

# Exemplo

Name	Debt	Salary	Married	Risk
JOSÉ	High	High	yes	High
ANA	Low	High	yes	Low
JOÃO	High	High	No	Low
MARIA	High	Low	yes	High
RUI	Low	Low	yes	High

seleciona-se os vizinhos e desses escolhe-se os K mais próximos

## a) K-Nearest-Neighbors Classifier, with k=3

1º converter para as variáveis categóricas para numéricas  
Name não conta para variável preditiva, pois é uma chave

Debt	Salary	Married	Distância	Risk
1	1	1	$\sqrt{2} \leftarrow *1$	High
0	1	1	$\sqrt{1} = 1 \leftarrow *1/*2$	Low
1	1	0	$\sqrt{1} = 1 \leftarrow *1/*2$	Low
1	0	1	$\sqrt{3}$	High
0	0	1	$\sqrt{2} \leftarrow *2$	High

distância euclidiana

atributo objetivo

$\frac{2}{5} \leftarrow \text{Low}$

como atributo categórico vai-se pela MODA

Moda = low, pois é o que aparece mais

Debt = 0 ; Salary = 1 ; Married = 0  $\Rightarrow$  Risk = Low

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots} = \sqrt{(1-0)^2 + (1-1)^2 + (1-0)^2} = \sqrt{1^2 + 0 + 1^2} = \sqrt{2} \leftarrow *1$$

Naive Bayes... [não é necessário usar variáveis numéricas]

Admite que as variáveis são independentes e calcula as frequências

TABELA DE FREQUÊNCIAS

TABELA DE PROBABILIDADES

Debt \ Risk	Low	High	Low	High
Low	1+1=2	1+1=2	1/2	2/5
High	1+1=2	2+1=3	1/2	3/5
Salary	Low	High	Low	High
Low	0+1=1	2+1=3	1/4	3/5
High	2+1=3	1+1=2	3/4	2/5
Married	Low	High	Low	High
yes	1+1=2	3+1=4	1/2	4/5
No	1+1=2	0+1=1	1/2	1/5

soma de todos  
Temos de usar a correção da estimativa Laplace=1 pois estamos perante valores nulos

b)

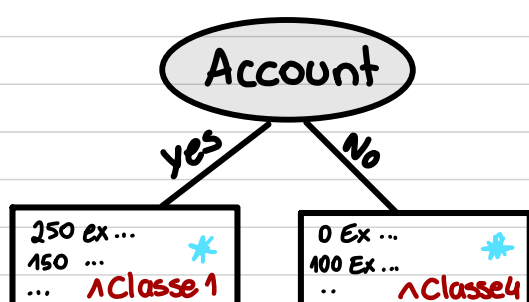
$$P(\text{Risk} = \text{Low} \mid \text{Debt} = \text{Low}, \text{Salary} = \text{High}, \text{Married} = \text{No}) = \frac{2}{5} \times \frac{1}{2} \times \frac{3}{4} \times \frac{1}{2} = \frac{6}{80},$$

$$P(\text{Risk} = \text{High} \mid \text{Debt} = \text{Low}, \text{Salary} = \text{High}, \text{Married} = \text{No}) = \frac{3}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{1}{5} = \frac{12}{375},$$

$P(\text{Risk} = \text{Low} \mid \dots) > P(\text{Risk} = \text{High} \mid \dots) \Rightarrow \text{Risk} = \text{Low}$

4 classes (1, 2, 3, 4)

250 instâncias para cada classe



\* só prevê classe 1 e 4

! classifica-se sempre pela classe majoritária

a)

		Prevista			
		^ classe 1	^ classe 2	^ classe 3	^ classe 4
Reais	classe 1	250	∅	∅	∅
	classe 2	150	∅	∅	100
	classe 3	150	∅	∅	100
	classe 4	50	∅	∅	200

prevê classe 1 e 2

c)

d)

b)  $\frac{\text{total de } \bar{n} \text{ bem classificados}}{\text{total}} = \text{taxa de erro} = 1 - \text{taxa de acerto} = 1 - 0.45 = 0.55$

taxa de acerto =  $\frac{\text{linha}}{\text{total}} = \frac{450}{1000} = 0.45$

muíto alta

c) Modelo não consegue prever a classe, pois valores nulos \*

d)

	$\wedge$ classe 1	$\wedge \neg$ classe 1 <sup><math>\rightarrow \bar{n}</math></sup>
classe 1	250 <sup>TP</sup>	$\emptyset$
$\neg$ classe 1	350 <sup>FP</sup> *	400 *

$$\text{Precision} = \frac{250}{600} = 41.7\%$$

$$\text{Recall} = 100\%$$

aceito todas as instâncias  
do classe 1 pelo que se  
escolheria este

	$\wedge$ classe 4	$\wedge \neg$ classe 4 <sup><math>\rightarrow n</math></sup>
classe 4	200 <sup>TP</sup>	50
$\neg$ classe 4	200 <sup>FP</sup>	550

↓  
precisão de cerca de 50%.

$$\text{Precision} = \frac{200}{400} = 50\%$$

$$\text{Recall} = \frac{200}{250} = 80\%$$

↓  
acaba por ser menor do  
que o certo

$$200 \rightarrow 200$$