $\mathbf{x}^T = [1, x_1, x_2, \dots, x_d]$

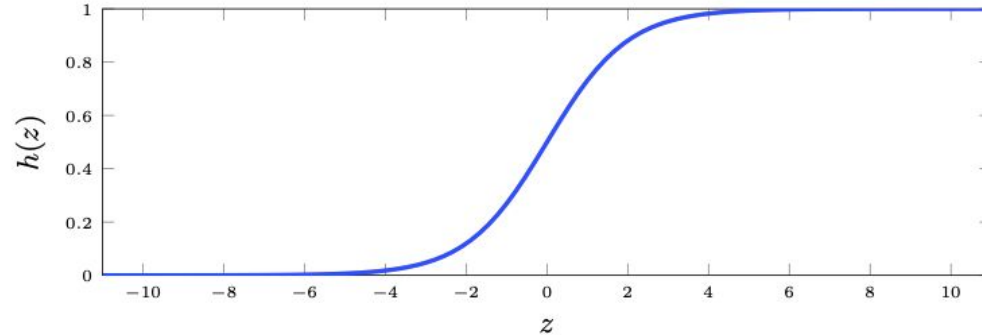
- Se considerarmos a saída \mathbf{Y} como um valor inteiro, podemos usar o modelo de regressão para realizarmos a classificação de dados.
- Vamos definir as probabilidade para o caso de duas classes:

$$\mathbf{p}(\mathbf{y}=1|\mathbf{x}) \text{ e } \mathbf{p}(\mathbf{y}=0|\mathbf{x})$$

Regressão Logística

- Vamos considerar a função logística:

$$h(z) = \frac{e^z}{1 + e^z}$$



- Essa função retorna valores no intervalo $[0,1]$.

Regressão Logística

- Usando a função logística, temos:

$$p(y = 1 | \mathbf{x}) = \frac{e^{\beta^T \mathbf{x}}}{1 + e^{\beta^T \mathbf{x}}} \quad p(y = 0 | \mathbf{x}) = 1 - p(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{\beta^T \mathbf{x}}}$$

- O aprendizado se resume a estimar o vetor de parâmetros β .
- Esse método é chamado regressão logística.

Regressão Logística

Para estimar os parâmetros do modelo, usamos o conjunto de treinamento $T = \{(x_i, y_i)\}_{i=1}^N$.

Usando estimação por máxima verossimilhança:

$$\hat{\beta} = \operatorname{argmax}_{\beta} p(\mathbf{y} | \mathbf{X}; \beta) .$$

Assim, a máxima verossimilhança:

$$\ell(\beta) = p(\mathbf{y} | \mathbf{X}; \beta) = \prod_{i=1}^N p(y_i | x_i; \beta) = \prod_{i: y_i=1} p(y = 1 | \mathbf{x}_i; \beta) \prod_{i: y_i=0} p(y = 0 | \mathbf{x}_i; \beta)$$

$$\ell(\beta) = \prod_{i: y_i=1} \frac{e^{\beta^T \mathbf{x}_i}}{1 + e^{\beta^T \mathbf{x}_i}} \prod_{i: y_i=0} \frac{1}{1 + e^{\beta^T \mathbf{x}_i}}$$

Regressão Logística

- Para estimar os parâmetros, precisamos otimizar essa função em relação à beta.

$$\ell(\beta) = \prod_{i:y_i=1} \frac{e^{\beta^T \mathbf{x}_i}}{1 + e^{\beta^T \mathbf{x}_i}} \prod_{i:y_i=0} \frac{1}{1 + e^{\beta^T \mathbf{x}_i}}$$

- Para facilitar os cálculos, podemos a função logaritmo, que é uma função monotônica e logo, o máximo é o mesmo que o da verossimilhança:

$$\log \ell(\beta) = \sum_{i:y_i=1} \left(\beta^T \mathbf{x}_i - \log(1 + e^{\beta^T \mathbf{x}_i}) \right) - \sum_{i:y_i=0} \log(1 + e^{\beta^T \mathbf{x}_i})$$

Regressão Logística

- Calculando a derivada e igualando a zero:

$$\nabla_{\beta} \log \ell(\beta) = \sum_{i=1}^n \mathbf{x}_i \left(y_i - \frac{e^{\beta^T \mathbf{x}_i}}{1 + e^{\beta^T \mathbf{x}_i}} \right) = 0$$

- Como essa equação é não-linear, devemos resolvê-la numericamente, por exemplo, usando o método de Newton–Raphson.
- Com isso, determinamos os parâmetros do modelo e realizamos a classificação em uma das duas classes.

Regressão Logística

- A superfície de decisão pode ser calculada usando:

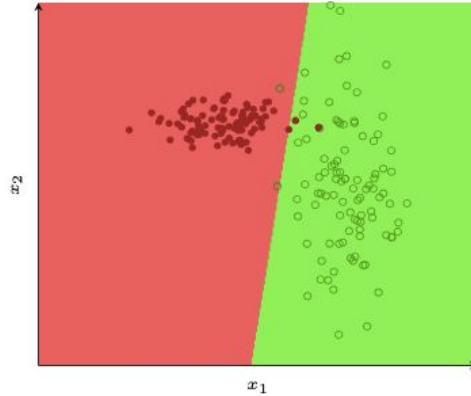
$$p(y = 1 | \mathbf{x}) = p(y = 0 | \mathbf{x})$$

$$\frac{e^{\beta^T \mathbf{x}}}{1 + e^{\beta^T \mathbf{x}}} = \frac{1}{1 + e^{\beta^T \mathbf{x}}}$$

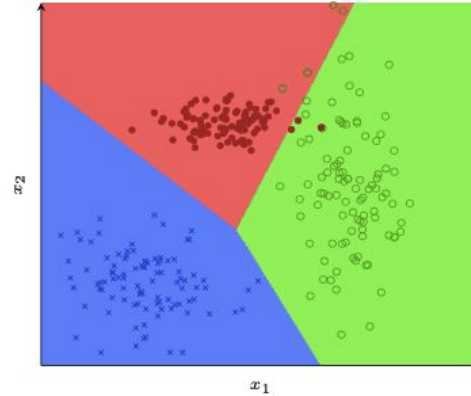
- Ou seja, basta resolvermos: $\beta^T \mathbf{x} = 0$
- A solução são hiperplanos. Logo, a superfície de separação são hiperplanos (lineares).

Regressão Logística

Duas classes



Três classes



- Para mais de duas classes, usamos one-hot-encoding e repetimos o mesmo procedimento.

http://www.it.uu.se/edu/course/homepage/sml/literature/lecture_notes.pdf

Naive Bayes

Naive Bayes

Teoria da decisão Bayesiana:

- Dada **M** classes $\omega_1, \omega_2, \dots, \omega_M$ e um padrão desconhecido **x**, determinar a probabilidade condicional $p(\omega_i|\mathbf{x})$ do padrão pertencer a cada classe **i**.

$$P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})}$$

OBJ

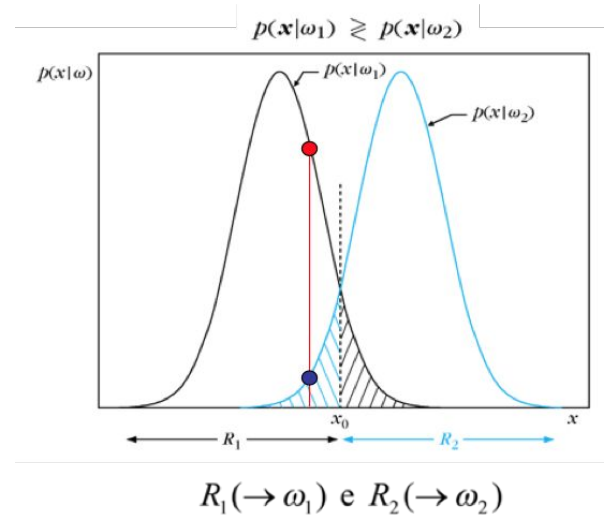
- Classificar de acordo com a classe mais provável.

$$\hat{y} = \underset{i \in \{1, \dots, M\}}{\operatorname{argmax}} \quad p(\omega_i | \mathbf{x})$$

Naive Bayes

Teoria da decisão Bayesiana:

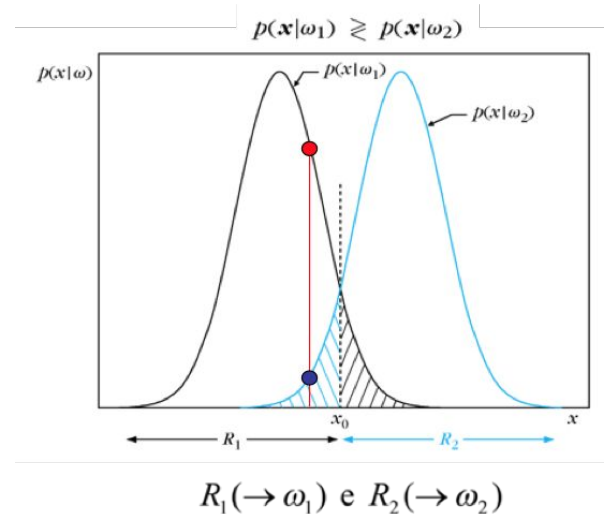
- Se a probabilidade condicional for conhecida para cada classe, o erro obtido é o menor possível (classificador ótimo).



Naive Bayes

Teoria da decisão Bayesiana:

- **Problema:** na maioria das vezes, não sabemos a distribuição de probabilidade conjunta e sua estimativa é bastante complicada.



Naive Bayes

Teoria da decisão Bayesiana:

- **Solução simples:** assumir que as variáveis que descreve os atributos são independentes.

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i) P(\omega_i)}{p(\mathbf{x})}$$

$$p(\mathbf{x} | \omega_i) = \prod_{j=1}^d p(x_j | \omega_i) \quad i = 1, 2, \dots, M$$

Naive Bayes

- Lembrem-se: Se duas variáveis X e Y aleatórias são independentes, então :

$$P(X|Y) = P(X) \longrightarrow P(X|Y) = \frac{P(X, Y)}{P(Y)} = P(X)$$

$$P(X, Y) = P(X)P(Y)$$

- Para variáveis aleatórias independentes, a distribuição de probabilidade conjunto é igual ao produto de suas marginais.

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^N P(X_i) = P(X_1)P(X_2)\dots P(X_n)$$

Naive Bayes

- Dada **M** classes $\omega_1, \omega_2, \dots, \omega_M$ e um padrão desconhecido \mathbf{x} , determinar a probabilidade condicional $p(\omega_i|\mathbf{x})$ do padrão pertencer a cada classe **i**.

[OBJ]

$$P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})}$$

Classificar de acordo com a classe mais provável.

$$\hat{y} = \operatorname{argmax}_{i \in \{1, \dots, M\}} p(\omega_i | \mathbf{x})$$

$$\hat{y} = \operatorname{argmax}_{i \in \{1, \dots, M\}} p(\omega_i) \prod_{j=1}^d p(x_j | \omega_i)$$

Naive Bayes

Exemplo:

Classifique:

**x = (Arrepios = Sim,
Nariz escorrendo = Não,
Dor de cabeça = média,
febre = Não)**

Arrepios	Nariz escorrendo	Dor de cabeça	febre	Gripe
Sim	Não	Média	Sim	Não
Sim	Sim	Não	Não	Sim
Sim	Não	Forte	Sim	Sim
Não	Sim	Média	Sim	Sim
Não	Não	Não	Não	Não
Não	Sim	Forte	Sim	Sim
Não	Sim	Forte	Não	Não
Sim	Sim	Média	Sim	Sim

Naive Bayes

Exemplo:

	Arrepios	
	Gripe = Sim	Gripe = Não
Arrepios = Sim	3/4	1/4
Arrepios = Não	2/4	2/4

	Nariz	
	Gripe = Sim	Gripe = Não
Escorrendo = Sim	4/5	1/5
Escorrendo = Não	1/3	2/3

	Cabeça	
	Gripe = Sim	Gripe = Não
Dor = Forte	2/3	1/3
Dor = Média	2/3	1/3
Dor = Não	1/2	1/2

	Febre	
	Gripe = Sim	Gripe = Não
Sim	4/5	1/5
Não	1/3	2/3

Podemos calcular todas as probabilidades conjuntas:

$$P(\text{Arrepios}=\text{Sim}, \text{Gripe} = \text{Sim}) = 3/4$$

$$P(\text{Arrepios}=\text{Não}, \text{Gripe} = \text{Sim}) = 2/4$$

$$P(\text{Dor}=\text{Forte}, \text{Gripe} = \text{Sim}) = 2/3$$

Naive Bayes

Exemplo:

- Vamos classificar os dados:

**x = (Arrepios = Sim, Nariz escorrendo = Não,
Dor de cabeça = média, febre = Sim)**

- Usando probabilidade condicional: **$P(A|B)=P(B|A)P(A)$**

$$P(\text{Gripe}=\text{Sim} | x) = P(x|\text{Gripe}=\text{Sim})P(\text{Gripe} = \text{Sim})$$

Naive Bayes

Exemplo:

- A probabilidade do paciente estar com gripe dados os sintomas:
 $x = (\text{Arrepios} = \text{Sim}, \text{Nariz escorrendo} = \text{Não}, \text{Dor de cabeça} = \text{média}, \text{febre} = \text{Não})$
- $P(\text{Gripe}=\text{Sim}) = 5/8$
- $P(\text{Arrepios}=\text{Sim}|\text{Gripe}=\text{Sim}) = 3/5$
- $P(\text{Nariz escorrendo}=\text{Não}|\text{Gripe}=\text{Sim}) = 1/5$
- $P(\text{Dor de cabeça}=\text{Média}|\text{Gripe}=\text{Sim}) = 2/5$
- $P(\text{febre}=\text{Não}|\text{Gripe}=\text{Sim}) = 1/5$
- $P(\text{Gripe}=\text{Sim}|x) = P(x|\text{Gripe}=\text{Sim})P(\text{Gripe}=\text{Sim}) =$
- $P(\text{Arrepios}=\text{Sim}|\text{Gripe}=\text{Sim}) * P(\text{Escorrendo}=\text{Não}|\text{Gripe}=\text{Sim}) * P(\text{Dor}=\text{média}|\text{Gripe}=\text{Sim}) * P(\text{Febre}=\text{Sim}|\text{Gripe}=\text{Sim}) * P(\text{Gripe}=\text{Sim}) =$
 $(3/5) * (1/5) * (2/5) * (1/5) * (5/8) = 0,006$

$$P(\text{Gripe}=\text{Sim}|x) = 0,006$$

Naive Bayes

Exemplo:

- A probabilidade do paciente não estar com gripe dados os sintomas:
 $x = (\text{Arrepios} = \text{Sim}, \text{Nariz escorrendo} = \text{Não}, \text{Dor de cabeça} = \text{média}, \text{febre} = \text{Não})$
 - $P(\text{Gripe}=\text{Não}) = 3/8$
 - $P(\text{Arrepios}=\text{Sim}|\text{Gripe}=\text{Não}) = 1/4$
 - $P(\text{Nariz escorrendo}=\text{Não}|\text{Gripe}=\text{Não}) = 2/3$
 - $P(\text{Dor de cabeça}=\text{Média}|\text{Gripe}=\text{Não}) = 1/3$
 - $P(\text{febre}=\text{Não}|\text{Gripe}=\text{Não}) = 2/3$
- $P(\text{Gripe}=\text{Não}|x) = P(x|\text{Gripe}=\text{Não})P(\text{Gripe}=\text{Não}) =$
 $P(\text{Arrepios}=\text{Sim}|\text{Gripe}=\text{Não}) * P(\text{Nariz escorrendo}=\text{Não}|\text{Gripe}=\text{Não}) * P(\text{Dor de cabeça}=\text{Média}|\text{Gripe}=\text{Não}) * P(\text{febre}=\text{Não}|\text{Gripe}=\text{Não}) * P(\text{Gripe}=\text{Não}) =$
 $(1/4) * (2/3) * (1/3) * (2/3) * (3/8) = 0,013$

$$P(\text{Gripe}=\text{Não}|x) = 0,013$$

Naive Bayes

Exemplo:

- Assim, temos:

$$P(\text{Gripe}=\text{Sim}|\mathbf{x}) = 0,006$$

$$P(\text{Gripe}=\text{Não}|\mathbf{x}) = 0,013$$

- Classificando de acordo com a classe mais provável.

$$\hat{y} = \underset{i \in \{1, \dots, M\}}{\operatorname{argmax}} \quad p(\omega_i | \mathbf{x})$$

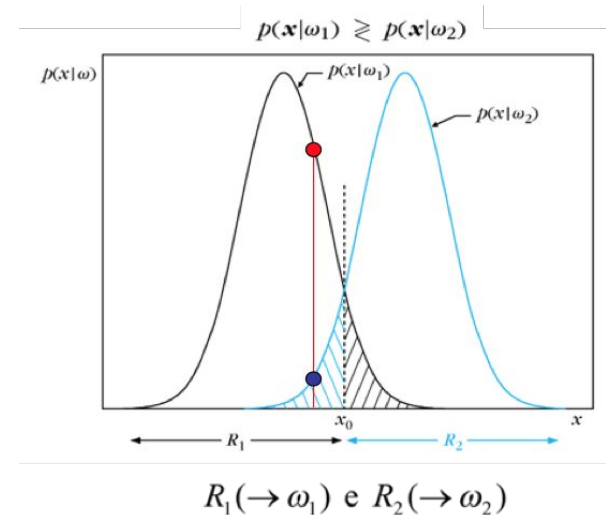
- Concluimos que o paciente não está gripado.**

\mathbf{x} = (Arrepios = Sim, Nariz escorrendo = Não, Dor de cabeça = média, febre = Sim)

Naive Bayes

Caso os atributos sejam contínuos, podemos assumir uma distribuição de probabilidade (geralmente Normal) e realizar a classificação maximizando a verossimilhança.

$$p(x_j | \omega_i) = \frac{1}{\sqrt{2\pi\sigma_{\omega_i}}} \exp\left(-\frac{(x_j - \mu_{\omega_i})^2}{2\sigma_{\omega_i}^2}\right)$$



Naive Bayes

Algoritmo (para atributos contínuos):

- Calcule a média e variância de cada atributo para cada classe.
- Calcule a verossimilhança para cada atributo dentro de cada classe.

[OBJ]

$$p(x_j | \omega_i) = \frac{1}{\sqrt{2\pi\sigma_{\omega_i}}} \exp\left(-\frac{(x_j - \mu_{\omega_i})^2}{2\sigma_{\omega_i}}\right)$$

- Assuma independência e calcule a distribuição conjunta.

[OBJ]

$$p(\mathbf{x} | \omega_i) = \prod_{j=1}^d p(x_j | \omega_i) \quad i = 1, 2, \dots, M$$

- Classifique de acordo com a classe mais provável.

[OBJ]

$$\omega_m = \arg \max_{\omega_i} \prod_{j=1}^d p(x_j | \omega_i), \quad i = 1, 2, \dots, M$$

Naive Bayes

Propriedades:

- Apesar da limitação em assumir independência dos atributos, o classificador Naive Bayes é robusto e apresenta boa performance para muitos dados reais.
- Todas as probabilidades exigidas podem ser calculadas dos dados de treinamento em uma passagem.
- Construção do modelo é bastante eficiente.
- Fácil de estender para incremental.
- Robusto a ruídos e atributos irrelevantes.

Sumário

- **k-vizinhos mais próximos,**
- **Regressão Logística,**
- **Naive Bayes.**

Leitura adicional

- Lindholm et al., Supervised Machine Learning, 2019.
http://www.it.uu.se/edu/course/homepage/sml/literature/lecture_notes.pdf
- James et al., Introduction to statistical learning with applications in R, 2014.
<https://www.ime.unicamp.br/~dias/Intoduction%20to%20Statistical%20Learning.pdf>