

**\*TÍTULO DO TRABALHO DE PROJECTO\***

**\*Fernando Pessoa\***

**\*Ricardo Reis\***

Licenciatura em Engenharia Informática e de Computadores  
Projecto e Seminário

Orientadores: **\*Álvaro de Campos\***

**\*Alberto Caeiro, SoftCompany\***

Apresentação \* \*

Maio de 2015

# Sumário

Introdução

O Problema

A Solução

Grande Ideia 1

Grande Ideia 2

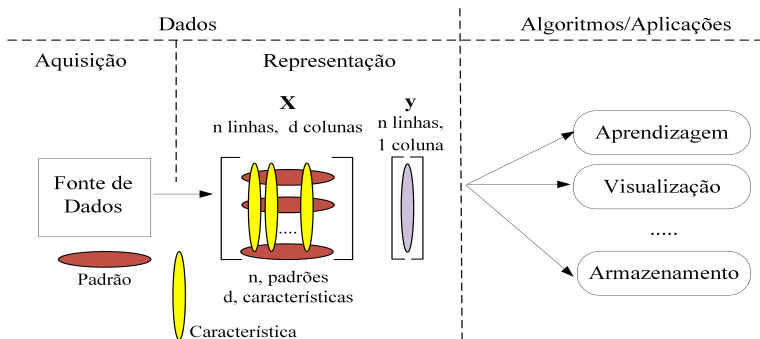
Resultados Obtidos

Conclusões

Trabalho Futuro

# Introdução

## Cenários típicos de manipulação de dados (figura)



Os dados (lista enumerada):

- organizados em  $n$  padrões, exemplos ou instâncias
- cada padrão tem  $d$  características (dimensionalidade)
- podem ser manipulados com diferentes propósitos/objetivos

## O Problema: conceitos

Exemplos de dados de baixa dimensionalidade e com AD (tabela):

- $d$  dimensões,  $c$  classes,  $n$  padrões
- em baixa dimensionalidade,  $n > d$  ou  $n \gg d$
- com AD, tipicamente tem-se  $d \gg n$

Conjunto	$d$	$c$	$n$	Tipo de Dados	Problema
Car	6	4	1728	Automóveis	Possibilidade de Compra
Pima	8	2	768	Clínicos	Deteção da Diabetes
Contraceptive	9	2	1473	Clínicos	Método de Contraceção
Wine	13	3	178	Químicos	Tipo de Vinho
Hepatitis	19	2	155	Clínicos	Deteção de Hepatite
Dermatology	34	6	358	Clínicos	Doença de Pele
Colon	2000	2	62	Exp. Genética	Deteção de Cancro
SRBCT	2309	4	83	Exp. Genética	Deteção de Cancro
TOX-171	5748	4	171	Exp. Genética	Deteção de Cancro
Example1	9947	2	2000	Texto	Assunto da Notícia
ORL10P	10304	10	100	Faces	Identificação
11-Tumors	12553	11	174	Exp. Genética	Deteção de Cancro
Lung-Cancer	12601	5	203	Exp. Genética	Deteção de Cancro
Dexter	20000	2	2600	Texto	Assunto da Notícia
GLI-85	22283	2	85	Exp. Genética	Deteção de Cancro
Dorothea	1000000	2	1950	Clínicos	Deteção de Composto

## Propostas para Seleção: Complexidade

Em termos de complexidade, tem-se para o RFS

$$C_{RFS} = \underbrace{O(nd)}_{\text{Relevancia}} + \underbrace{O(d \log d)}_{\text{Ordenacao}}$$

RRFS apresenta a complexidade adicional de calcular as semelhanças entre pares

$$C_{RRFS} = \underbrace{O(nd)}_{\text{Relevancia}} + \underbrace{O(d \log d)}_{\text{Ordenacao}} + \underbrace{O(nm)}_{\text{Redundancia}}$$

- RRFS é mais rápido do que outros filtros (FCBF e MRMR)
- Nalguns casos, RRFS é o mais rápido com menos erro
- CFS é o mais lento (inadequado para dados com AD)
- Medidas MAD e MM adequadas para todos os tipos de dados
- Medida AMGM é mais adequada para dados esparsos

## Representação de Dados

Assim, identifica-se a necessidade de:

- encontrar novas formas de representação dos dados
- representar de forma independente da tarefa a jusante
- facilitar a visualização e análise de dados com AD

As técnicas de **seleção** (*feature selection* - FS) e **discretização** (*feature discretization* - FD) realizam essa representação

- **Seleção**

→ escolha de sub-conjuntos de características adequados

- **Discretização**

→ representações discretas de características numéricas  
→ com informação suficiente para aprendizagem  
→ ignora ruído e flutuações irrelevantes

## Seleção: Taxonomia e Opções Tomadas

Algoritmos de seleção são categorizados como:

- i) filtros (*filters*)
- ii) envolvimento (*wrapper*)
- iii) embebidos (*embedded*)
- iv) híbridos (*hybrid*)

Escolha inicial → filtros não supervisionados e supervisionados

- Baixa complexidade, eficiência e interpretabilidade
- Independentes da tarefa de mineração de dados
- Alguns filtros existentes são:
  - ineficientes (tempo) em dados com  $AD^1$
  - sensíveis ao problema de elevado  $d$ , baixo  $n^2$

---

<sup>1</sup>CFS-Correlation-based Feature Selection, MRMR-Maximum Relevance Minimum Redundancy, RELIEF-Recursive Elimination of Features

<sup>2</sup>FCBF - Fast Correlation-Based Filter





## Avaliação Experimental: Relevância

Algoritmo *relevance FS* (RFS):

- guarda as características com maior *relevância*
- relevância medida pela *dispersão*

Para o caso não supervisionado  $\rightarrow$  medidas de *relevância @rel*:

- *Mean absolute difference*  $MAD_i = \frac{1}{n} \sum_{j=1}^n |X_{ij} - \bar{X}_i|$
- *Mean-median*  $MM_i = |\bar{X}_i - \text{median}(X_i)|$
- *Arithmetic mean geometric mean*  
$$AMGM_i = \frac{1}{n} \sum_{j=1}^n \exp(X_{ij}) / \left( \exp \left( \sum_{j=1}^n X_{ij} \right) \right)^{\frac{1}{n}}$$

Critério de relevância cumulativa  $\sum_{f=1}^m r_{if} / \sum_{i=1}^d r_i = c_m / c_d \geq L$ ,  
 $L \in [0, 8; 0, 95] \rightarrow$  escolha do número de características  $m (< d)$

## Propostas para Seleção: Resultados Experimentais

- Classificação supervisionada (SVM linear), validação cruzada (10-fold)
- Percentagem de erro de generalização com filtros **supervisionados**

	RRFS, $M_S = 0.8$			Filtros Supervisionados <sup>3</sup>					Base
Conjunto	MM	FiR	MI	RF	CFS	FCBF	FiR	MRMR	Sem FS
Colon	24.2	22.6	24.2	<b>19.4</b>	25.8	22.6	<b>19.4</b>	21.0	21.0
Lymphoma	<b>2.2</b>	<b>2.2</b>	<b>2.2</b>	<b>2.2</b>	N/A	3.3	<b>2.2</b>	22.8	<b>2.2</b>
Leukemia1	5.6	<b>2.8</b>	6.9	6.9	N/A	5.6	4.2	9.7	5.6
B-Tumor1	13.3	12.2	13.3	11.1	N/A	<b>18.9</b>	11.1	<b>25.6</b>	<b>10.0</b>
Leukemia	<b>2.8</b>	12.5	<b>2.8</b>	<b>2.8</b>	N/A	4.2	4.2	8.3	<b>2.8</b>
Example1	2.3	2.2	2.2	3.7	N/A	6.3	<b>2.1</b>	28.3	2.4
B-Tumor2	34.0	<b>22.0</b>	30.0	<b>22.0</b>	N/A	<b>36.0</b>	24.0	<b>42.0</b>	26.0
P-Tumor	<b>7.8</b>	<b>5.9</b>	<b>4.9</b>	<b>7.8</b>	N/A	<b>9.8</b>	<b>7.8</b>	<b>12.7</b>	<b>8.8</b>
L-Cancer	5.9	6.4	4.9	4.9	N/A	6.4	5.4	11.8	5.9
Dexter	6.7	<b>6.0</b>	7.7	9.3	N/A	15.3	6.7	18.0	6.3
Dorothea	<b>25.0</b>	26.0	<b>25.0</b>	N/A	N/A	N/A	<b>25.0</b>	N/A	<b>25.0</b>

<sup>3</sup>Mutual Information (MI), RELIEF (RF), Correlation-based Feature Selection (CFS), Fast Correlation-based Filter (FCBF), Fishers Ratio (FiR) e Maximum Relevance Minimum Redundancy (MRMR)

# Conclusões

- O projeto consistiu em....
- Atingiu-se uma solução ...
- Os métodos propostos foram avaliados:
  - sobre dados de domínio público
  - comparativamente com métodos existentes em implementações de domínio público
- Os métodos propostos complementam os existentes, podendo ser combinados entre si

# Trabalho Futuro

Perspetivam-se as direções de trabalho futuro:

- Melhorar a proposta de ...
- Explorar o afinamento dos parâmetros ...
- ..

Slides elaborados em  $\text{\LaTeX}$  com o package BEAMER.