

# Classificadores Bayesianos

## Mineração de Dados

Universidade Federal do ABC

Introdução

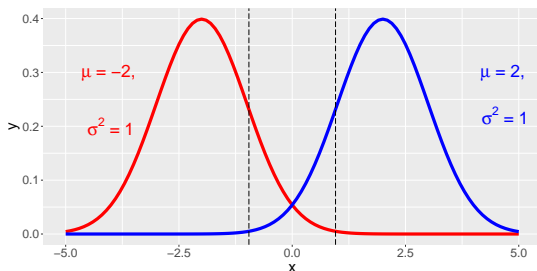
Discriminantes

Náive Bayes

# INTRODUÇÃO

# INTRODUÇÃO

- Podemos pensar o problema que discutimos do ponto de vista probabilístico
  - Temos duas classes que correspondem a funções de densidade de probabilidade distintas ( $p(x|c_1)$  e  $p(x|c_2)$ )



# INTRODUÇÃO

- ▶ Lembrando as aulas de probabilidade
  - ▶ [Lei da Probabilidade Total]

$$p(x) = \sum_{c \in \mathcal{C}} p(x, c) = \sum_{c \in \mathcal{C}} p(x|c)p(c)$$

- ▶ [Teorema de Bayes]

$$p(c|x) = \frac{p(c, x)}{p(x)} = \frac{p(x|c)p(c)}{p(x)} = \frac{p(x|c)p(c)}{\sum_{c \in \mathcal{C}} p(x|c)p(c)}$$

- ▶  $\text{posterior} = \text{likelihood} \times \text{prior} \times \text{evidence}^{-1}$

# INTRODUÇÃO

- Em um problema de classificação queremos encontrar

$$\arg \max_{c \in \mathcal{C}} p(c|x)$$

$$\arg \max_{c \in \mathcal{C}} \frac{p(x|c)p(c)}{\sum_{c \in \mathcal{C}} p(x|c)p(c)}$$

# INTRODUÇÃO

- ▶ Queremos encontrar uma regra de decisão que minimiza o erro

$$p(\text{erro}|x) = \begin{cases} p(c_1|x) & \text{se escolhermos } c_2 \\ p(c_2|x) & \text{se escolhermos } c_1 \end{cases}$$

- ▶ Logo chegamos a regra de decisão de Bayes
  - ▶ Classe  $c_1$  se  $p(c_1|x) > p(c_2|x)$ , e  $c_2$  caso contrário
  - ▶ Classe  $c_1$  se  $p(x|c_1)p(c_1) > p(x|c_2)p(c_2)$ , e  $c_2$  caso contrário
- ▶ Sob essa regra, temos  $p(\text{erro}|x) = \min(p(c_1|x), p(c_2|x))$
- ▶ Classificador *ótimo* de Bayes
  - ▶ Como  $p(x|c_1)$  é definido?

# INTRODUÇÃO

- ▶ Função discriminante
  - ▶  $g_i(x) > g_j(x) \forall j \neq i \rightarrow c_i$
- ▶ Formas equivalentes do ponto de vista de classificação
  - ▶  $g_i(x) = p(c_i|x) = \frac{p(x|c_i)p(c_i)}{\sum_{c_j \in C} p(x|c_j)p(c_j)}$
  - ▶  $g_i(x) = p(x|c_i)p(c_i)$
  - ▶  $g_i(x) = \log(p(x|c_i)) + \log(p(c_i))$
- ▶ Caso de duas classes (dicotomizador)
  - ▶ Uma única função discriminante
  - ▶  $g(x) = g_1(x) - g_2(x)$
  - ▶  $g(x) = \log\left(\frac{p(x|c_1)}{p(x|c_2)}\right) + \log\left(\frac{p(c_1)}{p(c_2)}\right)$



Introdução

Discriminantes

Naïve Bayes

# FUNÇÕES DISCRIMINANTES - CASO GAUSSIANA MULTIVARIADA

- ▶ Abordagem paramétrica
- ▶ Vamos assumir que temos duas classes
  - ▶  $p(\mathbf{x}|c_1) = \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$
  - ▶  $p(\mathbf{x}|c_2) = \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2)$
  - ▶  $p(c_i)$  estimada pela proporção de objetos da classe  $i$
  - ▶  $\boldsymbol{\mu}_i$  é a média ( $D$ -dimensões) da classe  $i$
  - ▶  $\Sigma_i$  é a matriz ( $D, D$ ) de covariância da classe  $i$

$$\mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i) = \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right]$$

$$g_i(\mathbf{x}) = -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_i|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln(p(c_i))$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

- ▶  $\sigma^2 = \frac{(n_1-1)\sigma_1^2 + \dots + (n_C-1)\sigma_C^2}{N-C}$  [ $\sigma_c^2$  é a variância estimada nos dados da classe  $c$ , *desvio padrão combinado*]
- ▶  $|\Sigma_i| = \sigma^{2d}$  [Determinante de  $\text{diag}(a_1, \dots, a_n) = a_1 \cdot \dots \cdot a_n$ ]
- ▶  $\Sigma^{-1} = \frac{1}{\sigma^2} \mathbf{I}$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

- Desconsiderando os termos iguais para todas as classes

$$g_i(\mathbf{x}) = -\frac{(\mathbf{x} - \boldsymbol{\mu}_i)^T(\mathbf{x} - \boldsymbol{\mu}_i)}{2\sigma^2} + \ln(p(c_i))$$

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} - 2\boldsymbol{\mu}_i^T \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i) + \ln(p(c_i))$$

$[\mathbf{x}^T \mathbf{x}]$  é uma constante aditiva]

$$g_i(\mathbf{x}) = \left(\frac{\boldsymbol{\mu}_i}{\sigma^2}\right)^T \mathbf{x} + \left(-\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \ln(p(c_i))\right)$$

- temos uma função discriminante **linear**

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

$$g_i(\mathbf{x}) = \left(\frac{\boldsymbol{\mu}_i}{\sigma^2}\right)^T \mathbf{x} + \left(-\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \ln(p(c_i))\right)$$

- ▶ Desconsiderando os termos iguais para todas as classes
  - ▶ se  $p(c_i)$  igual para todas as classes temos o **classificador de distância mínima**
    - ▶ novo objeto classificado para a classe com média mais próxima

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

## ► Exemplo

$$\text{► } \boldsymbol{\mu}_1^T = [1 \ 2] \quad \boldsymbol{\mu}_2^T = [4 \ 6] \quad \boldsymbol{\mu}_3^T = [-2 \ 4]$$

$$\text{► } p(c_1) = p(c_2) = \frac{1}{4} \quad p(c_3) = \frac{1}{2}$$

$$\text{► } \Sigma = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$

$$g_i(\mathbf{x}) = \left( \frac{\boldsymbol{\mu}_i}{\sigma^2} \right)^T \mathbf{x} + \left( -\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \ln(p(c_i)) \right)$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

$$g_1(\mathbf{x}) = \left( \frac{\begin{bmatrix} 1 & 2 \end{bmatrix}}{3} \right) [x_1 \ x_2]^T + \left( -\frac{5}{6} - 1, 38 \right)$$

$$g_2(\mathbf{x}) = \left( \frac{\begin{bmatrix} 4 & 6 \end{bmatrix}}{3} \right) [x_1 \ x_2]^T + \left( -\frac{52}{6} - 1, 38 \right)$$

$$g_3(\mathbf{x}) = \left( \frac{\begin{bmatrix} -2 & 4 \end{bmatrix}}{3} \right) [x_1 \ x_2]^T + \left( -\frac{20}{6} - 0, 69 \right)$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

- ▶ Exemplo
  - ▶ Agora basta verificar casos em que  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \forall j \neq i$
  - ▶ Podemos verificar as **fronteiras de decisão**
    - ▶ Regiões definidas por  $g_i(\mathbf{x}) = g_j(\mathbf{x})$

$$g_1(\mathbf{x}) = g_2(\mathbf{x}) \rightarrow -x_1 - \frac{4}{3}x_2 = -\frac{47}{6}$$

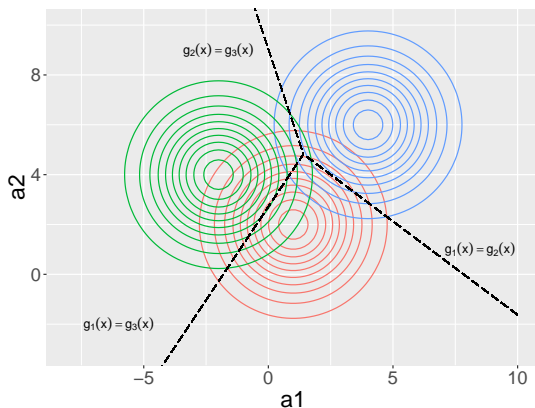
$$g_2(\mathbf{x}) = g_3(\mathbf{x}) \rightarrow 2x_1 + \frac{2}{3}x_2 = 6.02$$

$$g_1(\mathbf{x}) = g_3(\mathbf{x}) \rightarrow x_1 - \frac{2}{3}x_2 = -1.81$$



# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERESFÉRICAS

## ► Exemplo



# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

- ▶ Caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)
  - ▶  $\Sigma = \frac{(n_1-1)\Sigma_1 + \dots + (n_C-1)\Sigma_C}{N-C}$  [ $\Sigma_c$  é a matriz de covariância estimada nos dados da classe  $c$ , *matriz de covariância combinada*]
  - ▶  $|\Sigma_i|$  igual para todas as classes

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln(p(c_i))$$

- ▶ Se  $p(c_i)$  igual para todas as classes ainda temos um classificador de distância mínima
  - ▶ de acordo com a distância de **Mahalanobis**

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

- Caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln(p(c_i))$$

$$g_i(\mathbf{x}) = -\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2}\boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln(p(c_i))$$

$[\mathbf{x}^T \Sigma^{-1} \mathbf{x}$  constante aditiva]

$$g_i(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2}\boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln(p(c_i))$$

- Novamente temos um caso de discriminante linear
  - método conhecido como *análise de discriminante linear* (LDA)

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

- ▶ Exemplo caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)

- ▶  $\boldsymbol{\mu}_1^T = [1 \ 2] \quad \boldsymbol{\mu}_2^T = [-1 \ 5] \quad \boldsymbol{\mu}_3^T = [-2 \ 4]$

- ▶  $p(c_1) = p(c_2) = \frac{1}{4} \quad p(c_3) = \frac{1}{2}$

- ▶  $\Sigma = \begin{bmatrix} 1 & -1.5 \\ -1.5 & 4 \end{bmatrix}$

$$g_i(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln(p(c_i))$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

- ▶ Exemplo caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)
  - ▶ Vale a pena manipular a equação para encontrar as fronteiras de decisão

$$g_i(\mathbf{x}) = g_j(\mathbf{x})$$

$$\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln(p(c_i)) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j + \ln(p(c_j))$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

$$\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln(p(c_i)) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j + \ln(p(c_j))$$

$$\left( \boldsymbol{\mu}_i^T \Sigma^{-1} - \boldsymbol{\mu}_j^T \Sigma^{-1} \right) \mathbf{x} = \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i - \ln(p(c_i)) - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j + \ln(p(c_j))$$

$$\left( \boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T \right) \Sigma^{-1} \mathbf{x} = \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j + \ln \left( \frac{p(c_j)}{p(c_i)} \right)$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIAS I

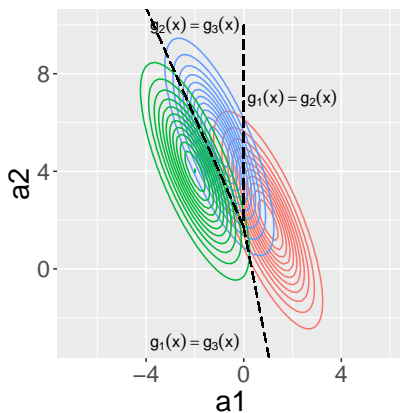
- ▶ Exemplo caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)

$$\left(\boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T\right) \Sigma^{-1} \mathbf{x} = \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j + \ln \left( \frac{p(c_j)}{p(c_i)} \right)$$

- ▶ Caso  $g_1(x) = g_2(x) \rightarrow x_1 = 0$
- ▶ Caso  $g_2(x) = g_3(x) \rightarrow 3,14x_1 + 1,4x_2 = 2,41$
- ▶ Caso  $g_1(x) = g_3(x) \rightarrow 5,14x_1 + 1,43x_2 = 2,41$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIS I

- Exemplo caso em que  $\Sigma_i = \Sigma$  (classes hiperelipsoidais, mesma forma)





# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIIS II

- ▶ Caso em que  $\Sigma_i$  são arbitrárias (classes hiperelipsoidais)
  - ▶ Não é possível remover muitos termos
  - ▶ Resulta em fronteiras **quadráticas**
    - ▶ método chamado de *análise de discriminante quadrática* (QDA)

$$g_i(\mathbf{x}) = -\frac{1}{2} \ln(|\Sigma_i|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln(p(c_i))$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \Sigma_i^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \Sigma_i^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^T \Sigma_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln(|\Sigma_i|) + \ln(p(c_i))$$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIIS II

- ▶ Exemplo caso em que  $\Sigma_i$  são arbitrárias (classes hiperelipsoidais)

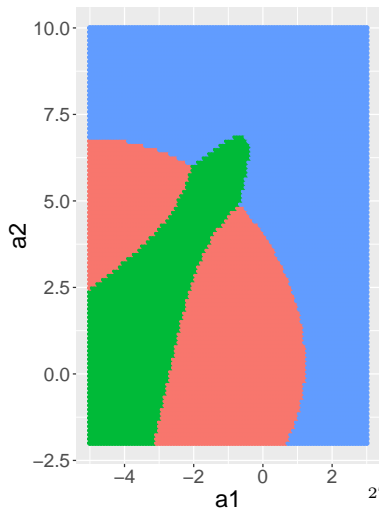
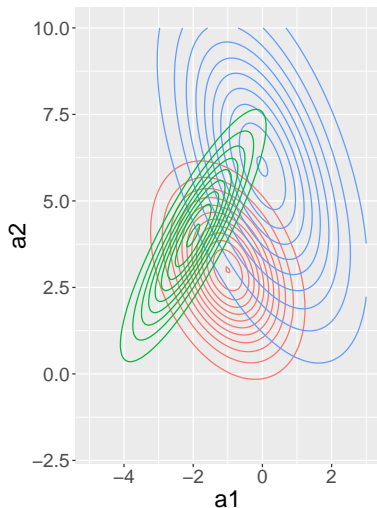
- ▶  $\mu_1^T = [-1 \ 3] \quad \mu_2^T = [0 \ 6] \quad \mu_3^T = [-2 \ 4]$

- ▶  $p(c_1) = p(c_2) = \frac{1}{4} \quad p(c_3) = \frac{1}{2}$

- ▶  $\Sigma_1 = \begin{bmatrix} 1 & -0,5 \\ -0,5 & 2 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 2 & -2 \\ -2 & 7 \end{bmatrix} \quad \Sigma_3 = \begin{bmatrix} 1 & 1,5 \\ 1,5 & 3 \end{bmatrix}$

# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIS II

- Exemplo caso em que  $\Sigma_i$  são arbitrárias



# CASO GAUSSIANA MULTIVARIADA - CLASSES HIPERELIPSOIDAIIS II

- ▶ Temos um classificador ótimo quando as premissas são válidas
  - ▶ Normalmente a premissa de normalidade não é válida
- ▶ Apesar disso o método apresenta bons resultados
  - ▶ Mesmo no caso restrito em que se assume matrizes de covariância iguais
- ▶ Em alta dimensionalidade o custo de inverter a matriz de covariância é alto —  $O(D^3)$ 
  - ▶ Como podemos melhorar?

Introdução

Discriminantes

Naïve Bayes

# NAÏVE BAYES

- ▶ De volta ao básico, Teorema de Bayes:

- ▶ [Teorema de Bayes]

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

- ▶ Ao invés de assumirmos que  $p(\mathbf{x}|c) \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , podemos assumir que os atributos são **condicionalmente independentes dado a classe**
- ▶ O modelo obtido a partir dessa premissa é chamado de *Naïve Bayes*

# NAÏVE BAYES

- ▶ Independência condicional
  - ▶  $p(a, b|c) = p(a|c)p(b|c)$
  - ▶ Note que não necessariamente  $a$  e  $b$  são independentes
    - ▶ Dois eventos são independentes se  $p(a, b) = p(a)p(b)$
  - ▶ Exemplo:
    - ▶ Temos duas moedas e uma delas é enviesada (sempre cara)
    - ▶ Pegamos uma moeda aleatoriamente e jogamos para cima 2x
    - ▶ Seja  $A$  o evento do primeiro lançamento ser cara
    - ▶ Seja  $B$  o evento do segundo lançamento ser cara
    - ▶ Seja  $C$  o evento de termos escolhido a moeda enviesada
    - ▶ Se não soubermos nada além de que  $A$  ocorreu, a probabilidade de  $B$  aumenta
    - ▶ Se soubermos que  $C$  ocorreu,  $A$  ocorrer não influencia a probabilidade de  $B$  ocorrer

# NAÏVE BAYES

- ▶ Independência condicional
  - ▶ No nosso contexto, os atributos são considerados independentes dado a classe
    - ▶ Premissa forte e normalmente inválida
    - ▶ Costuma funcionar bem



## NAÏVE BAYES

$$p(c|\mathbf{x}) = \frac{p(c, x_1, x_2, \dots, x_D)}{p(\mathbf{x})} = \frac{p(x_1|c, x_2, \dots, x_D)p(c, x_2, \dots, x_D)}{p(\mathbf{x})}$$

$$p(c|\mathbf{x}) = \frac{p(x_1|c, x_2, \dots, x_D)p(x_2|c, x_3, \dots, x_D)p(c, x_3, x_4, \dots, x_D)}{p(\mathbf{x})}$$

$$p(c|\mathbf{x}) = \frac{p(x_1|c)p(x_2|c) \dots p(x_D|c)p(c)}{p(\mathbf{x})} = \frac{p(c) \prod_{i=1}^D p(x_i|c)}{p(\mathbf{x})}$$

# NAÏVE BAYES

- ▶ Vantagens
  - ▶ Fácil de incorporar atributos de diferentes tipos
    - ▶ Cada atributo pode ser modelado por uma distribuição específica (não necessariamente Gaussiana)
  - ▶ Estimação de parâmetros sempre feita considerando apenas 1 dimensão
  - ▶ Modelo pode ser aprendido de forma incremental
- ▶ Cuidados
  - ▶ Atributos redundantes

# NAÏVE BAYES - MODELANDO OS ATRIBUTOS

- ▶ Considere  $\mathcal{X}^c$  como o conjunto de objetos da classe  $c$
- ▶ Atributos booleanos:
  - ▶ Distribuição de Bernoulli:
    - ▶  $p(x_i|c) = \theta^{x_i}(1 - \theta)^{1-x_i} \quad [x_i \in \{0, 1\}]$
    - ▶ MLE de  $\theta = \frac{1}{|\mathcal{X}^c|} \sum_{\mathbf{x} \in \mathcal{X}^c} x_i \quad [\text{e se todos os valores forem 1?}]$
- ▶ Atributos nominais:
  - ▶ Distribuição categórica:
    - ▶  $p(x_i|c) = \prod_{j=1}^V p_j^{x_i=j} \quad [x_i \in \{1, \dots, V\}]$
    - ▶ MLE de  $p_j = \frac{1}{|\mathcal{X}^c|} \sum_{\mathbf{x} \in \mathcal{X}^c} \mathbb{1}_{[x_i=j]} \quad [\text{e se todos os valores forem iguais?}]$

# NAÏVE BAYES - MODELANDO OS ATRIBUTOS

- ▶ Atributos contínuos
  - ▶ Mais comum Distribuição Normal
    - ▶  $p(x_i|c) = \mathcal{N}(\mu, \sigma^2)$
    - ▶ MLE de  $\mu$  e  $\sigma^2$  foram apresentados na aula de *Parzen Window*

## NÁÏVE BAYES - EXEMPLO

Chuta2Pes	Altura	Peso	HsTreino	Empresario	Sucesso
N	(1,5-1,6]	67	2	Fulano	N
S	(1,8-1,9]	80	8	Fulano	N
N	(1,8-1,9]	92	2	Ciclano	N
N	(1,7-1,8]	79	4	Ciclano	N
S	(1,7-1,8]	67	5	Beltrano	S
N	(1,5-1,6]	50	9	Ciclano	S
S	(1,6-1,7]	58	6	Ciclano	S
S	(1,7-1,8]	73	3	Fulano	S
N	(1,8-1,9]	90	10	Fulano	S
S	(1,5-1,6]	63	5	Ciclano	S
S	(1,7-1,8]	77	6	Beltrano	S
S	(1,7-1,8]	60	1	Beltrano	S

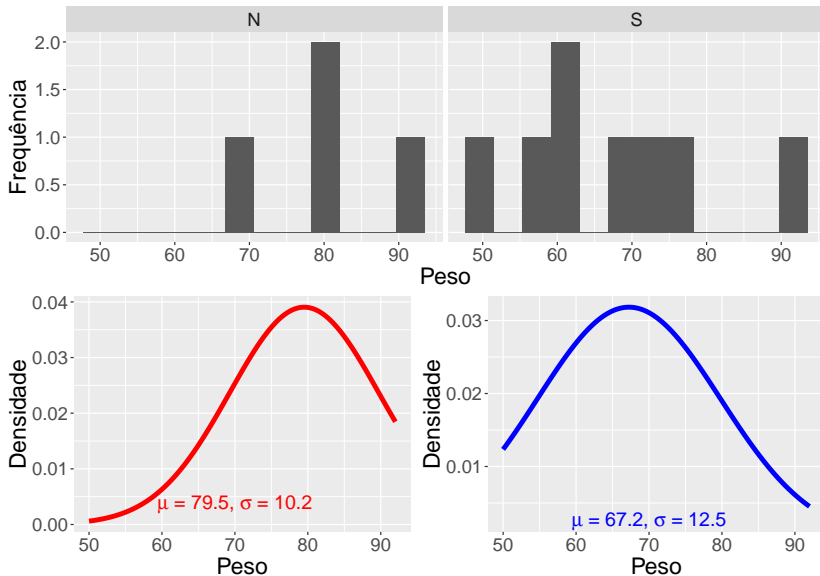
# NAÏVE BAYES - EXEMPLO

Sucesso/Chuta2Pes	N	S
N	3	1
S	2	6

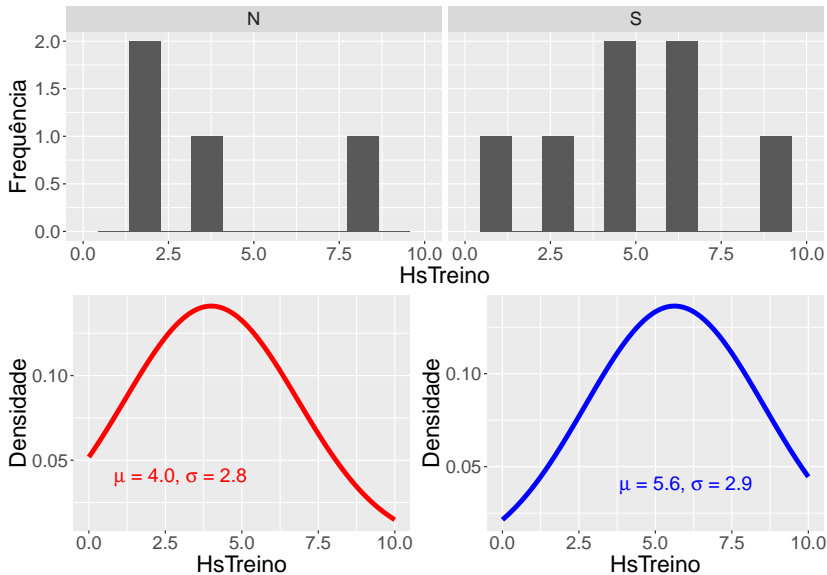
Sucesso/Empresario	Beltrano	Ciclano	Fulano
N	0	2	2
S	3	3	2

Sucesso/Altura	(1,5-1,6]	(1,6-1,7]	(1,7-1,8]	(1,8-1,9]
N	1	0	1	2
S	2	1	4	1

# NAİVE BAYES - EXEMPLO



# NAÏVE BAYES - EXEMPLO





# NAÏVE BAYES - EXEMPLO

- ▶  $p(\text{Sucesso} = S) = 8/12$   $p(\text{Sucesso} = N) = 4/12$
- ▶ Peso
  - ▶  $p(\text{Peso} = z | \text{Sucesso} = S) = \mathcal{N}(z | 67.25, 157.07)^*$
  - ▶  $p(\text{Peso} = z | \text{Sucesso} = N) = \mathcal{N}(z | 79.50, 104.33)^*$
- ▶ HsTreino
  - ▶  $p(\text{HsTreino} = z | \text{Sucesso} = S) = \mathcal{N}(z | 5.62, 8.55)^*$
  - ▶  $p(\text{HsTreino} = z | \text{Sucesso} = N) = \mathcal{N}(z | 4.00, 8.00)^*$
- ▶ Chuta2Pes
  - ▶  $p(\text{Chuta2Pes} = 1 | \text{Sucesso} = S) = 6/8$
  - ▶  $p(\text{Chuta2Pes} = 1 | \text{Sucesso} = N) = 1/4$
- ▶ Empresário
  - ▶  $p(\text{Empresario} = \text{Beltrano} | \text{Sucesso} = S) = 3/8$
  - ▶  $p(\text{Empresario} = \text{Ciclano} | \text{Sucesso} = S) = 3/8$
  - ▶  $p(\text{Empresario} = \text{Beltrano} | \text{Sucesso} = N) = 0/4$
  - ▶  $p(\text{Empresario} = \text{Ciclano} | \text{Sucesso} = N) = 2/4$

# NAÏVE BAYES - EXEMPLO

## ► Altura

- $p(\text{Altura} = (1, 5 - 1, 6] | \text{Sucesso} = S) = 2/8$
- $p(\text{Altura} = (1, 6 - 1, 7] | \text{Sucesso} = S) = 1/8$
- $p(\text{Altura} = (1, 7 - 1, 8] | \text{Sucesso} = S) = 4/8$
- $p(\text{Altura} = (1, 8 - 1, 9] | \text{Sucesso} = S) = 1/8$
- $p(\text{Altura} = (1, 5 - 1, 6] | \text{Sucesso} = N) = 1/4$
- $p(\text{Altura} = (1, 6 - 1, 7] | \text{Sucesso} = N) = 0/4$
- $p(\text{Altura} = (1, 7 - 1, 8] | \text{Sucesso} = N) = 1/4$
- $p(\text{Altura} = (1, 8 - 1, 9] | \text{Sucesso} = N) = 2/4$

# NAÏVE BAYES - EXEMPLO

- ▶ Modelo pronto
- ▶ Podemos fazer a classificação de novos dados
- ▶ Considere  $\mathbf{x}_t$  o seguinte objeto de teste:

Altura	Peso	HsTreino	Chuta2Pes	Empresario
(1,7-1,8]	73	6	S	Beltrano

## NAÏVE BAYES - EXEMPLO

- ▶  $p(\textit{Sucesso} = S) = 0.67$
- ▶  $p(\textit{Sucesso} = N) = 0.33$
- ▶  $p(\textit{Altura} = (1, 7 - 1, 8] | \textit{Sucesso} = S) = 0.50$
- ▶  $p(\textit{Altura} = (1, 7 - 1, 8] | \textit{Sucesso} = N) = 0.25$
- ▶  $p(\textit{Peso} = 73 | \textit{Sucesso} = S) = 0.029$
- ▶  $p(\textit{Peso} = 73 | \textit{Sucesso} = N) = 0.032$
- ▶  $p(\textit{HsTreino} = 6 | \textit{Sucesso} = S) = 0.135$
- ▶  $p(\textit{HsTreino} = 6 | \textit{Sucesso} = N) = 0.110$
- ▶  $p(\textit{Chuta2Pes} = 1 | \textit{Sucesso} = S) = 0.75$
- ▶  $p(\textit{Chuta2Pes} = 1 | \textit{Sucesso} = N) = 0.25$
- ▶  $p(\textit{Empresario} = \textit{Beltrano} | \textit{Sucesso} = S) = 0.38$
- ▶  $p(\textit{Empresario} = \textit{Beltrano} | \textit{Sucesso} = N) = 0.00$

## NAÏVE BAYES - EXEMPLO

$$p(Sucesso = S | \mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 0.67 \cdot 0.50 \cdot 0.029 \cdot 0.135 \cdot 0.75 \cdot 0.38$$

$$p(Sucesso = N | \mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 0.33 \cdot 0.25 \cdot 0.032 \cdot 0.110 \cdot 0.25 \cdot 0.00$$

# NAÏVE BAYES - LAPLACE SMOOTHING

- ▶ O zero observado em  $p(\textit{Empresario} = \textit{Beltrano} | \textit{Sucesso} = N)$  torna nula a probabilidade da classe ser “N”
  - ▶ Definitivamente um problema
  - ▶ Valor não visto no treinamento
- ▶ Para evitar esse problema em atributos discretos utilizamos *Laplace smoothing*
  - ▶ Adicionamos uma probabilidade pequena de um valor não visto ocorrer

# NAÏVE BAYES - LAPLACE SMOOTHING

- ▶ Para evitar esse problema em atributos discretos utilizamos *Laplace smoothing*
  - ▶ Adicionamos uma probabilidade pequena de um valor não visto ocorrer
  - ▶ Neste caso, temos:

$$p_j = \frac{1 + \sum_{\mathbf{x} \in \mathcal{X}^c} \mathbb{1}_{[x_i=j]}}{|\mathcal{X}^c| + V} \quad \text{V é o número de valores possíveis do atributo}$$

# NAÏVE BAYES - EXEMPLO

- ▶ Corrigindo nossa tabela do atributo Empresario temos:
  - ▶  $p(\text{Empresario} = \text{Beltrano} | \text{Sucesso} = S) = 4/11$
  - ▶  $p(\text{Empresario} = \text{Ciclano} | \text{Sucesso} = S) = 4/11$
  - ▶  $p(\text{Empresario} = \text{Beltrano} | \text{Sucesso} = N) = 1/7$
  - ▶  $p(\text{Empresario} = \text{Ciclano} | \text{Sucesso} = N) = 3/7$



## NAÏVE BAYES - EXEMPLO

$$p(Sucesso = S|\mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 0.67 \cdot 0.50 \cdot 0.03 \cdot 0.14 \cdot 0.75 \cdot 0.36$$

$$p(Sucesso = S|\mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 3.5239281 \times 10^{-4}$$

$$p(Sucesso = N|\mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 0.33 \cdot 0.25 \cdot 0.03 \cdot 0.11 \cdot 0.25 \cdot 0.14$$

$$p(Sucesso = N|\mathbf{x}_t) = \frac{1}{p(\mathbf{x}_t)} 1.0428371 \times 10^{-5}$$

# NAÏVE BAYES - EXEMPLO

- ▶  $p(Sucesso = S|\mathbf{x}_t) > p(Sucesso = N|\mathbf{x}_t)$ 
  - ▶ Já sabemos qual seria a classe predita
  - ▶ Não temos as probabilidades de cada classe! Não normalizamos por  $p(\mathbf{x}_t)$

$$\begin{aligned} p(\mathbf{x}_t) &= p(\mathbf{x}_t|Sucesso = S)p(Sucesso = S) \\ &+ p(\mathbf{x}_t|Sucesso = N)p(Sucesso = N) \end{aligned}$$

# NAÏVE BAYES - EXEMPLO

- Logo, já temos todos os valores para obter as probabilidades

$$p(Sucesso = S | \mathbf{x}_t) = \frac{0.000352}{(0.000010 + 0.000352)} = 0.9712575$$

$$p(Sucesso = N | \mathbf{x}_t) = \frac{0.000010}{(0.000010 + 0.000352)} = 0.0287425$$

# NAÏVE BAYES (NB)

- ▶ Estimar parâmetros para atributos contínuos com poucos objetos pode ser um problema
  - ▶ Qual característica do NB perdemos se usarmos *Parzen Window* para estimar densidade?
- ▶ É comum considerar a versão discretizada de atributos contínuos
  - ▶ Pode-se utilizar uma das técnicas discutidas na aula anterior (discretização por largura fixa ou frequência fixa)
  - ▶ Na próxima aula falaremos de outro método de discretização que considera os rótulos de classe
- ▶ Sabe-se que as estimativas de probabilidade obtidas pelo NB não são muito boas
  - ▶ Independente disso, o classificador possui ótimos resultados em diversas aplicações

# NAÏVE BAYES (NB)

- ▶ A fronteira de decisão obtida pelo NB depende das distribuições consideradas
  - ▶ No caso da distribuição Gaussiana a fronteira obtida é quadrática (relação com QDA)
- ▶ É possível incrementar o modelo considerando relação entre alguns atributos
  - ▶ Aumentamos o número de parâmetros
  - ▶ Podemos obter modelos mais realistas
- ▶ Classificadores desse tipo são conhecidos como Redes Bayesianas
  - ▶ Sendo o NB uma instância desse conjunto

# REFERÊNCIAS

R. Duda, P. Hart e D. Stork. Pattern Classification. **Seção 2.1, 2.4 e 2.6**

Slides Profa. Olga Veksler

[http://www.csd.uwo.ca/~olga/Courses/CS434a\\_541a/Lecture4.pdf](http://www.csd.uwo.ca/~olga/Courses/CS434a_541a/Lecture4.pdf)

C. Bishop. Pattern Recognition and Machine Learning. **Seção 2.1, 2.2, 2.3**

T. Hastie, R. Tibshirani e J. Friedman. The Elements of Statistical Learning. **Seção 4.3**