

# Mineração de Dados em Biologia Molecular

## Árvores de Decisão

André C. P. L. F. de Carvalho  
Monitor: Valéria Carvalho



## Principais tópicos

- Introdução
- Algoritmo de Hunt
- Medidas para escolha de atributos
- Ponto de referência
- Critério de parada
- Espaço de hipóteses

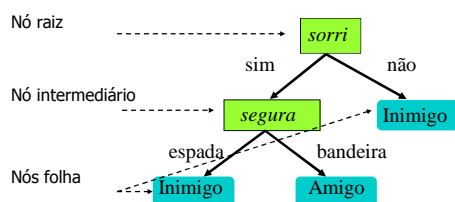
25/10/2012

André de Carvalho - ICMC/USP

2

## Introdução

- Alguns algoritmos de AM induzem ADs a partir de um conjunto de dados

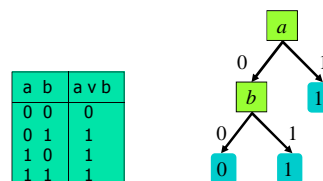


25/10/2012

André de Carvalho - ICMC/USP

3

## Outro exemplo simples



Nós internos e raiz: atributos preditivos  
Nós externos (folhas): classes

25/10/2012

André de Carvalho - ICMC/USP

4

## Exercício

- Encontrar árvore de decisão para:
  - a AND b
  - a XOR b
  - (a AND b) OR (b AND c)

25/10/2012

André de Carvalho - ICMC/USP

5

## Indução de AD

- Existem vários algoritmos
  - Algoritmo de Hunt
    - Um dos primeiros
    - Base de vários algoritmos atuais
  - CART
  - ID3, C4.5
  - SLIQ, SPRINT

25/10/2012

André de Carvalho - ICMC/USP

6

## Algoritmo de Hunt

- Seja  $D_t$  o conjunto de objetos de treinamento que atingem o nó  $t$

Se todos os objetos de  $D_t \in$  a mesma classe  $y_t$   
 Então  $t$  é um nó folha rotulado como  $y_t$   
 Se os objetos de  $D_t \in$  a mais de uma classe  
 Então Selecionar um atributo teste para dividi-los  
 Dividir exemplos em subconjuntos pelo atributo teste  
 Aplicar algoritmo a cada subconjunto gerado

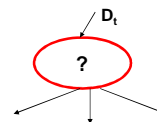
25/10/2012

André de Carvalho - ICMC/USP

7

## Algoritmo de Hunt

Emprego	Estado	Renda	Classe
Sim	Solteiro	9500	Não
Não	Casado	8000	Não
Não	Solteiro	7000	Não
Sim	Casado	12000	Sim
Não	Divorciado	9000	Sim
Não	Casado	6000	Não
Sim	Divorciado	4000	Não
Não	Solteiro	8500	Sim
Não	Casado	7500	Não
Não	Solteiro	9000	Sim



25/10/2012

André de Carvalho - ICMC/USP

8

## Algoritmo de Hunt

Emprego	Estado	Renda	Classe
Sim	Solteiro	9500	Não
Não	Casado	8000	Não
Não	Solteiro	7000	Não
Sim	Casado	12000	Não
Não	Divorciado	9000	Sim
Não	Casado	6000	Não
Sim	Divorciado	4000	Não
Não	Solteiro	8500	Sim
Não	Casado	7500	Não
Não	Solteiro	9000	Sim



Classe default

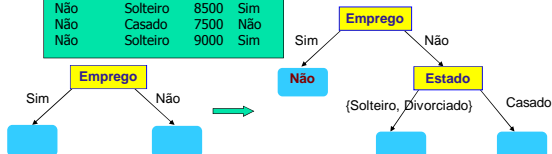
25/10/2012

André de Carvalho - ICMC/USP

9

## Algoritmo de Hunt

Emprego	Estado	Renda	Classe
Sim	Solteiro	9500	Não
Não	Casado	8000	Não
Não	Solteiro	7000	Não
Sim	Casado	12000	Não
Não	Divorciado	9000	Sim
Não	Casado	6000	Não
Sim	Divorciado	4000	Não
Não	Solteiro	8500	Sim
Não	Casado	7500	Não
Não	Solteiro	9000	Sim



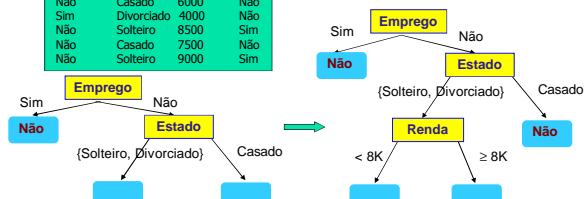
25/10/2012

André de Carvalho - ICMC/USP

10

## Algoritmo de Hunt

Emprego	Estado	Renda	Classe
Sim	Solteiro	9500	Não
Não	Casado	8000	Não
Não	Solteiro	7000	Não
Sim	Casado	12000	Não
Não	Divorciado	9000	Sim
Não	Casado	6000	Não
Sim	Divorciado	4000	Não
Não	Solteiro	8500	Sim
Não	Casado	7500	Não
Não	Solteiro	9000	Sim



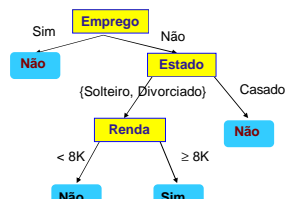
25/10/2012

André de Carvalho - ICMC/USP

11

## Algoritmo de Hunt

Emprego	Estado	Renda	Classe
Sim	Solteiro	9500	Não
Não	Casado	8000	Não
Não	Solteiro	7000	Não
Sim	Casado	12000	Não
Não	Divorciado	9000	Sim
Não	Casado	6000	Não
Sim	Divorciado	4000	Não
Não	Solteiro	8500	Sim
Não	Casado	7500	Não
Não	Solteiro	9000	Sim



André de Carvalho - ICMC/USP

12

## Indução de ADs

- Geralmente usa estratégia gulosa de divisão e conquista
  - Divide progressivamente objetos baseado em um atributo de teste
    - Escolhido para otimizar algum critério
- Decisões importantes
  - Como dividir os objetos?
    - Método para escolha do atributo de teste
  - Quando parar de dividir os objetos?

25/10/2012

André de Carvalho - ICMC/USP

13

## Como dividir os objetos?

- Depende do tipo do atributo
  - Binário
  - Simbólico (mais que dois valores)
    - Nominal
    - Ordinal
  - Numérico (discreto ou contínuo)
- Depende do número de divisões
  - 2 divisões
  - Mais que 2 divisões

25/10/2012

André de Carvalho - ICMC/USP

14

## Atributos binários

- Teste mais simples
  - Apenas dois possíveis resultados



25/10/2012

André de Carvalho - ICMC/USP

15

## Atributos simbólicos

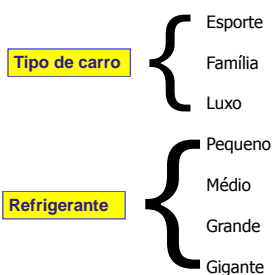
- Duas formas de condição de teste
  - Fazer #ramos = #possíveis valores
  - Agrupar parte dos valores em cada ramo
    - Ordinais:
      - Valores agrupados não devem violar relação de ordem
    - Nominais:
      - Objetos em um nó filho podem estar associados a um grupo de valores

25/10/2012

André de Carvalho - ICMC/USP

16

## Atributos simbólicos

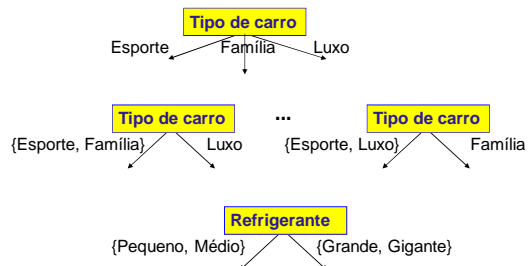


25/10/2012

André de Carvalho - ICMC/USP

17

## Atributos simbólicos



25/10/2012

André de Carvalho - ICMC/USP

18

## Atributos discretos ou contínuos

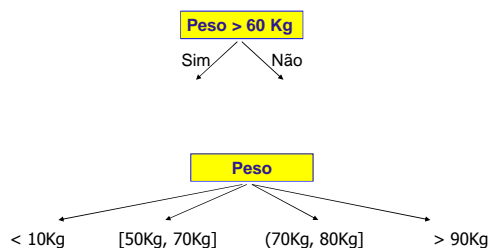
- Condição de teste pode ser expressa por:
  - Comparação simples ( $A < \text{valor}$ )
    - Escolher posição (valor) que gera melhor partição
      - Ponto de referência
  - Intervalos ( $\text{valor}_{\text{inf}} < A < \text{valor}_{\text{sup}}$ )
    - Considerar todos os possíveis intervalos
      - Alguns intervalos adjacentes pode ser agregados

25/10/2012

André de Carvalho - ICMC/USP

19

## Atributos contínuos



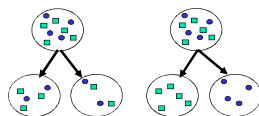
25/10/2012

André de Carvalho - ICMC/USP

20

## Medidas para escolha de atributo

- Existem várias medidas para determinar o atributo que melhor divide os dados
  - Geram diferentes partições dos dados
  - Medidas de impureza
    - Distribuição de classes dos dados após divisão
      - Quanto mais balanceadas as classes em uma partição, pior



25/10/2012

André de Carvalho - ICMC/USP

21

## Medidas de impureza

- Baseadas no grau de impureza dos nós filhos
  - Quando maior, pior
- Exemplos de medidas de impureza
  - Entropia
  - Gini
  - Erro de classificação
  - Qui-quadrado

25/10/2012

André de Carvalho - ICMC/USP

22

## Medidas para escolha de atributo

$$\text{Entropia}(v) = - \sum_{i=1}^C p(i/v) \log_2 p(i/v)$$

$$\text{Gini}(v) = 1 - \sum_{i=1}^C [p(i/v)]^2$$

$$\text{ErroClass}(v) = 1 - \max_i [p(i/v)]$$

Onde:

$P(i/v)$  = fração de dados pertencente a classe  $i$  em um nó  $v$

$C$  = número de classes

Considera-se que  $0 \log_2 0 = 0$

25/10/2012

André de Carvalho - ICMC/USP

23

## Exemplo

- Calcular a medida de impureza Gini para os dados abaixo:

$$\text{Gini}(v) = 1 - \sum_{i=1}^C [p(i/v)]^2$$

C1	0
C2	6
Gini=?	

C1	1
C2	5
Gini=?	

C1	2
C2	4
Gini=?	

C1	3
C2	3
Gini=?	

25/10/2012

André de Carvalho - ICMC/USP

24

## Exemplo

$$Gini(v) = 1 - \sum_{i=1}^C [p(i/v)]^2$$

$P(C1) = 0/6 = 0$      $P(C2) = 6/6 = 1$   
 $Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$   
 $P(C1) = 1/6$      $P(C2) = 5/6$   
 $Gini = 1 - (1/6)^2 - (5/6)^2 = 0.278$   
 $P(C1) = 2/6$      $P(C2) = 4/6$   
 $Gini = 1 - (2/6)^2 - (4/6)^2 = 0.444$   
 $P(C1) = 3/6$      $P(C2) = 3/6$   
 $Gini = 1 - (3/6)^2 - (3/6)^2 = 0.500$

C1	0
C2	6
Gini=0.000	

C1	1
C2	5
Gini=0.278	

C1	2
C2	4
Gini=0.444	

C1	3
C2	3
Gini=0.500	

25/10/2012

André de Carvalho - ICMC/USP

25

## Exercício

- Fazer os mesmos cálculos para as medidas de entropia e de erro de classificação

$$Entropia(v) = - \sum_{i=1}^C p(i/v) \log_2 p(i/v)$$

$$ErroClass(v) = 1 - \max_i [p(i/v)]$$

C1	0
C2	6
E=?	
C1	0
C2	6
Class=?	

C1	1
C2	5
E=?	
C1	1
C2	5
Class=?	

C1	2
C2	4
E=?	
C1	2
C2	4
Class=?	

C1	3
C2	3
E=?	
C1	3
C2	3
Class=?	

25/10/2012

André de Carvalho - ICMC/USP

26

## Medida Gini média ponderada

- Usada pelos algoritmos CART, SLIQ, SPRINT
- Quando um nó pai é dividido em  $k$  filhos, a impureza da divisão é definida por:

$$Gini_{divisão} = \sum_{f=1}^k \frac{N(v_f)}{N(v_p)} Gini(v_f)$$

Média ponderada

Onde:

$N(v_f)$ : número de objetos no filho  $f$  ( $v_f$ )

$N(v_p)$ : número de objetos no nó pai ( $v_p$ )

25/10/2012

André de Carvalho - ICMC/USP

27

## Medidas para selecionar divisão

- Escolha da medida influencia seleção da condição de teste
- Avaliação de qualidade de uma condição de teste
  - Comparar grau de impureza antes e após a divisão
    - Quanto maior a diferença, melhor a condição
    - Comparação pode se dar pela medida de ganho
      - Algoritmo ID3

25/10/2012

André de Carvalho - ICMC/USP

28

## Medida de ganho

$$\Delta = I(v_p) - \sum_{f=1}^k \frac{N(v_f)}{N(v_p)} I(v_f)$$

Onde:

$I(v)$ : mede o grau de impureza do nó  $v$

$N(v_f)$ : número de objetos no filho  $f$  ( $v_f$ )

$N(v_p)$ : número de objetos no nó pai ( $v_p$ )

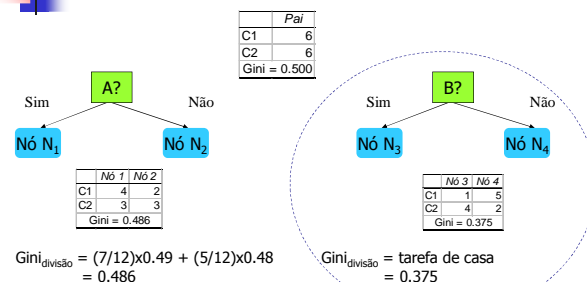
Quando a medida de impureza é entropia,  $\Delta$  mede o ganho de informação ( $\Delta_{info}$ )

25/10/2012

André de Carvalho - ICMC/USP

29

## Divisão de atributos binários



25/10/2012

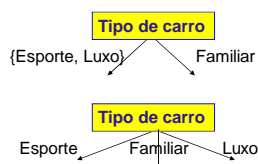
André de Carvalho - ICMC/USP

30

## Divisão de atributos nominais

### ■ Duas alternativas

- Divisão binária
- Divisão múltipla



25/10/2012

André de Carvalho - ICMC/USP

31

## Divisão de atributos nominais

### ■ Divisão binária

- Similar ao uso de atributos binários
  - Encontrar melhor binarização (ponto de referência)
- Índice de impureza é calculado para os 2 subconjuntos

### ■ Divisão múltipla

- Índice de impureza é calculado para cada divisão
- Resulta em subconjuntos em geral mais puros que a divisão binária

25/10/2012

André de Carvalho - ICMC/USP

32

## Exercício

- Definir a melhor divisão considerando divisão binária e divisão múltipla para:

Tipo de Carro			
	Familiar	Esporte	Luxo
C1	1	2	1
C2	4	1	1
Gini <sub>div</sub>	???		

Tipo de Carro		
	{Esporte, Luxo}	{Familiar}
C1	3	1
C2	2	4
Gini <sub>div</sub>	???	

Tipo de Carro		
	{Esporte}	{Familiar, Luxo}
C1	2	2
C2	1	5
Gini <sub>div</sub>	???	

25/10/2012

André de Carvalho - ICMC/USP

33

## Exercício

- Definir a melhor divisão considerando divisão binária e divisão múltipla para:

Tipo de Carro			
	Familiar	Esporte	Luxo
C1	1	2	1
C2	4	1	1
Gini <sub>div</sub>	0.393		

Tipo de Carro		
	{Esporte, Luxo}	{Familiar}
C1	3	1
C2	2	4
Gini <sub>div</sub>	0.400	

Tipo de Carro		
	{Esporte}	{Familiar, Luxo}
C1	2	2
C2	1	5
Gini <sub>div</sub>	0.419	

25/10/2012

André de Carvalho - ICMC/USP

34

## Divisão de atributos contínuos

- Várias possíveis escolhas para o ponto de referência
  - # possíveis divisões = # valores distintos
- Cada ponto de referência tem uma matriz de contagens associada a ele
  - Contagens das classes em cada uma das partições
- O mesmo vale para atributos discretos

25/10/2012

André de Carvalho - ICMC/USP

35

## Critério de parada

### ■ Diversas alternativas:

- Os objetos do nó atual têm a mesma classe
- Os objetos do nó atual têm valores iguais para os atributos de entrada, mas classes diferentes
- O número de objetos no nó é menor que um dado valor
- Todos os atributos já foram incluídos no caminho

25/10/2012

André de Carvalho - ICMC/USP

36

## Espaço de hipóteses

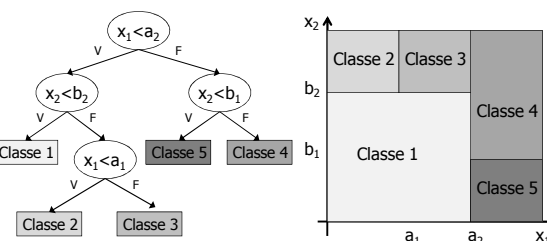
- Cada percurso da raiz a um nó folha representa uma regra de classificação
- Cada folha
  - Está associada a uma classe
  - Corresponde a uma região do espaço de soluções
    - Hiper-retângulo
      - Interseção de hiper-retângulos é um conjunto vazio
      - União é o espaço total

25/10/2012

André de Carvalho - ICMC/USP

37

## Espaço de hipóteses



25/10/2012

André de Carvalho - ICMC/USP

38

## Exemplo

- Sejam os dados abaixo referentes a solicitações de crédito bancário
  - Construir uma árvore de decisão que classifica aplicação para cartão de crédito

Idade	Renda	Classe
20	2000	Sim
30	5200	Não
60	5000	Sim
40	6000	Não
...	...	...

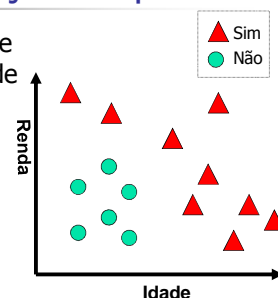
25/10/2012

André de Carvalho - ICMC/USP

39

## Busca no espaço de hipóteses

- Construir uma AD que classifica solicitante de cartão de crédito
  - Aprova (Sim)
  - Não aprova (Não)

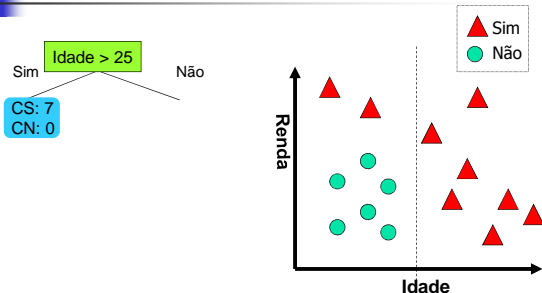


25/10/2012

André de Carvalho - ICMC/USP

40

## Busca no espaço de hipóteses

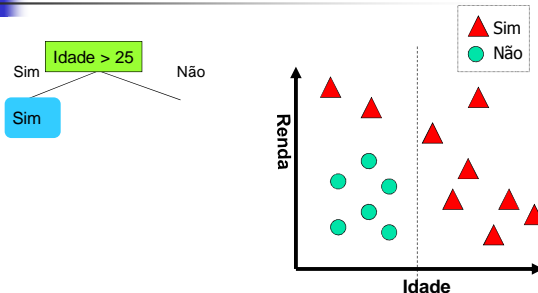


25/10/2012

André de Carvalho - ICMC/USP

41

## Busca no espaço de hipóteses

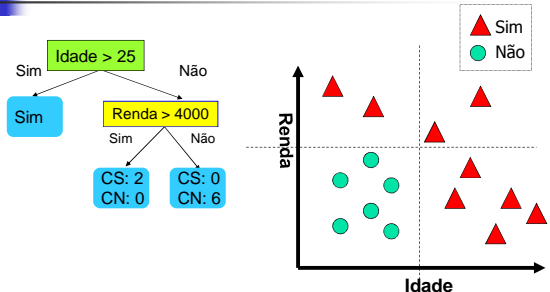


25/10/2012

André de Carvalho - ICMC/USP

42

## Busca no espaço de hipóteses

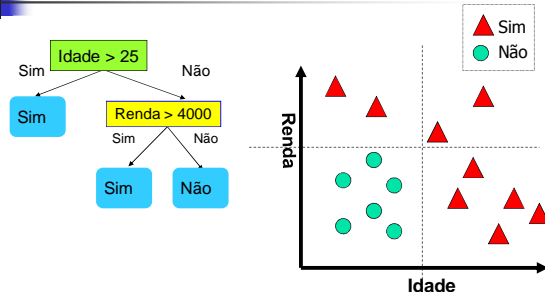


25/10/2012

André de Carvalho - ICMC/USP

43

## Busca no espaço de hipóteses



25/10/2012

André de Carvalho - ICMC/USP

44

## Algoritmo C4.5

- Proposto por Quinlan em 1993 como extensão do ID3
  - J48
  - C5.0
- Usa ganho de informação
- Pós-poda
- Todos os dados precisam caber na memória principal
  - Inadequado para grandes conjuntos de dados

25/10/2012

André de Carvalho - ICMC/USP

45

## Exercício

- Seja o seguinte cadastro de pacientes:

Nome	Febre	Enjôo	Manchas	Dores	Diagnóstico
João	sim	sim	pequenas	sim	doente
Pedro	não	não	grandes	não	saudável
Maria	sim	sim	pequenas	não	saudável
José	sim	não	grandes	sim	doente
Ana	sim	não	pequenas	sim	saudável
Leila	não	não	grandes	sim	doente

25/10/2012

André de Carvalho - ICMC/USP

46

## Exercício

- Usando medida de entropia,
  - Induzir uma árvore de decisão capaz de distinguir:
    - Pacientes potencialmente saudáveis
    - Pacientes potencialmente doentes
  - Testar a árvore para novos casos
    - (Luis, não, não, pequenas, sim)
    - (Laura, sim, sim, grandes, sim)

25/10/2012

André de Carvalho - ICMC/USP

47

## Conclusão


- Introdução
- Algoritmo de Hunt
- Medidas para selecionar divisão de atributos
- Critério de parada
- Espaço de hipóteses
- Variações

25/10/2012

André de Carvalho - ICMC/USP


48





## Perguntas

---



25/10/2012

André de Carvalho - ICMC/USP

49