

Data Mining com o programa Weka

Alípio Jorge



LIAAD LABORATÓRIO DE INTELIGÊNCIA ARTIFICIAL
E APOIO À DECISÃO



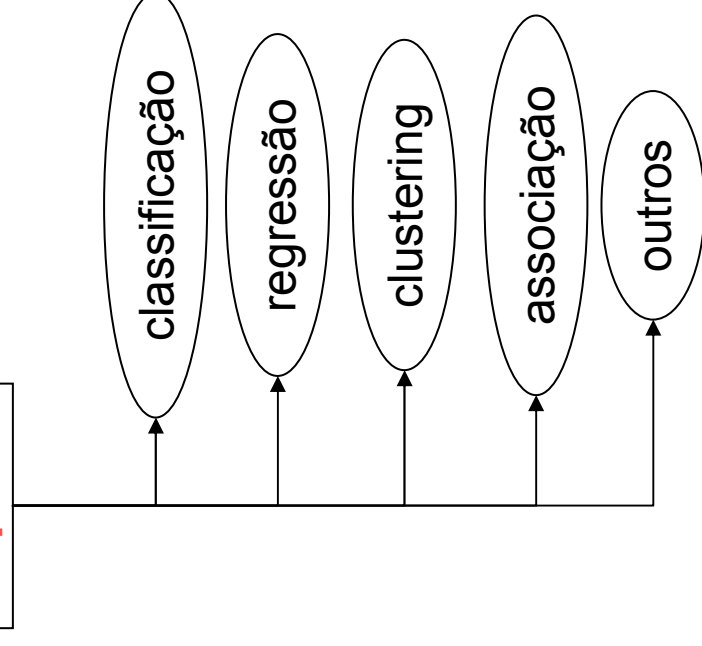
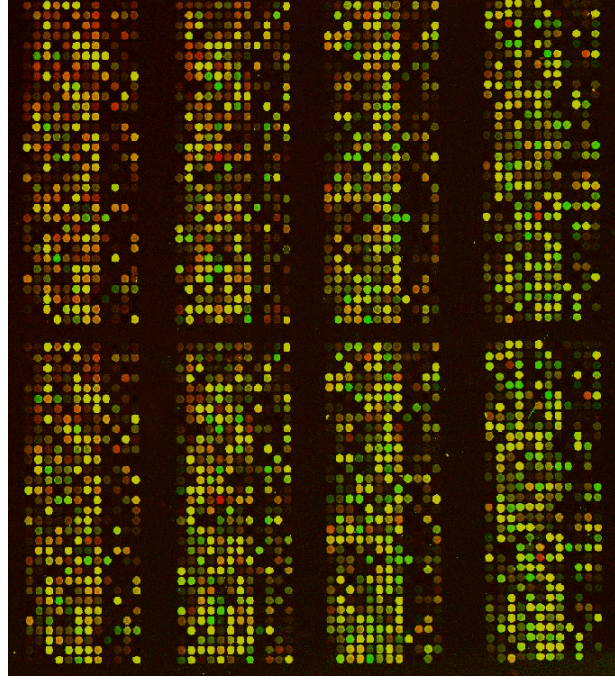
Sumário

- ❑ **O que é Data Mining?**
- ❑ **Dados**
- ❑ **Classificação**
- ❑ **Regressão**
- ❑ **Clustering**
- ❑ **Associação**
- ❑ **Processo de Data Mining**

O que é Data Mining?

❑ Procura de **padrões úteis** em grandes quantidades de dados

- padrão: motivo que se repete com alguma frequência
- útil: o padrão deve servir para resolver um **problema**



Dados

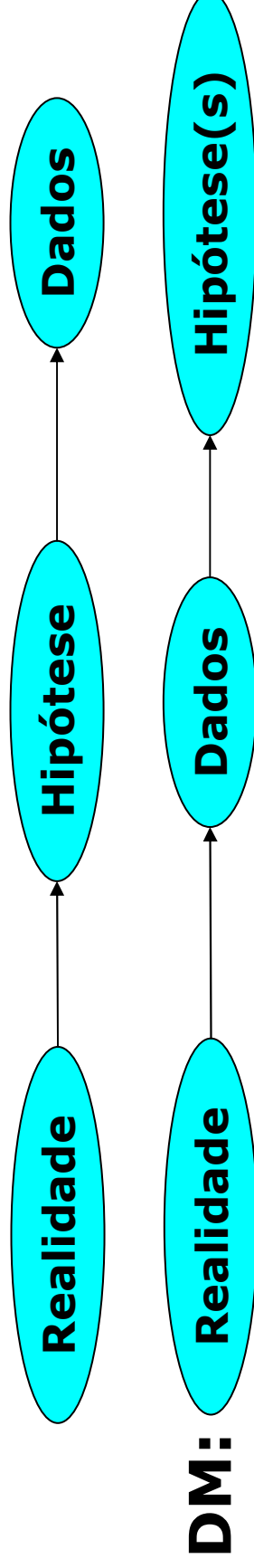
☐ Tipicamente

– tabela com N indivíduos e M atributos (depois de processados)

sexo	IDA	CIV	ESCOL	PROF	HDORM	ACTIV	DESP	TAB	ALC	CAF	Peso	ALT	IMC	Colest
m	40anos	c	EstSuperiores	sup	8ha10h	pouca	nao	nao	bebe	sim	70a60Kg	m160	normal	alto
f	40anos	c	12ano	int	6ha8h	pouca	sim	nao	bebe	sim	70a60Kg	m150	excessopeso	baixo
f	50anos	s	9Classe	sup	6ha8h	pouca	nao	nao	nao	sim	50a60Kg	m150	normal	baixo
m	60anos	c	4Classe	semi-qual	6ha8h	pouca	nao	ex	bebe	nao	mais80	m160	excessopeso	medio
f	60anos	c	4Classe	sem-prof	menos6h	alguma	nao	nao	nao	sim	50a60Kg	m150	excessopeso	medio
f	50anos	c	EstSuperiores	sup	8ha10h	pouca	sim	nao	ocas	sim	50a60Kg	m150	normal	medio
m	40anos	c	4Classe	esp-man	mais10h	alguma	nao	ex	bebe	sim	mais80	m170	excessopeso	baixo
m	40anos	c	EstSuperiores	sup	6ha8h	nenhuma	sim	ex	bebe	sim	70a60Kg	m170	normal	baixo
m	40anos	c	4Classe	esp-n-man	6ha8h	pouca	sim	nao	bebe	sim	80a70kg	m160	excessopeso	medio
m	60anos	c	EstSuperiores	sup	8ha10h	nenhuma	sim	ex	bebe	sim	mais80	m170	excessopeso	medio
m	60anos	c	4Classe	semi-qual	8ha10h	pouca	nao	ex	ex	sim	70a60Kg	m180	normal	alto
m	50anos	c	9Classe	esp-n-man	8ha10h	pouca	nao	nao	bebe	sim	70a60Kg	m150	excessopeso	medio
f	40anos	v	4Classe	esp-n-man	8ha10h	nenhuma	nao	nao	nao	sim	50a60Kg	m160	normal	baixo
m	40anos	c	9Classe	esp-n-man	6ha8h	nenhuma	sim	fuma	bebe	sim	mais80	m160	obesidade	alto
f	50anos	c	12ano	int	6ha8h	alguma	sim	ex	bebe	sim	70a60Kg	m150	excessopeso	medio
m	50anos	c	12ano	int*	mais8h	pouca	sim	nao	bebe	sim	80a70kg	m170	norm	medio
	40anos	d	EstSuperiores		pouca	pouca	sim				mais80	m160		

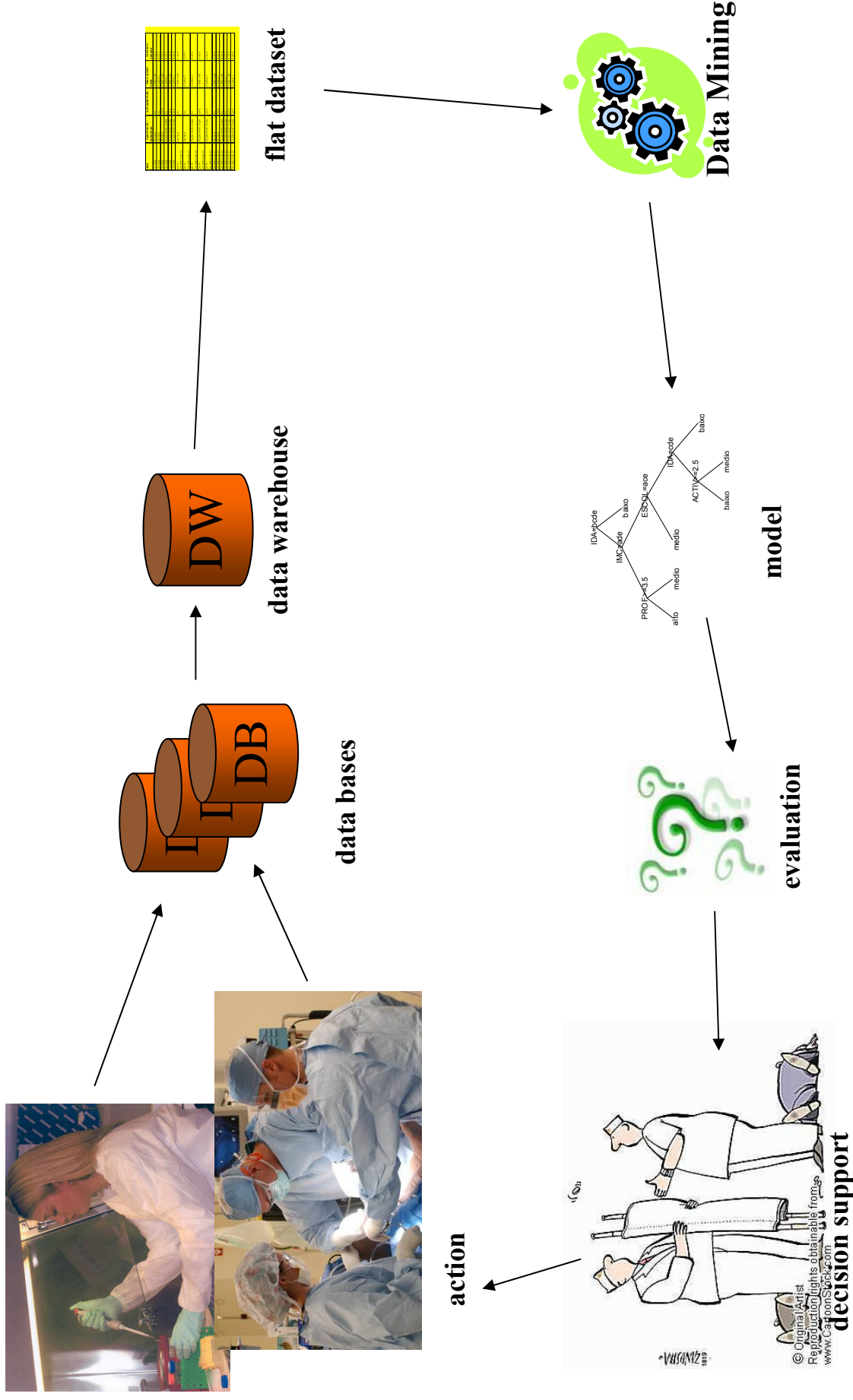
Dados em DM

- ❑ dados observados (não experimentais)
 - recolhidos para outros fins
 - os dados não são recolhidos em função de uma hipótese
 - experiência de Galileu



- ❑ grandes quantidades de dados
 - scale-up dos métodos
 - preparação dos dados

Data mining: o ciclo



Fontes de dados

☐ **Actividade clínica (transaccional)**

- **consultas**
- **exames**
- **tratamentos**
- **diagnósticos**

☐ **Laboratório**

- **genética, microarrays**
- **proteínas, compostos químicos**

☐ **Textos**

- **abstracts**
- **receitas médicas**
- **...**

O output: Padrões, Hipóteses, Modelos,...

- ❑ **O sumário derivado pelos algoritmos de data mining:**
- ❑ **Padrão**
 - afirmação sobre as variáveis envolvidas, válida para uma parte do espaço das variáveis (local)
- ❑ **Modelo**
 - afirmação sobre as variáveis, válida para todo o espaço (global)
- ❑ **Hipótese**
 - Modelo ou padrão candidato

Tipos de problemas de DM

- ❑ Classificação
- ❑ Regressão
- data mining dirigido

- ❑ Clustering
- ❑ Associação
- data mining exploratório

- ❑ Outros
 - o conjunto de tipos de problemas não é fechado

Problemas de Classificação

□ Dados

- N casos de classes conhecidas

□ Descobrir

- descobrir uma relação funcional $f(X) = y$
- que a um novo caso X, faça corresponder uma classe y

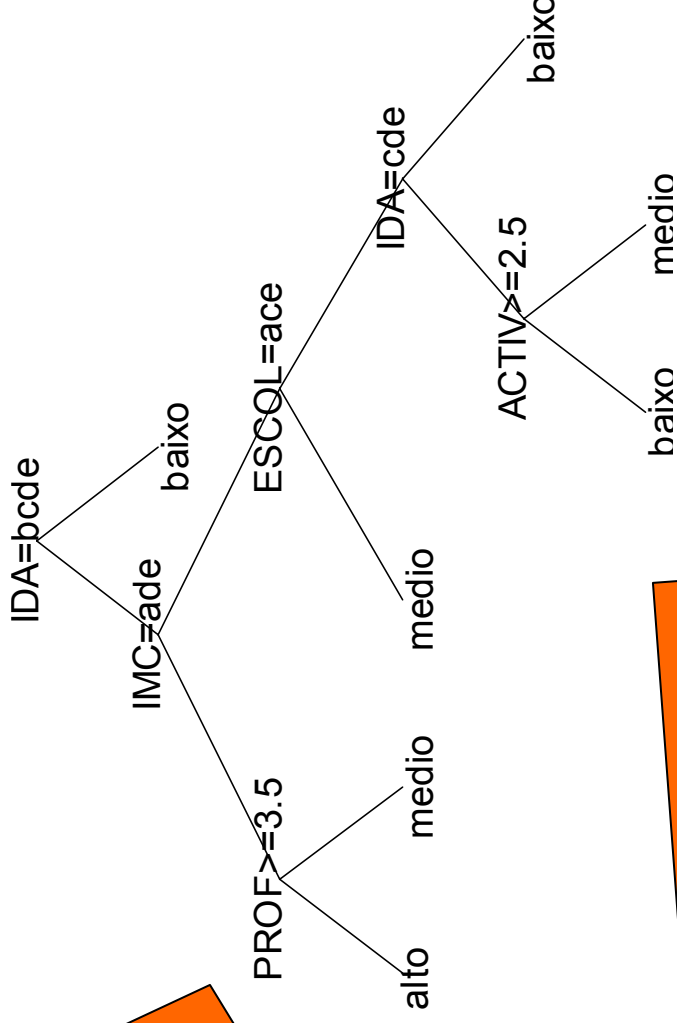
3 classes:
- alto
- médio
- baixo

sexo	IDA	CIV	ESCOL	PROF	HDORM	ACTIV	DESP	TAB	ALC	CAF	Peso	ALT	IMC	Colest
m	40anos	c	EstSuperiores	sup	8ha10h	pouca	nao	nao	bebe	sim	70a60Kg	m160	normal	alto
f	40anos	c	12ano	int	6ha8h	pouca	sim	nao	bebe	sim	70a60Kg	m150	excessopeso	baixo
f	50anos	s	9Classe	sup	6ha8h	pouca	nao	nao	nao	sim	50a60Kg	m150	normal	baixo
m	60anos	c	4Classe	semi-qual	6ha8h	pouca	nao	ex	bebe	nao	mais80	m160	excessopeso	medio
f	60anos	c	4Classe	sem-prof	menos6h	alguma	nao	nao	nao	sim	50a60Kg	m150	excessopeso	medio
f	50anos	c	EstSuperiores	sup	8ha10h	pouca	sim	nao	ocas	sim	50a60Kg	m150	normal	medio
m	40anos	c	4Classe	esp-man	mais10h	alguma	nao	ex	bebe	sim	mais80	m170	excessopeso	baixo
m	40anos	c	EstSuperiores	sup	6ha8h	nenhuma	sim	ex	bebe	sim	70a60Kg	m170	normal	baixo
m	40anos	c	4Classe	esp-n-man	6ha8h	pouca	sim	nao	bebe	sim	80a70kg	m160	excessopeso	medio
m	60anos	c	EstSuperiores	sup	8ha10h	nenhuma	sim	ex	bebe	sim	mais80	m170	excessopeso	medio
m	60anos	c	4Classe	semi-qual	8ha10h	pouca	nao	ex	ex	sim	70a60Kg	m180	normal	alto
m	50anos	c	9Classe	esp-n-man	8ha10h	pouca	nao	nao	bebe	sim	70a60Kg	m150	excessopeso	medio
f	40anos	v	4Classe	esp-n-man	8ha10h	nenhuma	nao	nao	nao	sim	50a60Kg	m160	normal	baixo
m	40anos	c	9Classe	esp-n-man	6ha8h	nenhuma	sim	fuma	bebe	sim	mais80	m160	obesidade	alto
f	50anos	c	12ano	int	6ha8h	alguma	sim	ex	bebe	sim	70a60Kg	m150	excessopeso	medio
m	50anos	c	12ano	int	6ha8h	pouca	sim	nao	bebe	sim	80a70kg	m170	normal	medio
m	40anos	d	EstSuperiores	sup	6ha8h	pouca	sim	fuma	bebe	sim	70a60Kg	m160	normal	alto

Classificação

sexo	IDA	CV	ESCOL	PROF	HOORI	ACTIV	DESP	TAB	JALC	CAF	Paro	ALT	IMC	Colist
m	40anos	t	EstSuperores	sup	brat0h	pouca	nao	nao	bebe	sim	70a0Kg	m160	normal	alto
f	40anos	t	2ano	inf	brat0h	pouca	sim	nao	bebe	sim	70a0Kg	m150	excessosop	baixo
f	50anos	s	3Classe	sup	brat0h	pouca	nao	nao	nao	sim	50a0Kg	m150	normal	baixo
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	ex	bebe	nao	nao	m160	excessosop	medio
f	60anos	t	4Classe	semqual	brat0h	pouca	nao	nao	nao	nao	nao	m150	excessosop	medio
m	60anos	t	4Classe	semqual	brat0h	pouca	sim	nao	nao	nao	nao	m150	normal	alto
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	ex	bebe	sim	nao	m170	excessosop	baixo
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	nao	nao	nao	nao	m170	normal	baixo
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	ex	bebe	sim	nao	m150	excessosop	medio
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	ex	bebe	sim	nao	m160	normal	alto
m	60anos	t	4Classe	semqual	brat0h	pouca	nao	ex	bebe	sim	nao	m160	excessosop	baixo
f	40anos	t	3Classe	esp-nmea	brat0h	nao	nao	nao	nao	nao	nao	m160	normal	baixo
f	40anos	t	3Classe	esp-nmea	brat0h	nao	nao	nao	nao	nao	nao	m160	obsusop	alto
f	50anos	t	2ano	inf	brat0h	nao	nao	nao	nao	nao	nao	m150	excessosop	medio
m	50anos	t	2ano	inf	brat0h	pouca	sim	nao	nao	nao	nao	m170	normal	medio
m	40anos	t	EstSuperores	sup	brat0h	pouca	sim	nao	bebe	sim	70a0Kg	m160	normal	alto

MINING



DEPLOYMENT

Previsão:
Colect=alto

Novo caso:
Sexo: m
IDA: 39
...

Weka

- ❑ **WEKA (da Universidade de Waikato, Nova Zelândia)**
 - <http://www.cs.waikato.ac.nz/ml/weka/>
 - é um suite de DM
 - tem uma vasta colecção de métodos disponíveis
 - faz pré-processamento
 - faz (alguma) visualização
 - permite organizar experiências
 - é gratuito
 - é open source (java)

Weka

Weka 3.5.6 - Explorer

Program Applications Tools Visualization Windows Help

Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose **None** Apply

Current relation: iris
Relation: iris
Instances: 150

Attributes: 5

Selected attribute: Name: sepalength
Missing: 0 (0%)
Distinct: 35
Type: Numeric
Unique: 9 (6%)

Statistic	Value
Minimum	4.3
Maximum	7.9
Mean	5.843
StdDev	0.828

Class: class (Nom) Visualize All

No.	Name	Value
1	sepalength	
2	sepalwidth	
3	petallength	
4	petalwidth	
5	class	

Remove

Status: OK

Weka: dados

- ❑ ARFF (Attribute-Relation File Format)
 - o formato nativo do Weka

```
@RELATION iris
@ATTRIBUTE sepallength REAL
@ATTRIBUTE sepalwidth REAL
@ATTRIBUTE petallength REAL
@ATTRIBUTE petalwidth REAL
@ATTRIBUTE class {Iris-setosa,Iris-versicolor,Iris-virginica}

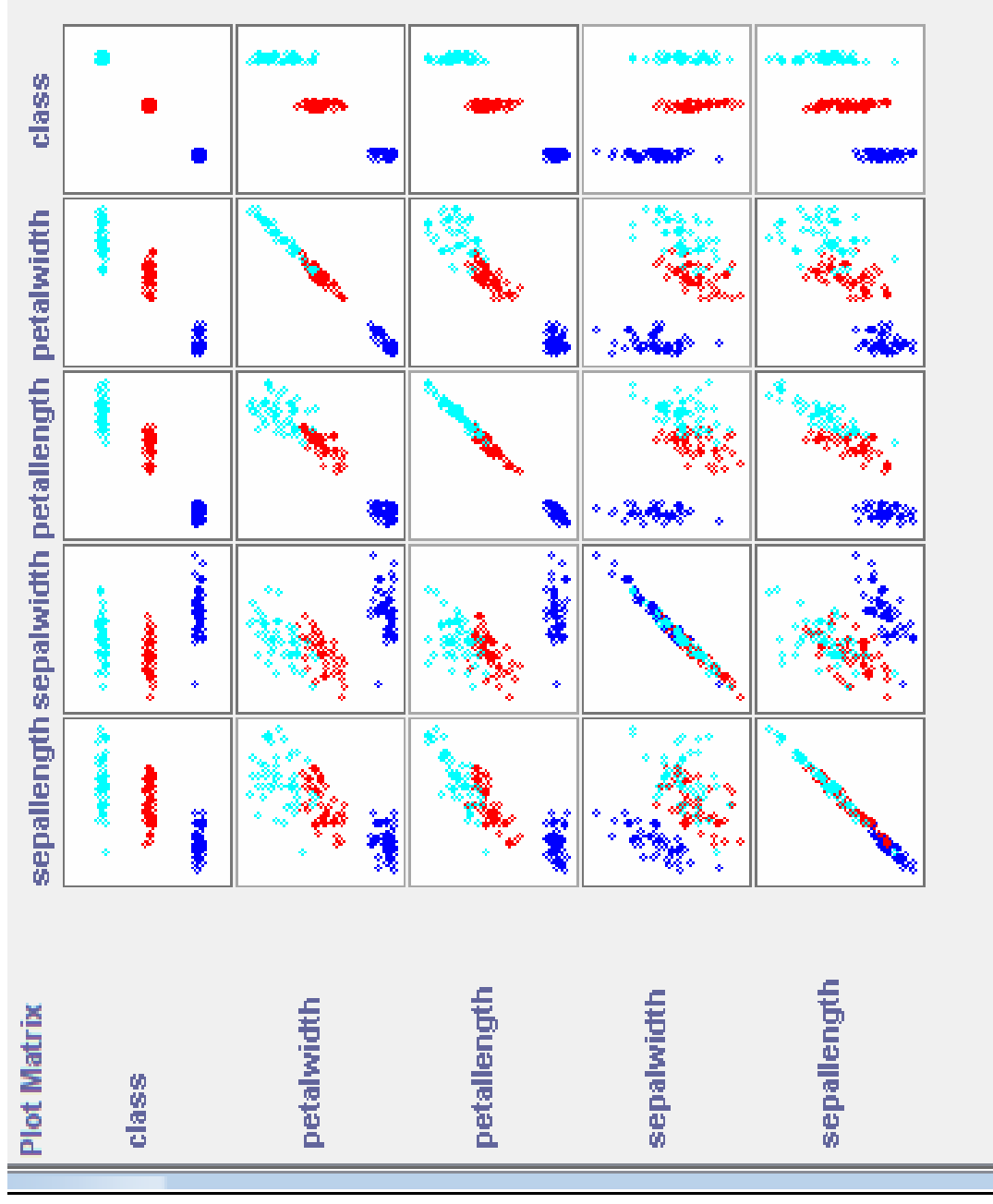
@DATA
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
...
```

- ❑ CSV (Comma Separated Value)
 - tem conversor para este formato e outros formatos populares

Actividade : Iris : compreensão dos dados

- ❑ **O conjunto de dados Iris**
 - **descreve 150 plantas de três tipos diferentes. O que distingue esses 3 tipos?**
- ❑ **Carregar os dados Iris para o Weka**
- ❑ **Observar as distribuições dos valores**
- ❑ **procurar visualmente atributos correlacionados com o target (class / Espécie)**
- ❑ **interpretar o plot em Visualize**
 - **jitter: espalha os pontos**
 - **class colour**
 - **plot size: escala do gráfico**

Actividade : Iris : visualize



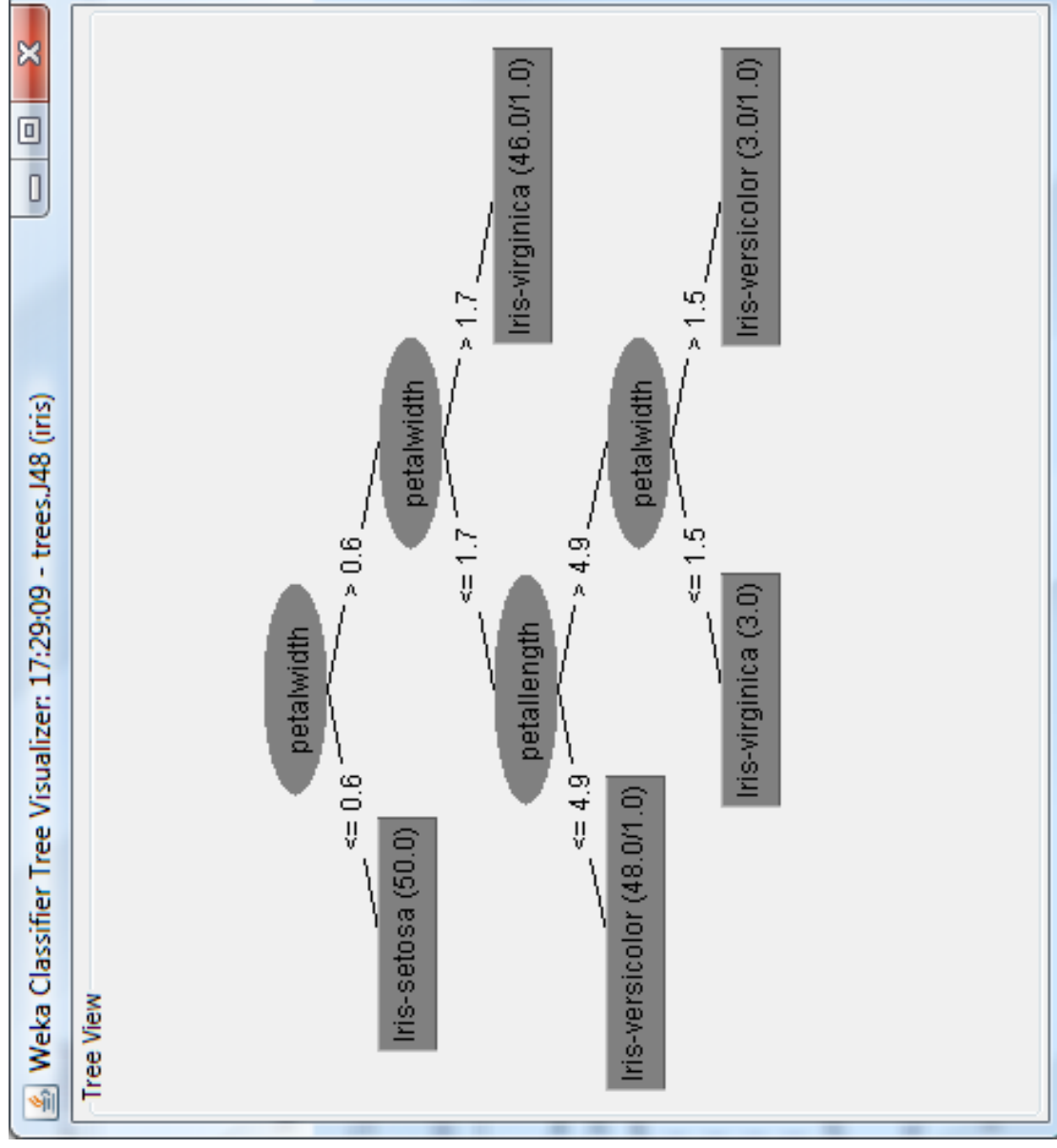
Actividade : Iris : Classificação : OneR

- ❑ **Método OneR (um classificador muito simples)**
 - **correr o método com treino e teste: split de 66%**
 - **qual o modelo obtido?**
 - **quantos respostas erradas no teste?**
 - **qual a percentagem de erro?**
 - **qual a classe mais fácil?**
 - **quais as mais difíceis?**
 - **quais as classes que se confundiram no teste?**
 - **relacione a matriz de confusão com a visualização**
 - **correr com cross-validation: mudou a sua opinião sobre a qualidade do modelo?**

Actividade : Iris : Classificação : J48

- ❑ **Método J48 (árvore de decisão)**
 - **usar cross validation:**
 - **qual a árvore obtida? Como se lê?**
 - **visualize a árvore (botão direito sobre Result list)**
 - **qual a percentagem de erro?**
 - **como se compara com o OneR?**

Visualizar a árvore



Interpretar o output

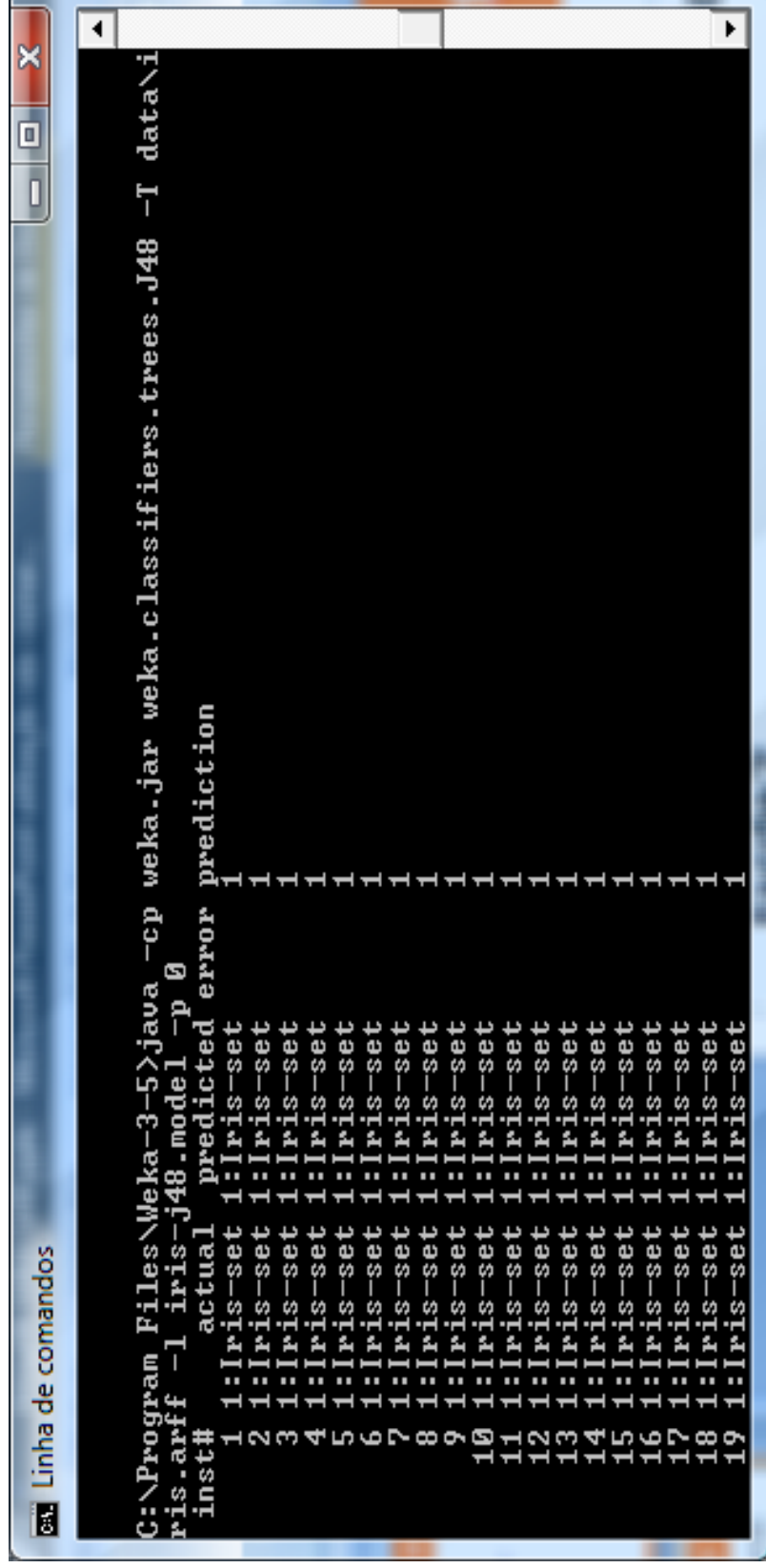
- ❑ **mean absolute error**
 - **SOMA (|prob(prev) – verdade| / N)**
- ❑ **TP rate (true positives) da classe C**
 - **#previstos correctamente da classe C / #exemplos da classe C**
- ❑ **FP rate (false positives) da classe C**
 - **#previstos como da classe C / #exemplos de outras classes**
- ❑ **Precision da classe C**
 - **#previstos correctamente da classe C / #previstos na classe C**
- ❑ **Recal da classe C**
 - **#previstos correctamente da classe C / #exemplos da classe C**
- ❑ **F-measure**
 - **2*Prec*Rec / (Prec+Rec)**

Actividade : Iris : Deployment : J48

- ❑ **Como se faz para guardar o modelo e o usar mais tarde?**
 - com o botão direito do rato sobre o item da **result list**
 - save model (ex. iris-j48.model)
 - o modelo é agora um ficheiro autónomo
- ❑ **Novos dados**
 - são colocados no formato ARFF
- ❑ **Usar o modelo**
 - na linha de comando

```
java -cp weka.jar weka.classifiers.trees.J48 -T  
iris-novos.arff -l iris-j48.model -p 0
```

Actividade : Iris : Deployment : J48



```
C:\Program Files\Meka-3-5>java -cp weka.jar weka.classifiers.trees.J48 -T data\iris.arff -l iris-j48.model -p 0
inst# actual predicted error prediction
1 1:iris-set 1:iris-set 1
2 1:iris-set 1:iris-set 1
3 1:iris-set 1:iris-set 1
4 1:iris-set 1:iris-set 1
5 1:iris-set 1:iris-set 1
6 1:iris-set 1:iris-set 1
7 1:iris-set 1:iris-set 1
8 1:iris-set 1:iris-set 1
9 1:iris-set 1:iris-set 1
10 1:iris-set 1:iris-set 1
11 1:iris-set 1:iris-set 1
12 1:iris-set 1:iris-set 1
13 1:iris-set 1:iris-set 1
14 1:iris-set 1:iris-set 1
15 1:iris-set 1:iris-set 1
16 1:iris-set 1:iris-set 1
17 1:iris-set 1:iris-set 1
18 1:iris-set 1:iris-set 1
19 1:iris-set 1:iris-set 1
```

Actividade : Iris : Deployment : J48

- ❑ De uma forma mais rápida (mas menos flexível)
 - usar Supplied test set, com novos exemplos, classe "?"
 - carregar o modelo (ou aprendê-lo)
 - com o botão direito sobre a Result list, Re-evaluate model on current test set.

Classifier output

Size of the tree : 9

Time taken to build model: 0.01 seconds

=== Evaluation on test split ===
=== Summary ===

Correctly Classified Instances	49	96.0784 %
Incorrectly Classified Instances	2	3.9216 %
Kappa statistic	0.9408	
Mean absolute error	0.0396	
Root mean squared error	0.1579	
Relative absolute error	8.8979 %	
Root relative squared error	33.4051 %	
Total Number of Instances	51	

Accuracy By Class ===

Class	Count	Precision	Recall	F-Measure	ROC Area
Iris-setosa	1	1	1	1	1
Iris-versicolor	63	0.905	1	0.95	0.969
Iris-virginica	1	0.862	0.862	0.838	0.967

Matrix ===

Test options

Use training set

Supplied test set Set...

Cross-validation Folds: 10

Percentage split %: 66

More options...

(Nom) Colest

Start Stop

Result list (right-click for options)

- 17:26:01 - rules.OneR
- 17:29:09 - trees.J48
- 17:31:16 - trees.J48
- 17:31:25 - trees.J48
- 17:31:29 - trees
- 17:31:30 - trees
- 17:31:31 - trees
- 17:31:32 - trees
- 17:31:36 - trees
- 00:19:39 - rules
- 00:20:01 - rules
- 00:24:01 - bayes
- 00:27:16 - meta
- 00:27:43 - meta

View in main window

View in separate window

Save result buffer

Delete result buffer

Load model

Save model

Re-evaluate model on current test set

Actividade : Iris : Classificação : PART

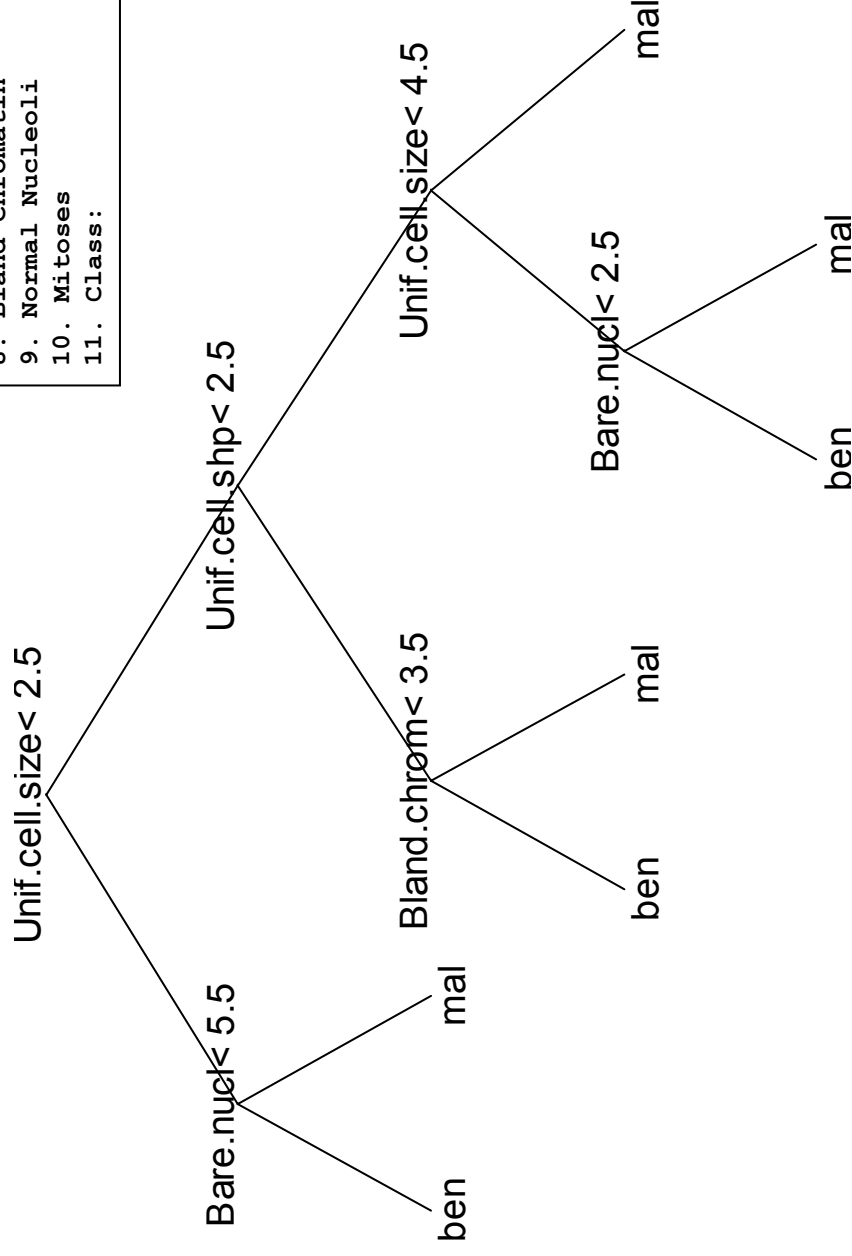
- **PART**
 - obtém regras a partir de árvores de decisão parciais

- **regras**
 - petalwidth \leq 0.6: Iris-setosa (50.0)
 - petalwidth \leq 1.7 AND
 - petalwidth \leq 1.7 AND petalwidth \leq 4.9: Iris-versicolor (48.0/1.0)
 - : Iris-virginica (52.0/3.0)

Classificação: breast cancer > diagnóstico

☐ Breast cytology

– Análise de 699 exames



Attributes

id	number
1.	Sample code number
2.	Clump Thickness
3.	Uniformity of Cell Size
4.	Uniformity of Cell Shape
5.	Marginal Adhesion
6.	Single Epithelial Cell Size
7.	Bare Nuclei
8.	Bland Chromatin
9.	Normal Nucleoli
10.	Mitoses
11.	Class:
	benign, malignant

Actividade : breast cancer

- ❑ **Obtenha e avalie classificadores para os dados**
 - como são os dados? distribuição das classe
 - OneR
 - J48 (veja a árvore)
 - erro?
 - clas
- ❑ **Como se compara com o conjunto de dados Iris?**

Classificação: Caso

- ❑ **Fertilização “in vitro”**
 - **recolha de óvulos na mulher**
 - **fertilização externa -> produção de vários embriões**
 - **selecção de embriões**
 - **colocação de embriões no útero**
- ❑ **O problema**
 - **seleccionar os embriões mais promissores.**
 - **60 variáveis descritivas**
 - **dados históricos**
 - **limites cognitivos do embriologista**

Classificação: Caso

□ Análise:

- **Situação operacional:**
 - recolha de vários óvulos de uma mulher, para fertilização externa e posterior reimplantação no útero da mulher.
- **Problema de decisão:**
 - Quais os embriões a seleccionar (quais os que têm mais hipóteses de sobreviver?)
- **Falta de conhecimento sobre o processo de decisão:**
 - não há teoria (ou é incompleta) sobre como seleccionar os embriões.
- **Necessidade de obtenção do conhecimento explícito:**
 - Validação dos critérios de decisão obtidos automaticamente
- **Necessidade de automação do processo:**
 - demasiadas variáveis e conhecimento relacionado, limites cognitivos do embriologista.
- **Existência de dados:**
 - Registos históricos.
- **Utilização de uma técnica de DM:**
 - classificação.

Classificação

- **Aplicações**
 - **Microarrays**
 - quais os genes importantes para esta doença?
 - **Farmacologia**
 - qual o padrão que corresponde a um químico carcinogénico?
 - **Diagnóstico a partir de imagens (mamografia)**
 - benigno ou maligno?
 - ...

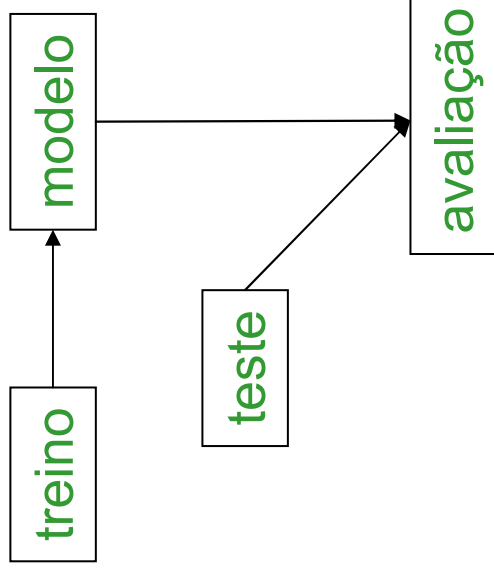
Classificação

□ Avaliação

- **Que confiança podemos ter na resposta do modelo de classificação?**
- **Importa estimar o erro cometido pelo modelo**
 - assumindo que dados futuros mantêm a distribuição presente
- **O que medir?**
 - Error rate: percentagem de respostas erradas
 - Accuracy / acerto: percentagem de respostas certas
 - Recall: percentagem de casos de uma classe identificados
 - Precision: acerto numa classe
- **Esse estimador deve ser fiável**
 - pode ser medido de várias formas

Classificação: Avaliação

- ❑ **Exemplos de treino**
 - para construir o modelo
 - *training-set*
- ❑ **Exemplos de teste**
 - para avaliar o modelo
 - *test-set*
- ❑ **Distribuição dos exemplos**
 - geralmente assume-se dist. treino = dist. teste
 - *concept drift* (deriva conceptual)
- ❑ **Avaliação mais robusta**
 - validação cruzada
 - *leave-one-out*



Classificação

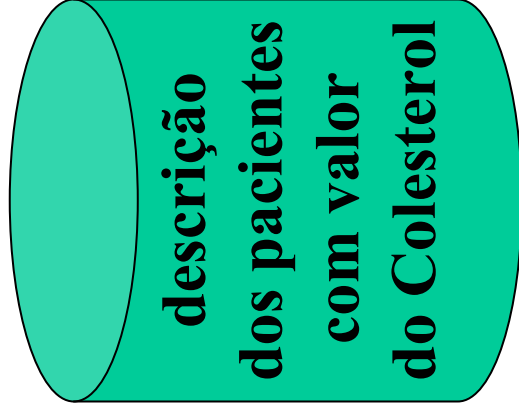
- **Técnicas**
 - **Árvores de decisão**
 - **Sistemas de regras**
 - **Support Vector Machines (SVM)**
 - **Redes Neurais (NN)**
 - **Discriminante linear**
 - **Regressão logística**
 - **Abordagens Bayesianas**
 - **Baseadas em instâncias**
 - ...

Actividade : Comparação

- **Com o dataset Colestlab**
 - **comparar métodos**
 - OneR
 - J48
 - Bagging com OneR
 - Bagging com J48
 - **Sugestão: usar o experimenter do Weka**

Problemas de Regressão

- **Dados**
 - N casos descritos por M variáveis, sendo uma delas a var. objectivo
- **Descobrir**
 - uma relação funcional $f(X) = y$
 - que a um novo caso X, faça corresponder um novo valor y



$$\text{CHOL} = a_1 * \text{IDA} + a_2 * \text{EST} \dots$$

Problemas de Regressão

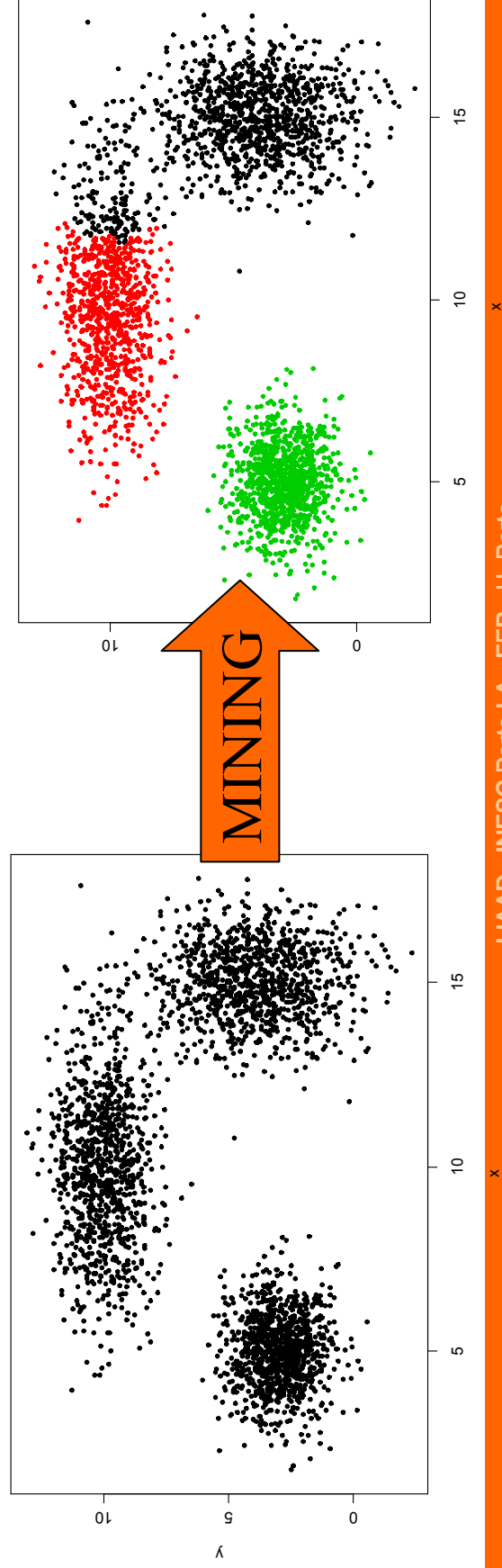
- **Aplicações**
 - Previsão numérica
 - Séries temporais
 - Planeamento logístico
- **Avaliação**
 - erro: distância “média” entre previsão e verdade
 - treino e teste, validação cruzada, etc.
- **Técnicas**
 - Regressão linear
 - Árvores de regressão
 - Support vector machines
 - ...

DM: dirigido vs. exploratório

- ❑ **Dirigido (Classificação, Regressão,...)**
 - há um objectivo bem definido
 - é facilmente avaliável
 - presta-se à automação das decisões
- ❑ **Exploratório (Clustering, Associação, Visual DM,...)**
 - procuram-se padrões interessantes, accionáveis
 - é difícil de avaliar
 - o perito é muito importante no ciclo de descoberta
 - a visualização é importante

Problemas de Clustering

- ❑ **Dados**
 - N casos, sem classe conhecida
- ❑ **Descobrir**
 - Uma partição dos casos que maximize a semelhança intra-grupo e a dissemelhança entre grupos distintos
- ❑ **Aplicações**
 - descobrir grupos de genes / pacientes / compostos químicos



Clustering: breast cancer

- ❑ **Vamos tentar redescobrir a classe**
 - Usamos clustering hierárquico
 - Pedimos 2 clusters
 - Fornecemos a tabela com todos os atributos menos Diag
 - Confrontamos com os verdadeiros valores de classe
 - Sobreposição de 91% (avaliação?)

	cluster	ben	mal
1		452	50
2		6	191

Actividade : Weka : Clustering : Breast Cancer

- ❑ **Carregue os dados Breast Cancer**
 - use o **SimpleKMeans**
 - use **class** para avaliar
 - veja os resultados

Clustering: breast cancer > visualização

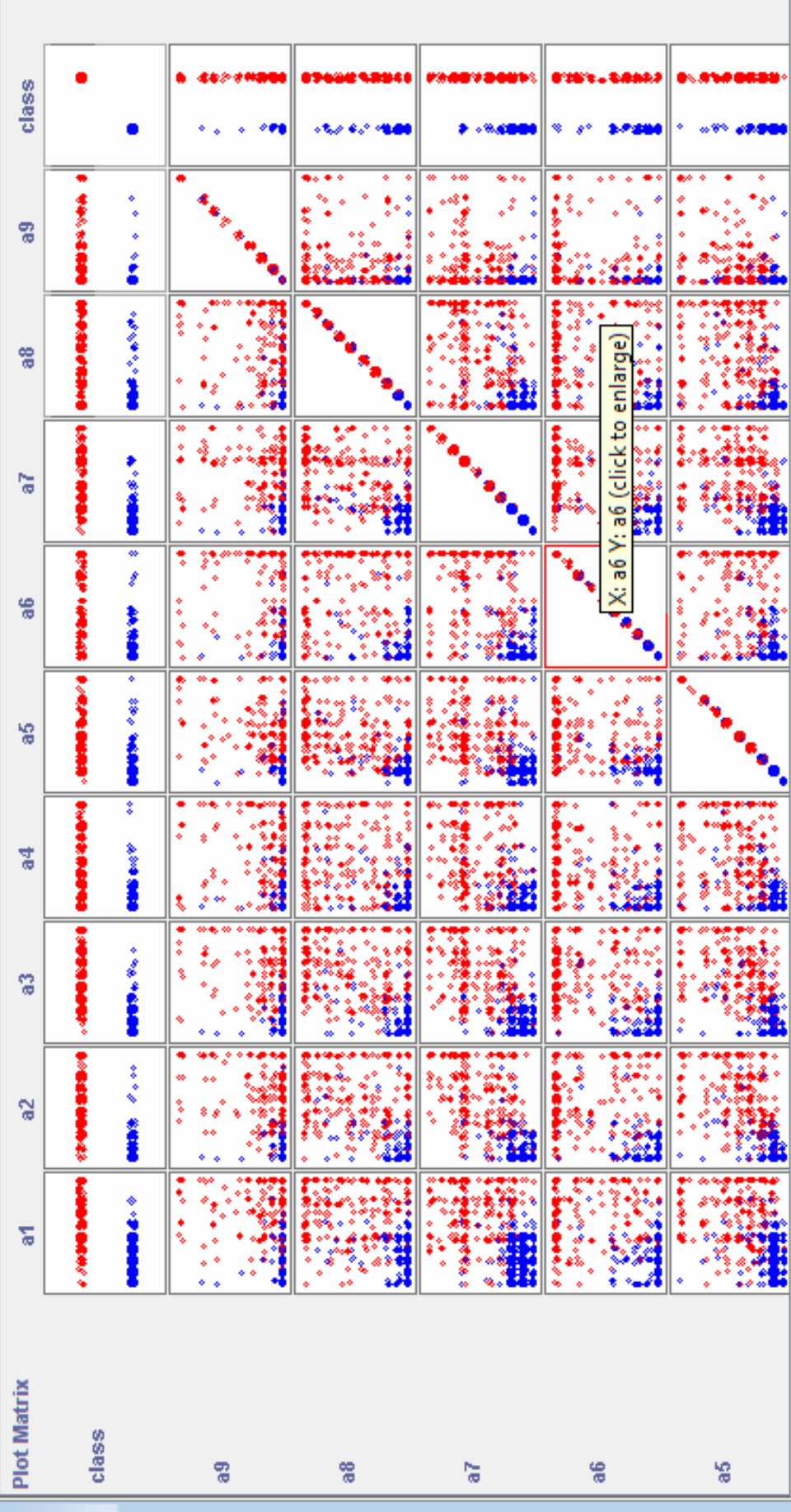
Weka 3.5.6 - Explorer

Program Applications Tools Visualization Windows Help

Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Visualize



Problemas de Associação

- **Dados**
 - N casos descritos por M variáveis
- **Descobrir**
 - Co-ocorrências frequentes das variáveis
- **Aplicações**
 - Prognóstico / prevenção / planeamento
 - Descoberta de interacção entre factores
- **Técnicas**
 - Regras de associação
 - ...

Regras de Associação

Rules	Support	Confidence
IMC=obesidade, DESP=nao, ALC=bebe -> Colest=alto	0,03714	0,5
IMC=obesidade, DESP=nao, CIV=c, ALC=bebe -> Colest=alto	0,03429	0,52174
IMC=obesidade, DESP=nao, CIV=c, CAF=sim, ALC=bebe -> Colest=alto	0,03333	0,51471
DESP=nao, CAF=sim, IDA=60anos, TAB=nao -> Colest=alto	0,02952	0,51667
CAF=sim, IDA=60anos, TAB=nao, ALC=bebe -> Colest=alto	0,02476	0,53061
Peso=80a70kg, ALT=m150, DESP=nao, CIV=c -> Colest=alto	0,02381	0,5102
HDORM=8ha10h, ESCOL=4Classe, DESP=nao, CIV=c, CAF=sim, TAB=nao -> Colest=alto	0,02381	0,5
HDORM=8ha10h, ESCOL=4Classe, ALT=m150, CIV=c, CAF=sim -> Colest=alto	0,02286	0,6
HDORM=8ha10h, ESCOL=4Classe, ALT=m150, CIV=c -> Colest=alto	0,02286	0,55814
ACTIV=pouca, DESP=nao, IDA=60anos -> Colest=alto	0,02286	0,53333
CIV=c, IDA=60anos, TAB=nao, ALC=bebe -> Colest=alto	0,02286	0,53333
IMC=obesidade, DESP=nao, TAB=nao, ALC=bebe -> Colest=alto	0,02286	0,51064
HDORM=8ha10h, IMC=excessopeso, ALT=m150 -> Colest=alto	0,02286	0,51064
CIV=c, CAF=sim, IDA=60anos, TAB=nao, ALC=bebe -> Colest=alto	0,0219	0,62162
HDORM=8ha10h, CAF=sim, IDA=60anos, TAB=nao -> Colest=alto	0,0219	0,52273
Peso=80a70kg, sexo=f, ALT=m150, DESP=nao, CIV=c -> Colest=alto	0,0219	0,52273
HDORM=8ha10h, sexo=f, ESCOL=4Classe, CIV=c, CAF=sim -> Colest=alto	0,0219	0,52273
sexo=f, DESP=nao, CAF=sim, IDA=60anos, TAB=nao -> Colest=alto	0,0219	0,51111

Regras de Associação

□ Simbolicamente

$$\{ A1, A2 \} \Rightarrow \{ A3 \}$$

- Se observarmos um conjunto de itens A1 e A2 devemos também observar o item A3.
- suporte
 - percentagem dos cestos onde a co-ocorrência se observa
 - Estima $\text{Prob}(A1 \& A2 \& A3)$
- confiança
 - percentagem dos casos onde a ocorrência de $\{A1, A2\}$ “prevê” correctamente a ocorrência de $\{A3\}$
 - Estima $\text{Prob}(A3 \mid A1 \& A2)$

Em geral

- dado um conjunto de transacções D
- $D = \{t \mid t \text{ é um conjunto de items } i\}$

$D = \{t \mid t \text{ é um conjunto de items } i\}$

- uma regra de associação tem a forma $A \rightarrow B$
- A e B são conjuntos de items

$$\text{suporte}(A \rightarrow B) = \frac{\#(A \cup B)}{\#D}$$

$$\text{confiança}(A \rightarrow B) = \frac{\#(A \cup B)}{\#A}$$

Interesse de uma regra

- ❑ **Regra interessante** [Silberschatz & Tuzhilin]
 - **Inesperada: surpreendente para o utilizador**
 - medida de interesse: *desvio do esperado ou do acreditado*
 - **Útil (accionável)**
 - medida de interesse: *benefício estimado*
- ❑ **Subjectivo**
 - **interesse depende do conhecimento do utilizador**
 - um utilizador pode identificar uma oportunidade de aplicação
 - **envolve itens de interesse**
 - um utilizador pode estar mais interessado num determinado domínio
- ❑ **Objectivo**
 - **desvio da independência estatística**
 - **valores que se destacam**
 - e.g.: regras sintacticamente singulares

Regras interessantes

□ Tipicamente

– $A \rightarrow B$ é interessante se A e B não são estatisticamente independentes

– se A e B são estatisticamente independentes

$$\text{sup}(A \cup B) \approx \text{sup}(A) \times \text{sup}(B)$$

$$\text{conf}(A \rightarrow B) \approx \text{conf}(\emptyset \rightarrow B)$$

– $A \rightarrow B$ pode ter suporte elevado, confiança elevada e não ser interessante.

{ jornal } \rightarrow { combustível }

sup = 5 %

conf = 95 %

não é inesperada

não é útil

Regras interessantes

- ❑ O desvio da independência
 - Regra inesperada \approx que se desvia do pressuposto de independência
- ❑ Medidas de interesse
 - **quociente** $\frac{\text{probabilidade a posteriori}}{\text{probabilidade a priori}} > 1$
 - **diferença** $\text{probabilidade a posteriori} - \text{probabilidade a priori} > 0$ *leverage*
 - **teste de hipóteses** $\text{probabilidade a posteriori} \gg \text{probabilidade a priori?}$ χ^2

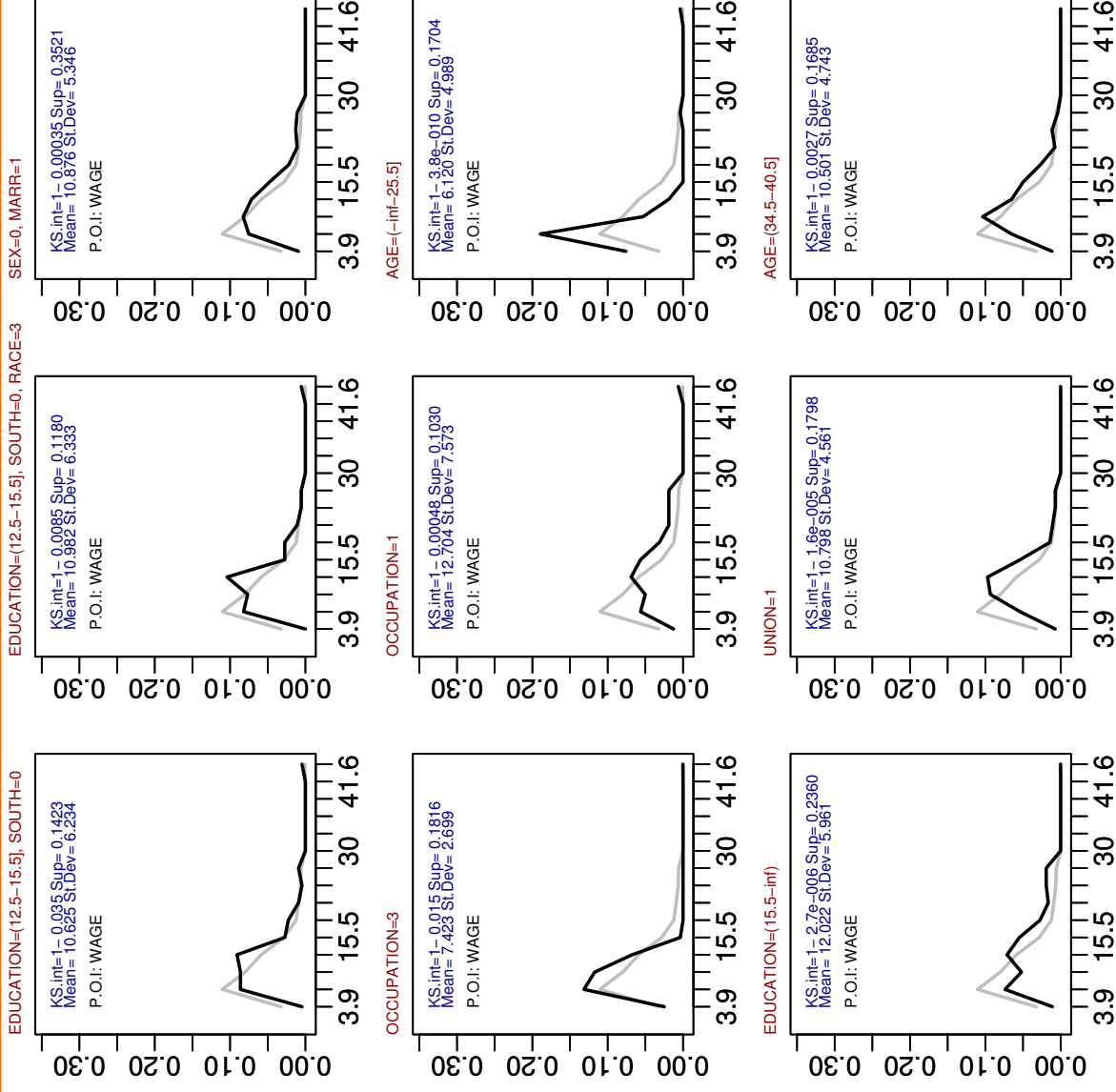
$$\text{lift}(A \rightarrow B) = \frac{\text{confiança}(A \rightarrow B)}{\text{suporte}(B)}$$

$$\text{conviction}(A \rightarrow B) = \frac{\text{suporte}(\neg B)}{\text{lift}(A \rightarrow \neg B)}$$

Actividade : Colestlab : Associação

- ❑ **Use o Weka para obter regras de associação**
 - **quais as regras mais interessantes?**

Regras de distribuição (consequente numérico)



Outros problemas/técnicas de DM

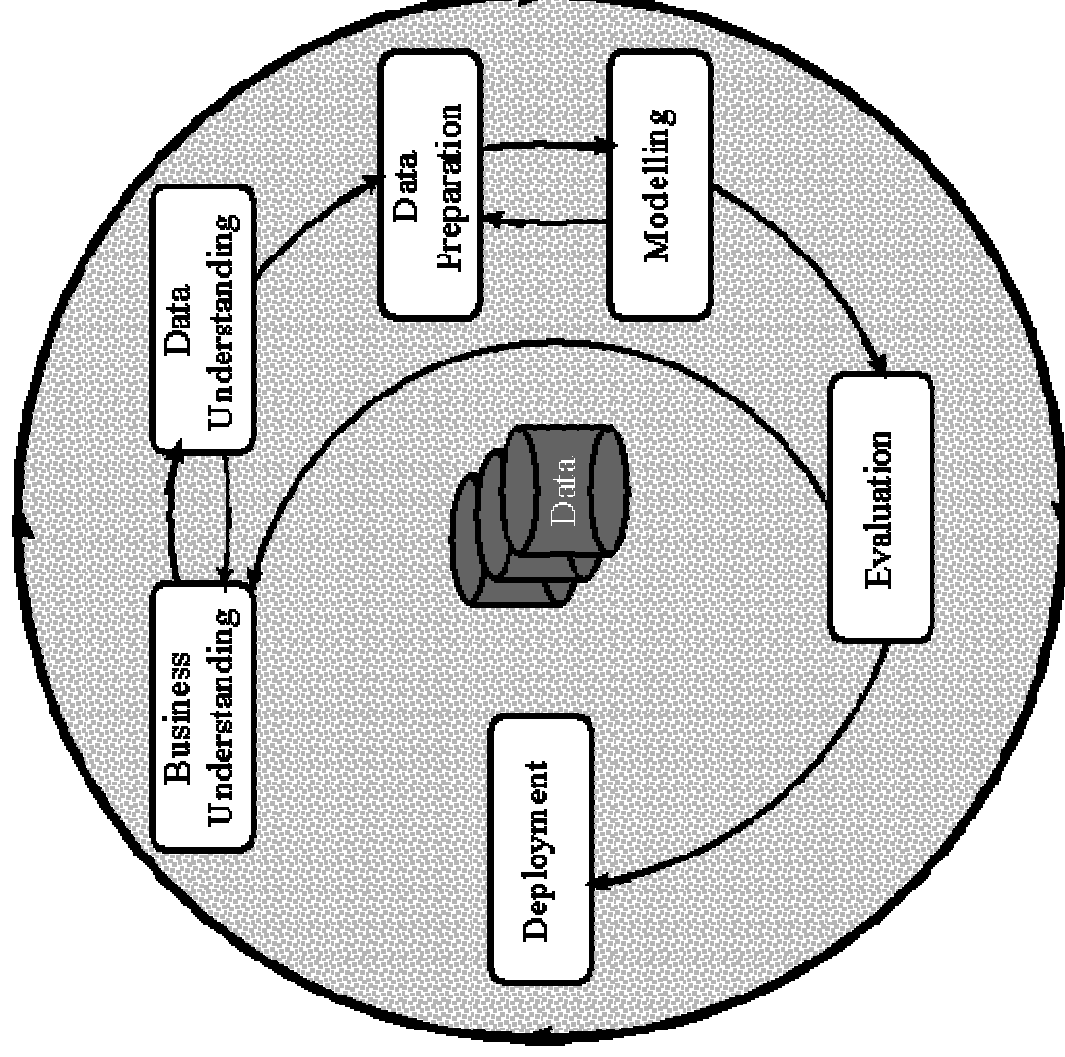
- ❑ **Detecção de outliers**
 - Previsão de fenómenos raros
- ❑ **Text e Web mining**
 - Sumariação automática de textos
 - Extracção de informação
- ❑ **Data mining multi-relacional**
 - Lidar com estruturas complexas (moléculas)
 - Descobrir padrões que explicam actividade
- ❑ **Pré-processamento dos dados**
 - Muito importante
 - DNA micro array pode ter milhares de variáveis

O Processo de Data Mining

- ❑ **Quais são os passos de um processo de DM?**
 - como chegar do problema à acção?
- ❑ **Existe alguma metodologia/ processo modelo de DM?**
 - pronta a ser seguida ?

O processo de Data Mining (versão 4)

CRISP-DM (NCR, SPSS, Daimler-Chrysler, OHRA) (1997)



O (modelo de) processo CRISP-DM

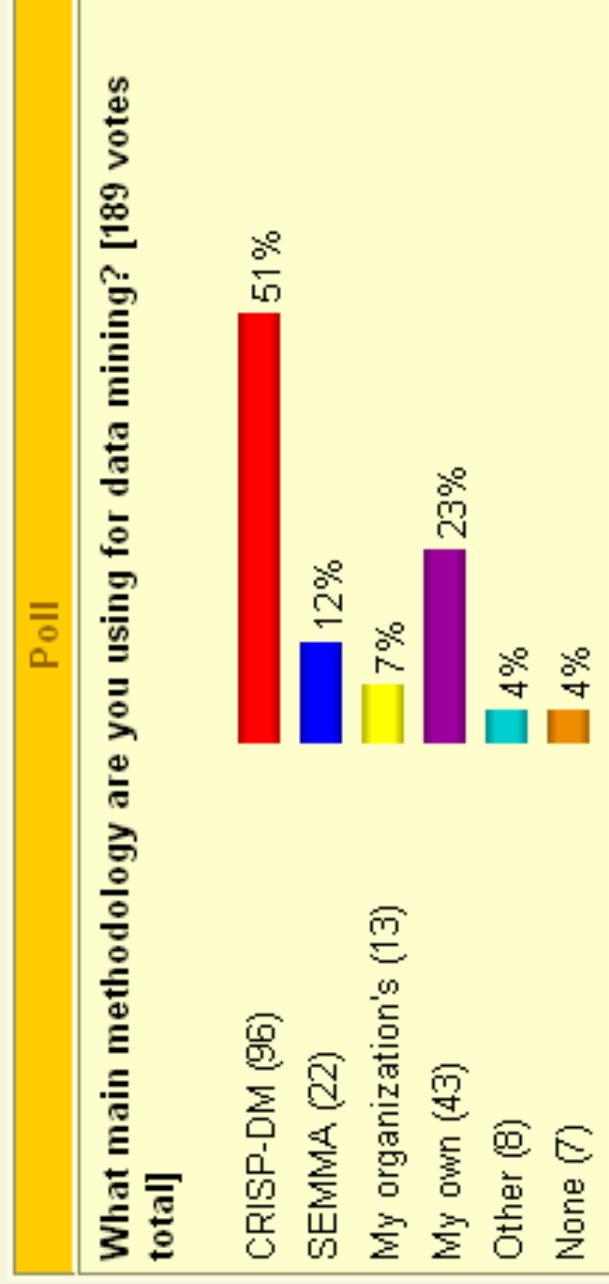


- ❑ Pretende ser um **standard**
 - **Cross-Industry Standard Process for Data Mining**
- ❑ Desenvolvido com base na **experiência de muitos profissionais**
- ❑ A partir do **Modelo de Processo CRISP-DM** podemos definir um **Processo de DM**

Fases do ciclo de vida de um projecto de DM

- ❑ **Compreensão da Actividade (*business understanding*)**
 - Os objectivos do projecto do ponto de vista da actividade
- ❑ **Compreensão dos dados**
 - Familiarização com os dados, qualidade dos dados
- ❑ **Preparação dos dados**
 - Construir o *data set* a utilizar
- ❑ **Modelação**
 - Experimentar e seleccionar técnicas de modelação
- ❑ **Avaliação**
 - Rever os passos seguidos para construir o modelo
- ❑ **Acção / Produção (Deployment)**
 - Utilização dos modelos criados em problemas reais

KDnuggets : Polls : What main methodology are you using for data mining?



Comments

- ◆ Uta Winter, Subject: CRISP-DM
As the ultimate reason of doing Data Mining is always to improve a business situation which is perceived as being not satisfactory it is only logical that the Crisp-DM methodology starts there: How do we tackle the business problem by transforming it into a DM-problem. It doesn't start with sampling (which is not always necessary anyway)! This is a big improvement over more technically focused approaches. The Crisp-Manual is well-written and especially non-tekkies can benefit from it. After having done some projects myself I can feel that this process model is not written by people who have never been out in a company and only discussing DM in theory but who exactly know all the issues which occur when doing a real project.

Referências

- ❑ www.kdnuggets.com
- ❑ Usama M. Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth: The KDD Process for Extracting Useful Knowledge from Volumes of Data. Commun. ACM 39(11): 27-34 (1996)
- ❑ David Hand, Heikki Mannila, Padhraic Smyth, Principles of Data Mining, MIT Press, 2001
- ❑ Igor Jurisica, John Mylopoulos, Janice I. Glasgow, Heather Shapiro, Robert F. Casper: Case-based reasoning in IVF: prediction and knowledge mining. Artificial Intelligence in Medicine 12(1): 1-24 (1998)
- ❑ Shearer, C., The CRISP-DM Model: The New Blueprint for Data Mining. In Journal Data Warehousing, Vol. 5, No. 4, (2000), 13—22
- ❑ Kim, J., Zhao, S., and Heber, S. 2007. Finding association rules of cis-regulatory elements involved in alternative splicing. In Proc. of the 45th Annual Southeast Regional Conference (2007). ACM.
- ❑ www.r-project.org
- ❑ www.liaad.up.pt/~amjorge/Projectos/Class/