

# Formulário [MAS]

## 1. Introdução

### Medidas de Diversidade d

▲ **Menos diversidade** associa-se a um **menor risco de previsão** (classificação ou regressão) pelo que durante a aprendizagem se procura explicar/reduzir essa diversidade

Para uma variável **nominal**:

- A  $d = 0$  quando todos os elementos de um conjunto de observações estão numa categoria
- A diversidade é **máxima** quando a distribuição é uniforme por categorias

Para uma variável **métrica**:

- A  $d = 0$  quando todos os elementos de um conjunto de observações são iguais
- A diversidade aumenta quando as observações tendem a ser diferentes

### Validação Cruzada

🌀 O recurso à validação cruzada permite obter estimativas mais realistas dos erros de aprendizagem.

**Método V-Fold:** a amostra original é particionada em  $V$  subamostras de dimensões iguais; cada submostra é excluída da aprendizagem à vez e os seus erros calculados com base no modelo estimado no treino:

$$E_1^{(V)} \dots E_{nv}^{(V)} \quad E_{nv}^{(V)} \text{ estimativa de erro para cada observação } n \text{ quando esta está na amostra de validação } v \text{ (e não está na amostra de treino).}$$

## Regressão Alvo Métrico

### Medidas de Diversidade

➤ Variância:  $Var = S^2 = \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{n}$

➤ Deviance:  $DEV^{met} = \sum_{i=1}^N (y_i - \bar{y})^2 = Var \times (n - 1)$

✓ **Residual Sum of Squares (RSS) or Sum of Squares error (SSE)**

$$\sum (y_i - \hat{y}_i)^2$$

- Comparar # resultados sobre o mesmo conjunto de dados
- Sensível à presença de outliers

✓ **Mean Absolute Error (MAE)**

$$\frac{\sum |y_i - \hat{y}_i|}{n}$$

- Comparar # resultados sobre o mesmo conjunto de dados
- Não tão sensível à presença de outliers como medidas anteriores (baseadas nos quadrados dos erros)

✓ **Mean Squared Percentage Error (MSPE)**

$$\frac{\sum \left( \frac{y_i - \hat{y}_i}{y_i} \right)^2}{n} \times 100\%$$

- O erro é comparado com a observação a que se refere

✓ **(pseudo) R-squared | R<sup>2</sup>**

$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

- Compara a diversidade residual com a diversidade inicial do alvo
- Compara o ajustamento obtido pelo modelo com o ajustamento proporcionado pela média
- X% → O modelo consegue explicar ... % da variabilidade do valor médio de  $X$
- Quando  $R^2$  é negativo, significa que o modelo é pior a prever do que a média simples.

Note: na Regressão Linear o mínimo é zero.

✓ **Relative Squared Error (RSE)**

$$\frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

✓ **Relative Absolute Error (RAE)**

$$\frac{\sum |y_i - \hat{y}_i|}{\sum |y_i - \bar{y}|}$$

- Fornece o erro absoluto das previsões relativamente a um **Modelo Nulho** que usa a média como previsão para qualquer observação.

## Classificação

Alvo Nominal

Legenda

- $K$  - nº de classes / categorias da variável nominal;
- $n_k$  - nº de observações da classe / categoria  $k$  da variável
- $n$  - nº total de observações da variável nominal

**cdg** - Medida de Desigualdade de uma distribuição ( $0 \leq G \leq \frac{1}{K}$ )

- Se  $G = 0$  então corresponde à completa igualdade
- Se  $G = 1/K$  então corresponde à completa desigualdade

➤ Coeficiente de Gini:  $G = 1 - \sum_{k=1}^K \left( \frac{n_k}{n} \right)^2$

➤ Deviance:  $DEV^{nom} = -2 \times \sum_{k=1}^K n_k \log \left( \frac{n_k}{n} \right) = 2nH$

➤ Entropia\*:

$$H = -\sum_{k=1}^K \frac{n_k}{n} \log \left( \frac{n_k}{n} \right)$$

**Entropia** (Medida de Dispersão para Dados Qualitativos)

- ▶ A Entropia de Shannon (H) permite avaliar em que medida a moda de uma variável qualitativa representa melhor ou pior os dados.
- Maior H, maior será o nº de estados possíveis, bem como a sua aleatoriedade.**

➤ Entropia Normalizada

- ▶ A Entropia (H) é a média da quantidade de informação que se associa a uma distribuição de frequências

A "quantidade de informação" associada a um acontecimento mais provável contém menos informação do que o de uma mais improvável

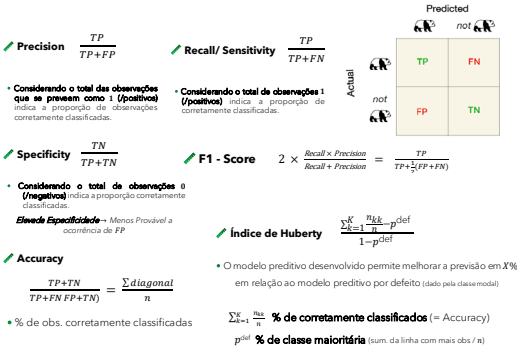
- ▶ Pode ser representado em **bits** - 2, **nats** - e ou **decits** - 10.

- ▶ Hn varia entre 0 **não há dispersão** : todos as obs. pertencem a 1 só classe/categoria (i.e., não há dispersão - "potência base prevista")

- Ao considerar um evento impossível, convencionou-se que  $0 \cdot \log(0) = 0$

1 **dispensão máxima** = há impureza, o que dificulta a obtenção de boas previsões

## Matriz de Classificação Confusion Matrix



## Classificador Naïve Bayes

💎 O classificador **Naïve Bayes** tem como base o Teorema de Bayes usando como pressuposto a **independência dos preditores**:

$$P(c_k | x) = \frac{P(x | c_k) P(c_k)}{P(x)} = \frac{P(c_k) \prod_{j=1}^p P(x_j | c_k)}{P(x)} \propto P(c_k) \prod_{j=1}^p P(x_j | c_k)$$

é proporcional a

▶ O classificador **Naïve Bayes** pode ser usado como referência para comparar o desempenho de outros métodos.

### Preditores Qualitativos

$$P(c_k | x) \propto P(c_k) \cdot \prod_{j=1}^p \prod_{i \in \{j_1, \dots, j_{L_j}\}} P(x_j = i | c_k)$$

Onde  $\prod_{i=1}^L P(x_j = i | c_k)$  é aproximado por produto de frequências relativas associadas a categorias de preditores intra-classe alvo  $c_k$

naive Bayes Classifier for Discrete Predictors

call:  
naiveBayes.default(x = x, y = y)

A-priori probabilities:

y  
0.6428571 0.3571429

conditional probabilities:

LOOK  
y  
0.4444444 0.3333333 0.2222222  
don't Play 0.0000000 0.4000000 0.6000000

$$P(\text{sunny} | \text{play})$$

### Preditores Contínuos

$$P(c_k | \underline{x}) \propto P(c_k) \prod_{j=1}^p P(x_j | c_k)$$

Onde  $\prod_{j=1}^p P(x_j | c_k)$  é aproximado por produto de Gaussians:

$$\prod_{j=1}^p \frac{1}{\sqrt{2\pi}\sigma_k} \exp \left( -\frac{(x_j - \bar{x}^k)^2}{2\sigma_k^2} \right)$$

A-priori probabilities:

x  
setosa versicolor virginica  
0.3333333 0.3333333 0.3333333

conditional probabilities:

sepal.width  
[.1] [.2]  
setosa 3.428 0.3790644  
versicolor 2.770 0.3137983  
virginica 2.974 0.3224966

Mean in col. 1  
Std. Dev in col. 2

### Função Densidade de Probabilidade da Normal

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

$$-\infty < x < \infty$$

$$E[X] = \mu \rightarrow \bar{x}$$

$$Var(X) = \sigma^2$$

## K-Nearest Neighbour (KNN)

### KNN (Não Paramétrico) VS Abordagem Paramétrica

#### VANTAGENS

Pode ser usado em **Regressão e Classificação**

É capaz de lidar com funções alvo complexas

É bastante preciso e simples de usar

Não considera pressupostos sobre os dados

Pode ser usado para uma ampla variedade de problemas.

#### DESVANTAGENS

É necessário parametrizar o valor  $K$

Custos elevados em termos computacionais: armazenamento e processamento

O KNN é muito sensível à presença de atributos irrelevantes

O KNN é naturalmente dependente da escala do conjunto de dados e das medidas de distância utilizadas

### KNN | Passo a Passo

1. Indicar  $K_i$  o número de vizinhos, e a medida de distância

O parâmetro  $K$  irá determinar o modo de aprendizagem - para uma nova observação, ele define o  $n^\circ$  de observações mais próximas que serão a base da previsão do valor do alvo.

No entanto a escolha de um valor  $K$  específico poderá ter um efeito drástico nos resultados obtidos.

Quando  $K$  tem um valor reduzido o método tende a ter baixo enviesamento, mas  $Var$  muito alta.

Conforme  $K$  cresce, o método torna-se menos flexível, apresentando uma variância reduzida, mas enviesamento elevado. Naturalmente  $k$  não pode ser superior a  $n - 1$ ...

A experimentação com diversos valores de  $K$  num contexto treino-teste ou de validação cruzada pode apoiar a seleção de  $K$ .

2. Calcular a distância entre um exemplo e os exemplos do conjunto de dados

As implementações do **KNN** no **R** recorrem, usualmente, à medida de distância **Euclidiana** (eventualmente a distâncias de **Minkowski** das quais a **Euclidiana** é um caso particular com  $p = 2$ )

$$d(x, y) = \sum_{i=1}^n \{|x_i - y_i| |p|^{1/p}} \rightarrow d(a, b) = \sqrt{(a_x - b_x)^2 + (a_y - b_y)^2}$$

**Note:** A integração de preditores qualitativos poderá ser feita mediante codificação **dummy** ( $n - 1$ ).

▶ Para que os preditores tenham protagonismo similar na previsão será necessário **NORMALIZAR**

min - max (comum - resulta em valores [0,1])

média - desvio padrão

normalize <- function(x){  
return ((x - min(x)) / (max(x) - min(x)))}

standardize <- function(x){  
return ((x - mean(x)) / sd(x))}

**Note:** Quando normalizamos temos de fazer para todas as variáveis, mesmo as binárias!

3. Ordenar as distâncias calculadas (ordem crescente)

4. Obter os valores do alvo para as  $K$  observações mais próximas do exemplo

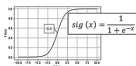
## Regressão Logística

### Métodos de Classificação

• No modelo de regressão logística para **classificação**

**binária** é usada a função **sigmóide** ou **logística**

➤ **Logistic Regression (binary response)**



• No modelo de regressão logística, para **classificação em múltiplas classes**, usa-se a função **softmax**

➤ **Multinomial Logistic Regression ou Softmax Regression**

💎 O objetivo de aprendizagem é

**Minimizar**  $ce$  (**Entropia Cruzada**)

$$ce = H(Y, \hat{y}) = -\sum_{i=1}^n \sum_{k=1}^K y_{ik} \log(\hat{y}_{ik})$$

ou, **Minimizar** a medida **Residual Deviance**

$$RESID DEV^{nom}(y_i, \hat{y}_i) = \sum_{i=1}^n 2 \left[ \sum_{k=1}^K y_{ik} \ln \left( \frac{y_{ik}}{\hat{y}_{ik}} \right) \right] = -2 \log(L) = 2 \cdot ce$$

ou, **Máximizar** o **logaritmo da Função de Verossimilhança**

$$L = \prod_{i=1}^n P(Y = y_i | X = \underline{x}_i) = \prod_{i=1}^n \pi(\underline{x}_i)^{y_i} [1 - \pi(\underline{x}_i)]^{1-y_i}$$

• Quando o **ajustamento** do modelo for **perfeito** estes **indicadores** serão **nulos**

### Função Sigmóide ou Logística

• Sendo  $Y$  uma var. dependente binária e  $X$  um preditor considere-se

$$P(x \in C_1 | x) = P(Y = 1 | x) = \pi(x)$$

$$P(x \in C_2 | x) = P(Y = 0 | x) = 1 - \pi(x)$$

O Modelo de Regressão Logística ou Modelo Logit resulta de uma relação não linear entre a probabilidade  $\pi(x)$  e  $x$ , modelada pela função **sigmóide**

$$\hat{\pi}(x) = sig(x) = \frac{1}{1 + e^{-[\beta_0 + \beta x]}} = \frac{e^{\beta_0 + \beta x}}{1 + e^{\beta_0 + \beta x}}$$

**Note:** a combinação linear  $[\beta_0 + \beta x]$  pode ser alargada a  $n$  preditores.

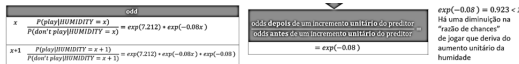
### Função Logit

A função logit é a **inversa da função logística**:

➤  $log\text{-odds}$  ou **logit** considerando  $C_1$  como **Classe de Referência**

➤ **logit** de  $\pi(x)$  pode ser qualquer número real (embora  $\pi(x) \in [0,1]$ )

**Odds Ratio** → **Razão de chances** (Comparando probabilidade de sucesso e probabilidade de fracasso)



### Escolha de Múltiplos Preditores

- É importante ter em conta que a omissão de preditores importantes (correlacionados) pode "confundir" a interpretação
- A questão da **Multicolinearidade** volta a colocar-se, tal como na regressão linear múltipla
- Notar que um preditor qualitativo deverá ser codificado com o auxílio de **variáveis dummy**.

$$< 0.4 \text{ - Baixa} \quad 0.4 \leq Corr \leq 0.69 \text{ - Moderada} \quad \geq 0.7 \text{ - Alta}$$

### Medidas de Multicolinearidade

➤ **Tolerance (TOL)** para  $X_i$ ;  $1 - R_i^2$  em que  $R_i^2$  é o o coeficiente de determinação relativo à RLM de  $X_i$  sobre os restantes preditores; assim, quanto maior a TOL (menor o  $R_i^2$ ) melhor;

→ É geralmente aceite que existe uma forte multicolinearidade se  $TOL_i < 0.1$ .

➤ **Variance Inflation Factor (VIF)**: o inverso de TOL, que também quantifica quanto a variação do estimador do coeficiente é inflacionada pela presença de multicolinearidade; assim, quanto menor o VIF, melhor;

→ Assume valores  $\geq 1$  e é geralmente aceite que há uma forte multicolinearidade se  $VIF \hat{\beta}_i > 10$ .

### Medidas de Avaliação da Regressão

▲ **Accuracy**

▲ **Índice de Huberty**  $\frac{accuracy - p_{def}}{1 - p_{def}}$

**Interpretação:** A capacidade preditiva desta regressão traduz uma melhoria relativa de X% face ao diferencial entre a capacidade preditiva de 100% e a capacidade preditiva de modelo que propõe classificar todas as observações na classe modal.

▲ **Regressão Linear**  $\rightarrow R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{SSE}{Deviance}$

▲ **R<sup>2</sup> de Hosmer and Lemeshow**  $\rightarrow R_{HL}^2 = \frac{(-2 \cdot \log(L^0)) - (-2 \cdot \log(L))}{-2 \cdot \log(L^0)}$  sendo  $L^0$  a verossimilhança sem preditores;

♦ A medida  $R_{HL}^2$  mede uma melhoria relativa na **Residual Deviance**

**Pretende-se então um valor máximo para  $R_{HL}^2$** , que indica uma capacidade preditiva máxima para o modelo considerado

A diferença  $(-2 \cdot \log(L^0)) - (-2 \cdot \log(L))$ :

- $-2 \cdot \log(L^0)$  é a **Residual Deviance** do modelo sem preditores (apenas com constantes)
- $-2 \cdot \log(L)$  é a **Residual Deviance** do modelo com preditor(es)

▲ Critério **AIC (Akaike's criterion)**  $\rightarrow AIC = -2 \log(L) + 2\theta$

em que  $\theta$  é o nº de parâmetros a estimar (nº de valores no output)

É um critério da Teoria de Informação que faz um trade-off entre o ajustamento pela Função de Verossimilhança e a Complexidade do Modelo:

♦ Pretende-se então um **valor mínimo para AIC**

# Algoritmo CART

## CART - Classification and Regression Trees

- A construção das **Árvores de Decisão** é baseada num procedimento recursivo que divide os dados para conquistar **greedy**: uma boa previsão - são modelos **explicativos** e **preditivos**.
- ✓ No **nó raiz** estão todas as observações.
- ✗ Nos **nós folha** são feitas as previsões.

# CART | Construção

- Um critério para decidir qual é a melhor **ramificação** de um nó
- Um método de **previsão** de cada nó da folha
- Regras para decidir **parar o processo de ramificação** num nó folha
- Medidas apropriadas de **desempenho do modelo**
- Critérios para  **cortes de ramos** da árvore (poda)

► No algoritmo **CART**, as **ramificações são binárias**:

- Se a variável explicativa for **métrica** ou tiver **K valores ordenados**, **K – 1 possíveis divisões** serão consideradas com base nos pontos médios entre esses valores ordenados.
- Se a variável explicativa **nominal** tiver **K categorias**, **2<sup>K</sup>-1 – 1 divisões possíveis** serão consideradas

► De entre todas as ramificações possíveis num nó, seleciona-se a que proporciona o **maior decréscimo de diversidade – Deviance (R e C)** ou **Coefficiente de Gini (C)**

# Exemplo: Play

## Árvore de Regressão

(com 2 preditores)

Ramificação do nó 5):

```

node), split, n, deviance, yval
* denotes terminal node
1) root 14 11580.0 36.79
  4) TEMPERATURE < 82.5 7 3843.0 52.14
    5) TEMPERATURE > 82.5 7 670.0 61.00
      10) TEMPERATURE < 22.5 2 12.5 67.50 *
      11) TEMPERATURE > 22.5 3 516.7 56.67 *
  
```

DEV(O)<sub>1</sub> = 11580.0  
 DEV(O)<sub>2</sub> = (50-61)\*\*2\*(7-61)\*\*2 + (75-61)\*\*2 + (45-61)\*\*2 + (65-61)\*\*2 = 2670  
 DEV(O)<sub>3</sub> = (70-67.5)\*\*2 + (65-67.5)\*\*2 + 12.5  
 DEV(O)<sub>4</sub> = (50-56.7)\*\*2 + (75-56.7)\*\*2 + (45-56.7)\*\*2 + 516.7

→ O maior decréscimo de diversidade no nó 5) foi obtido com esta ramificação

► O decréscimo de diversidade na ramificação de nó 5) (5/14) 670-[(2/14) 12.5 + (3/14) 516.7] = 126.786

# Exemplo: Cloudy

## Árvore de Classificação

(com ajustamento perfeito)

```

node), split, n, deviance, yval, (prob)
* denotes terminal node
1) root 14 18.250 Play (0.3571 0.6429)
  2) OUTLOOK: rain, sunny 10 13.860 Play (0.5000 0.5000)
    4) HUMIDITY < 82.5 5 5.004 Play (0.2000 0.8000)
      6) TEMPERATURE < 19.1 0.000 Don't Play (1.0000 0.0000) *
      9) TEMPERATURE > 19.4 0.000 Play (0.0000 1.0000) *
    5) HUMIDITY > 82.5 5 5.004 Don't Play (0.8000 0.2000)
      10) HUMIDITY > 95.5 4 0.000 Don't Play (1.0000 0.0000) *
      11) HUMIDITY > 95.5 1 0.000 Play (0.0000 1.0000) *
  3) OUTLOOK: cloudy 4 0.000 Play (0.0000 1.0000) *
  
```

Accuracy=1

O decréscimo de diversidade:

Na ramificação de nó 3):  
 (10/14)\* 13.863 - 2\*(5/14)\* 5.004 = 6.328

Na ramificação de nó 2):  
 (10/14)\* 13.863 - 2\*(5/14)\* 5.004 = 6.328

Na ramificação de nó 4):  
 (5/14)\* 5.004 - 2\*(2/14)\* 0.000 - 3\*(3/14)\* 5.004 = 5.004

Na ramificação de nó 5):  
 (5/14)\* 5.004 - 2\*(2/14)\* 0.000 - 3\*(3/14)\* 5.004 = 5.004

Na ramificação de nó 6):  
 (2/14)\* 0.000 - 2\*(1/14)\* 0.000 - 0\*(1/14)\* 1.000 = 0.000

Na ramificação de nó 9):  
 (1/14)\* 1.000 - 2\*(0/14)\* 0.000 - 1\*(1/14)\* 1.000 = 1.000

Na ramificação de nó 10):  
 (4/14)\* 0.000 - 2\*(2/14)\* 0.000 - 0\*(0/14)\* 1.000 = 0.000

Na ramificação de nó 11):  
 (1/14)\* 1.000 - 2\*(0/14)\* 0.000 - 1\*(1/14)\* 1.000 = 1.000

# 3. Regras de Paragem

- Um decréscimo **min** de diversidade da var. dependente
- Um **nº** **mín** de observações em nós "pais" e nós "filhos"
- Um **nº** **máx** de níveis na árvore

# 4. Medidas de Desempenho

## Regressão

Deviance:  $\sum_{i \in A_j} \sum_{i \in O} (y_i - \hat{y}_i)^2$

No final, a Residual Deviance será quantificada pela Deviance nos nós 0 pertencentes ao conjunto de nós folha ( $A_k$ )

## Classificação

Accuracy:

- Nº num 0:  $P(O_i) = \%$  obs. corretamente classificadas em 0
- Na **árvore A** (soma, ponderada pela frequência relativa dos nós, da accuracy em todos os nós folha/conjunto  $A_j$ ):

$P(A) = \sum_{i \in A_j} p(O_i)P(O_i)$

# Exame 2021

- $\text{knm.reg(Boston}_n, y = \text{Boston}_n\$medv, k=1, \text{algorithm} = \text{"brute"})$ 
    - Validação **one-hold-out** → **n** é o **nº de folds**
    - Como **"brute"** o KNN usa o vizinho mais próx. (e não a própria obs).
  - Calcule o **SSE – Sum of Square Error - (ou RSS – Residual Sum of Square)**

$$SSE = \sum (y_i - \hat{y}_i)^2 \quad R^2 = \frac{\sum y_i - \hat{y}_i}{\sum y_i - \hat{y}_i} = 1 - \frac{SSE}{Deviance} = \dots$$

$$\frac{x_i - \min(X)}{\max(X) - \min(X)}$$

Distância Euclidiana entre  $x = (x_1, \dots, x_n)$  e  $y = (y_1, \dots, y_n)$ :

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- A pertinência de normalizar no **KNN** é a **amplitude** das variáveis.

# 2.

- $H_n = \dots = 0.87 \rightarrow$  O valor está muito próximo de 1 indicando **dispersão elevada** pelas 3 categorias de origem.
- O modelo de **RL Multinomial** considera equações baseadas na função **softmax** para determinar as **probabilidades** das classes alvo;

Para obter o AIC soma-se 2 vezes **nº parâmetros estimados** à **Residual Dev.**

**NOTA:** O **intersept** é a classe de referência - Não aparece no **output**

# 3.

- A **Árvore de Classificação** aprende sobre todas as observações do **dataset** (basta ver se o **n** na **root** – onde está a **Deviance** da **variável alvo** tmb – se é igual ao **n** do **dataset**)
- $\text{Índice de Huberty} = \frac{\text{Accuracy} - \% \text{ Classe Modal}}{1 - \% \text{ Classe Modal}} = \dots \approx 0.497$

A capacidade preditiva desta árvore traduz uma melhoria relativa de 49.7% face ao diferencial entre a capacidade preditiva de 100% e a capacidade preditiva de modelo que propõe classificar todas as observações na classe modal.

# Exame 2022

- $DEV^{nom} = -2 \times \left[ n_{No} \times \log\left(\frac{n_{No}}{n}\right) + n_{Yes} \times \log\left(\frac{n_{Yes}}{n}\right) \right]$
  - No **KNN** ter atenção à multicolinearidade
  - A constituição das amostras de treino/teste **não** recorre a um processo **estratificado**

# Standardização $\mu/\sigma$

Para a observação 1 do conjunto de teste, a **aparência** previu, em meses, encontra-se acima da média (observada neste conjunto); o desvio face à média deste valor observado é menor que o desvio padrão.

Standardização  $\mu/\sigma$  do valor  $x_i$  da variável  $X$ :  $\frac{x_i - \bar{X}}{S_X}$

Para a observação 1 do conjunto de teste, a **aparência** previu, em meses, encontra-se acima da média (observada neste conjunto); o desvio face à média deste valor observado é menor que o desvio padrão.

**c) Decréscimo de Diversidade:**

N6 3) deviance= 22440 (nº obs = 71)  
 N6 6) deviance= 9439 (nº obs = 46)  
 N6 7) deviance= 4588 (nº obs = 25)  
 Decréscimo: (71/292) \* 22440 - ((46/292) \* 9439 + (25/292) \* 4588) = 3576.527

- Para efetuar uma ramificação **rtree\_large.employee** **segue-se** **valor= 0.0001** como decréscimo mínimo da deviance  
 $F = \text{mindev} \rightarrow 0.0001 \times \text{Deviance inicial} = 6 \text{ o mínimo valor de decréscimo da Deviance}$
- A **Residual Deviance** que se associa ao modelo obtido é 20760 **V**  
 As previsões em **rtree.employee** obtêm-se apenas em 6 nós **V** – é só nos nós folha  
 É no nó **1** que se observa o maior valor da deviance **F-3) não é um nó folha**

# 2.

- $\text{Índice de Huberty} = 0 \rightarrow$  A capacidade preditiva do modelo nada acrescenta à que se obtém mediante a mera afetação à classe modal **"No"**.
- Calcule a probabilidade de um indivíduo na categoria **"Manager"** e com **"Salary=40000"** estar na classe minoritária ("Yes" de "minority") de acordo com o modelo **Naive Bayes** estimado (e indique a classe em que será classificado de acordo com o mesmo modelo).

```

Naive Bayes Classifier for Discrete Predictors
call:
naiveBayes.default(x = x, y = y)
A-priori probabilities:
y      No      Yes
0.7945205 0.2054795

conditional probabilities:
y      jobcat      Administrative      Manager
No      0.741      0.034      0.214
Yes     0.050      0.117      0.033

salary
y      No      Yes
No      36511.64      17012.457
Yes     27261.67      8225.067
  
```

$P(\text{minority} = \text{Yes} \mid \text{jobcat} = \text{Manager} \text{ e } \text{Salary} = 40000) \propto$  **é proporcional a**

$P(m = \text{Yes}) \times P(j = \text{Manager} \mid m = \text{Yes}) \times P(S = 40000 \mid m = \text{Yes}) =$

$0.205 \times 0.033 \times 0.000022$

já que  $\phi(40000; \mu = 27261.67; \sigma = 8225.067) = \frac{1}{8225.067 \times 2\pi} \exp\left\{-\frac{1}{2} \left(\frac{40000 - 27261.67}{8225.067}\right)^2\right\} = 0.0000222$

$P(\text{minority} = \text{No} \mid \text{jobcat} = \text{Manager} \text{ e } \text{Salary} = 40000) \propto$  **é proporcional a**

$P(m = \text{No}) \times P(j = \text{Manager} \mid m = \text{No}) \times P(S = 40000 \mid m = \text{No}) =$

$0.795 \times 0.224 \times (40000; \mu = 36511.64; \sigma = 17012.457) = 0.795 \times 0.224 \times 0.000021$

já que  $\phi(40000; \mu = 36511.64; \sigma = 17012.457) = \frac{1}{17012.457 \times 2\pi} \exp\left\{-\frac{1}{2} \left(\frac{40000 - 36511.64}{17012.457}\right)^2\right\} = 0.000021$

**Logo:**  $P(\text{minority} = \text{Yes} \mid \text{jobcat} = \text{Manager} \text{ e } \text{Salary} = 40000) = \frac{0.205 \times 0.033 \times 0.000022}{0.205 \times 0.033 \times 0.000022 + 0.795 \times 0.224 \times 0.000021} = \frac{0.00008}{0.0708} = 0.038$

$P(\text{minority} = \text{No} \mid \text{jobcat} = \text{Manager} \text{ e } \text{Salary} = 40000) = \frac{0.795 \times 0.224 \times 0.000021}{0.205 \times 0.033 \times 0.000022 + 0.795 \times 0.224 \times 0.000021} = \frac{0.07}{0.00008 + 0.07} = 0.962$

► Pelo que a obs. se classifica em **minority = No**

# 3.

- Escolha dos Preditores no **KNN**
  - Ao procurar 2 preditores não correlacionados entre si tem-se em conta o pressuposto implícito no uso da **Distância Euclidiana** no **KNN**:
    - Pressuposto de dimensões não correlacionadas/ortogonais
  - Ter atenção à classe **positive**  $\Delta$ 

$$\text{Sensitivity} = \frac{TP}{TP + FN} = \frac{10}{10 + 17} = 0.370$$

A capacidade de prever corretamente casos que são positivos é 37%

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 32} = 0.238$$

De entre os casos que se preveem positivos registam-se 23.8% de previsões corretas.
  - A **sensitivity** é maior que a **precision**, i.e., o modelo favorece a previsão correta de casos efetivamente positivos face à capacidade de prever positivos com precisão
- Os dados normalizados entre 0 e 1 e entre a **média** e o **desvio padrão**, resultam em vizinhos mais próximos **diferentes**.
- Que valor da distância é considerado pelo **knn.employee**, entre a obs. 2 de teste e a primeira observação do conjunto de treino?

Considerando os preditores salary e jobtime (standardizados) temos os dados seguintes:

Obs.	Teste	salary	jobtime
2º	obs. de Teste	-0.742	1.621
1º	obs. de Treino	-0.521	-1.471

**Distância Euclidiana:**  $\sqrt{(-0.742 - (-0.521))^2 + (1.621 - (-1.471))^2} = 3.100$

- $\text{knn(employee\_train\_s[, predictors], employee\_test\_s[, predictors], employee\_train\_s$minority, 3, prob = TRUE)}$ 
  - A **distância Euclidiana** neste **KNN** é feita com uma obs. de treino e uma de teste.
  - Não** foi feita **one-hold-out cross-validation** → Só se aplica quando não se põe o conjunto de teste

# 4.

- $R^2 = 1 - \frac{\text{Deviance/RSS}_{\text{final}}}{\text{RSS}_{\text{inicial}} \times n}$
- $\text{rtree\_large.employee} <- \text{tree(salary[1000 - gender:educ:minority, data=employee\_train, control=tree.control(nrow(employee\_train), mfcut = 2, minsize = 4, mindev = 0.0001), split = "deviance")}$ 

**Ver como é que target está na tree**  $\Delta$

Para ver o valor estimado pela **Árvore de Regressão** temos de ver o **yval** → **corresponde** (se necessário)

# Exame 2023

- $\bullet$  Cálculo da **Entropia não se considera NAs**
- $\text{Accuracy} = (87 + 288) / (87 + 188 + 42 + 288) = 375 / 605 = 0.6198$ 

(Sim) Accuracy é maior que **Default accuracy** que resulta da afetação de todas as obs. à classe modal, cuja freq. relativa é 330/605 = 0.5454
  - Cálculos da probabilidade  $P(i \in Z.3 = 4 \text{ e } i \in Z.11 = 3 \mid i \in 1.0 = \text{"Não"})$  - **NB** (atendendo à **independência dos preditores, que se pressupõe**)
 
$$\dots = P(12.3 - 41 \mid 111 \text{ } 0.1 - \text{Não}) \times P(2.11 - 31 \mid 10.1 - \text{Não}) = **$$

Tendo em conta que as distribuições condicionais acima são modeladas por Gaussianas com fdp

$f(x) = \exp(-)$  com os seguintes parâmetros referidos à classe **Não** (indicados no modelo NB)

$$i \in Z.3 \quad \mu = 3.221 \mid \sigma = 0.71 \quad i \in Z.11 \quad \mu = 3.713 \mid \sigma = 0.548$$

$$** = \frac{1}{0.71 \sqrt{2\pi} \times 3.1416} \exp\left\{-\frac{1}{2} \left(\frac{4 - 3.221}{0.71}\right)^2\right\} \times \frac{1}{0.548 \sqrt{2\pi} \times 3.1416} \exp\left\{-\frac{1}{2} \left(\frac{3 - 3.713}{0.548}\right)^2\right\} =$$

$$= 0.3077851 \times 0.3122714 = 0.09611248$$
- $\bullet$  A diferença entre o **AUC** e **Residual Deviance** **não** corresponde ao **nº** de parâmetros estimados (corresponde a 2x esse número)
  - MRL** estimado
 
$$\log\left(\frac{P(\text{sim, investe em FIF})}{1 - P(\text{sim, investe em FIF})}\right) = -5.1025 + 0.5386i \in Z.3 + 0.9115i \in Z.11$$

Interpretação de 0.5386 : Quando  $i \in Z.3$  aumenta 1 unidade (sendo assim atribuído mais um grau de importância ao seguro de riscos de exploração quando se pondera a decisão de investir em FIF), há uma variação de  $\exp(0.5386) = 1.7136$  (i.e. um aumento) da razão de chances de investir em FIF face à de não investir
  - Probabilidade de a 1ª Observação pertencer à classe **"Não"**

$$P(\text{Não} \mid i \in Z.3 = 4, i \in Z.11 = 3) = 1 - \frac{1}{\exp(-[-5.1025 + 0.5386 + 4 + 0.9115 \times 3])} = 1 - 0.4468 = 0.5532$$
- Em que **nó folha** se classifica
    - A 1º obs. do conjunto de teste tem  $i \in 1.0 = \text{Não}$ ,  $2.3 = 4$ ,  $12.10 = 3$ ,  $12.11 = 3$ ,  $2.13 = 3$ ,  $2.14 = 1$
    - $i \in 2.11 = 3 < 3.5 \Rightarrow$  Segue para o nó 2)
    - $i \in 2.13 = 3 < 3.5 \Rightarrow$  Segue para o nó 4)
    - $i \in 2.3 = 4 > 2.5 \Rightarrow$  Segue para o nó 9)
    - $i \in 2.14 = 1 < 1.5 \Rightarrow$  Segue para o nó 18) (nó folha onde é classificada em **Sim**)
  - Originalmente  $i \in 1.0 = \text{Não}$  pelo que esta obs. é incorretamente classificada.
  - Calcula-se o **Índice de Huberty** para o conjunto treino e teste
 

Face à melhoria possível sobre a **Accuracy** obtida mediante a classificação na classe modal, a **AC** proporciona uma melhoria relativa de 24.4% na amostra de treino e de 7.46% na amostra de teste; é uma diferença substancial que indicia **sobreajustamento**.

# 5.

- KNN - Com Treino/ Teste**

Para cada observação da amostra de **teste** forma considerados os vizinhos mais próximos na amostra de **treino**

  - Não seria necessário já que todos os preditores se encontram na mesma escala (1 made importante... 4 made importante)
- Os dados dos preditores deveriam ter sido **normalizados** antes de proceder à análise do **KNN**
- RSE (Relative Squared Error)** correspondente às estimativas obtidas pelo **KNN** e interprete.
 
$$RSE = \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = \frac{105988}{21.65^2 \times (226-1)} = \frac{105988}{105462.6} = 1.004982$$

A soma dos erros quadráticos associados ao modelo é ligeiramente superior à soma dos erros quadráticos em torno da média (**RSE** ligeiramente acima de 1) o que indica que o **KNN** tem um desempenho pior que o simples modelo que adota como previsão a média dos valores alvo.

ii) Em igualdade de condições de risco e rentabilidade, com que probabilidade investiria num Fundo de Investimento Florestal?

```

knn(employee_train_s[, predictors], employee_test_s[, predictors],
employee_train_s$minority, 3, prob = TRUE)
  
```