

Aprendizado de Máquina

Árvores de Decisão



Prof. Dr. André C. P. L. F. de Carvalho
 Pós-doutorando: Isvani Frias-Blanco
 ICMC-USP

Principais tópicos

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

© André de Carvalho - ICMC/USP 2

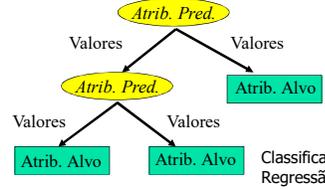
Introdução

- Explicação das decisões pode ser importante para algumas aplicações
 - RNAs, SVMs e RNPs são caixas pretas
- Modelos interpretáveis são gerados por algumas algoritmos de AM
 - Árvores de características (decisão)
 - Conjunto de regras
 - Redes Bayesianas

© André de Carvalho - ICMC/USP 3

Árvores de características

- Alguns algoritmos de AM particionam características (atributos) de forma hierárquica

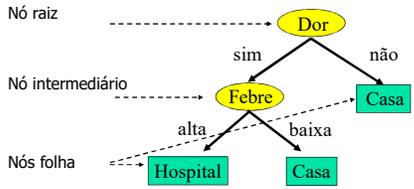


Classificação: decisão (AD)
 Regressão: regressão (AR)

© André de Carvalho - ICMC/USP 4

Algoritmo de indução de AD

- Induzem modelos representados por ADs

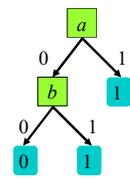


Nó raiz: Dor
 Nó intermediário: Febre
 Nós folha: Hospital, Casa

© André de Carvalho - ICMC/USP 5

Outro exemplo simples

a	b	a v b
0	0	0
0	1	1
1	0	1
1	1	1



Nós internos e raiz: atributos preditivos
 Nós externos (folhas): atributo alvo

© André de Carvalho - ICMC/USP 6

Exercício

- Encontrar árvore de decisão para:
 - a AND b
 - a XOR b
 - (a AND b) OR (b AND \bar{c})

© André de Carvalho - ICMC/USP 7

Algoritmo de indução de AD

- Existem vários, entre eles:
 - Algoritmo de Hunt
 - Um dos primeiros
 - Base de vários algoritmos atuais
 - CART
 - ID3
 - C4.5
 - VFDT

© André de Carvalho - ICMC/USP 8

Algoritmo de Hunt

- Seja X_t o conjunto de objetos de treinamento que atingem o nó t

*Se todos os objetos de $X_t \in$ a mesma classe y
Então O nó t é um nó folha rotulado pela classe y
Senão Selecionar um atributo preditivo teste para dividir X_t
Dividir X_t em subconjuntos usando valores desse atributo
Aplicar algoritmo a cada subconjunto gerado*

© 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	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

© 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	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

© 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	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

© André de Carvalho - ICMC/USP 12

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 13

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 14

Indução de ADs

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

© André de Carvalho - ICMC/USP 15

Escolha do atributo preditivo

- Atributo preditivo que melhor particiona conjunto atual de objetos
 - Para um mesmo atributo preditivo, diferentes partições podem ser geradas
 - Necessário escolher:
 - Atributo preditivo mais discriminativo
 - Melhor partição para esse atributo
 - Como dividir os objetos

© André de Carvalho - ICMC/USP 16

Como dividir os objetos

- Depende do tipo do atributo preditivo e do número de divisões a serem geradas
 - Atributo preditivo binário
 - Divisão binária
 - Árvore binária
 - Atributo preditivo n-ário
 - Divisão binária
 - Divisão n-ária (n > 2)

© André de Carvalho - ICMC/USP 17

Atributo binário

- Teste mais simples que existe
 - Define dois possíveis resultados (filhos)

© André de Carvalho - ICMC/USP 18

Atributo n-ário

- Divisão
 - Binária
 - N-ária
- Depende do tipo do atributo
 - Simbólico
 - Nominal
 - Ordinal
 - Numérico
 - Discreto
 - Contínuo

© André de Carvalho - ICMC/USP 19

Divisão binária para atributo n-ário

- Único teste com 2 possíveis resultados (filhos)
- Condição de teste é uma comparação
 - Ex.: $A < \text{valor}$, $A = \text{valor}$, $A \in \{\text{valores}\}$, ...
 - Escolher valor(es) que gera(m) melhor partição
 - Ponto de referência
 - Tipo simbólico: agrupar parte dos valores em cada ramo
 - Ordinais: valores agrupados não devem violar relação de ordem
 - Nominais: grupos devem fazer sentido

© André de Carvalho - ICMC/USP 20

Atributos simbólicos

Tipo de carro

}

Esporte

Família

Luxo

Refrigerante

}

Pequeno

Médio

Grande

Gigante

© André de Carvalho - ICMC/USP 21

Divisão binária para atributo n-ário

```

graph TD
    TC1[Tipo de carro] --> E1["{Esporte, Família}"]
    TC1 --> L1[Luxo]
    TC2[Tipo de carro] --> E2["{Esporte, Luxo}"]
    TC2 --> F2[Família]
    R[Refrigerante] --> P["{Pequeno, Médio}"]
    R --> G["{Grande, Gigante}"]
    P --> W["Peso > 60 Kg"]
    W --> S1[Sim]
    W --> N1[Não]
    G --> A["Idade = 18"]
    A --> S2[Sim]
    A --> N2[Não]
    
```

© André de Carvalho - ICMC/USP 22

Divisão n-ária para atributo n-ário

- Atributos simbólicos
 - Duas alternativas para definir número de resultados do teste
 - Fazer #ramos = #possíveis valores
 - Agrupar parte dos valores em cada ramo
 - Ordinais
 - Nominais

© André de Carvalho - ICMC/USP 23

Divisão n-ária para atributo n-ário

```

graph TD
    TC[Tipo de carro] --> E[Esporte]
    TC --> F[Família]
    TC --> L[Luxo]
    P[Peso] --> W1["< 10Kg"]
    P --> W2["[50Kg, 70Kg]"]
    P --> W3["(70Kg, 80Kg)"]
    P --> W4["> 90Kg"]
    R[Refrigerante] --> P1["{Pequeno, Médio}"]
    R --> P2["{Grande}"]
    R --> P3["{Gigante}"]
    
```

© André de Carvalho - ICMC/USP 24

Divisão n-ária para atributo n-ário

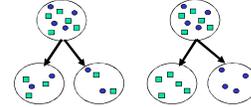
- Condição de teste formada por 2 ou mais comparações
 - Cada comparação pode ter 1 ou mais operadores relacionais
 - Um operador
 - Ex. $A < \text{valor}$, $A = \text{valor}$
 - Mais de um operador
 - Ex.: $\text{valor}_{\text{inf}} < A < \text{valor}_{\text{sup}}$
 - Escolher valores (pontos de referência)

© André de Carvalho - ICMC/USP

25

Medidas para escolha de atributo

- Seleccionam atributo que melhor discrimina os objetos atuais
 - Buscam partições mais puras após divisão
 - Quanto mais homogêneas as partições, mais puras
 - Medidas de impureza

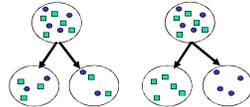


© André de Carvalho - ICMC/USP

26

Medidas de impureza

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



© André de Carvalho - ICMC/USP

27

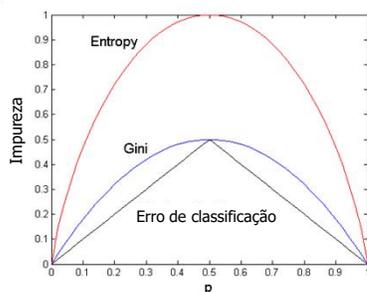
Propriedade de uma medida de impureza

- Seja uma tarefa de classificação binária
 - Deve ser dependente apenas da magnitude relativa das classes
 - Não muda se a magnitude das duas classes for multiplicada por uma constante
 - Proporção dos exemplos em cada classe
 - Não deve mudar se as classes positiva e negativa são trocadas
 - Deve ser 0, se a proporção de exemplos de uma das classes for 1 (da outra será 0)

© André de Carvalho - ICMC/USP

28

Comparação: duas classes



André de Carvalho - ICMC/USP

29

Medidas de impureza

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

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

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

Onde:

 $p(i/v)$ = fração de dados pertencente a classe i em um nó v C = número de classesConsidera-se que $0 \log_2 0 = 0$

© André de Carvalho - ICMC/USP

30

Exemplo

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

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

C1	0
C2	6
Gini=?	

C1	1
C2	5
Gini=?	

C1	2
C2	4
Gini=?	

C1	3
C2	3
Gini=?	

© André de Carvalho - ICMC/USP 31

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	

© André de Carvalho - ICMC/USP 32

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	1
C2	5
E=?	

C1	2
C2	4
E=?	

C1	3
C2	3
E=?	

C1	0
C2	6
Class=?	

C1	1
C2	5
Class=?	

C1	2
C2	4
Class=?	

C1	3
C2	3
Class=?	

© André de Carvalho - ICMC/USP 33

Medida Gini média ponderada

- Usada pelos algoritmos CART, SLIQ, SPRINT
- Quando um nó pai possui 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 (v_f)
- $N(v_p)$: número de objetos no nó pai (v_p)

© André de Carvalho - ICMC/USP 34

Divisão de atributos binários

Atrib. A?

Sim → Nó N₁
Não → Nó N₂

Nó 1		Nó 2	
C1	4	C1	2
C2	3	C2	3
Gini =			

Gini_{divisão} =

Atrib. B?

Sim → Nó N₃
Não → Nó N₄

Nó 3		Nó 4	
C1	1	C1	5
C2	4	C2	2
Gini =			

Gini_{divisão} =

© André de Carvalho - ICMC/USP 35

Divisão de atributos binários

Atrib. A?

Sim → Nó N₁
Não → Nó N₂

Nó 1		Nó 2	
C1	4	C1	2
C2	3	C2	3
Gini = 0.486			

Gini_{divisão} = (7/12)x0.49 + (5/12)x0.48 = 0.486

Atrib. B?

Sim → Nó N₃
Não → Nó N₄

Nó 3		Nó 4	
C1	1	C1	5
C2	4	C2	2
Gini = 0.375			

Gini_{divisão} = tarefa de casa = 0.375

© André de Carvalho - ICMC/USP 36

Exercício

- Divisão de atributos n-ários
 - Duas alternativas
 - Divisão binária
 - Igual a divisão de atributos binários
 - Divisão n-ária

Tipo de carro

{Esporte, Luxo} Familiar

```

    graph TD
      A[Tipo de carro] --> B[Esporte, Luxo]
      A --> C[Familiar]
      B --> D[Esporte]
      B --> E[Luxo]
      C --> F[Familiar]
      
```

Tipo de carro

Esporte Familiar Luxo

© André de Carvalho - ICMC/USP 37

Exercício

- Divisão n-ária
 - Índice de impureza é calculado para cada um dos n subconjuntos
 - Necessário definir valor de n
 - Resulta em subconjuntos, em geral, mais puros que a divisão binária
 - Porém, nós com menos exemplos = menor confiança

© André de Carvalho - ICMC/USP 38

Exercício

- Definir a melhor divisão considerando divisão binária e divisão n-ária para:

Tipo de carro			
Classe	Família	Esporte	Luxo
C1	1	2	1
C2	4	1	1
Gini _{div}			

Tipo de carro			
Classe	Família	Esporte e Luxo	
C1	1	3	
C2	4	2	
Gini _{div}			

Tipo de carro		
Classe	Esporte	Família e Luxo
C1	2	2
C2	1	5
Gini _{div}		

© André de Carvalho - ICMC/USP 39

Exercício

- Definir a melhor divisão considerando divisão binária e divisão n-ária para:

Tipo de carro			
Classe	Família	Esporte	Luxo
C1	1	2	1
C2	4	1	1
Gini _{div}	0.393		

Tipo de carro		
Classe	Família	Esporte e Luxo
C1	3	1
C2	2	4
Gini _{div}	0.400	

Tipo de carro		
Classe	Esporte	Família e Luxo
C1	2	2
C2	1	5
Gini _{div}	0.419	

© André de Carvalho - ICMC/USP 40

Atributos n-ários

- Várias possíveis posições de referência
- Cada posição tem uma matriz de contagens associada a ela
 - Contagens das proporções das classes em cada uma das partições

© André de Carvalho - ICMC/USP 41

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 do nó é menor que um dada quantidade
 - Todos os atributos preditivos já foram incluídos no caminho atual

© André de Carvalho - ICMC/USP 42

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
...		

© André de Carvalho - ICMC/USP 43

Busca no espaço de hipóteses

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

© André de Carvalho - ICMC/USP 44

Busca no espaço de hipóteses

© André de Carvalho - ICMC/USP 45

Busca no espaço de hipóteses

© André de Carvalho - ICMC/USP 46

Busca no espaço de hipóteses

© André de Carvalho - ICMC/USP 47

Busca no espaço de hipóteses

© André de Carvalho - ICMC/USP 48

Busca no espaço de hipóteses

```

graph TD
    A[Idade > 25] -- Sim --> B[Aprova]
    A -- Não --> C[Renda > 4000]
    C -- Sim --> D[Aprova]
    C -- Não --> E[Não aprova]
        
```

© André de Carvalho - ICMC/USP 49

Algoritmo de indução de AD

- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido

© André de Carvalho - ICMC/USP 50

Algoritmo de indução de AD

- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido

© André de Carvalho - ICMC/USP 51

Algoritmo de indução de AD

- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido

© André de Carvalho - ICMC/USP 52

Algoritmo de indução de AD

- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido

© André de Carvalho - ICMC/USP 53

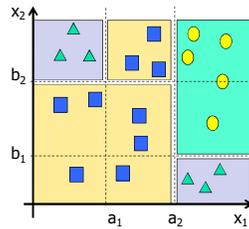
Algoritmo de indução de AD

- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido

© André de Carvalho - ICMC/USP 54

Algoritmo de indução de AD

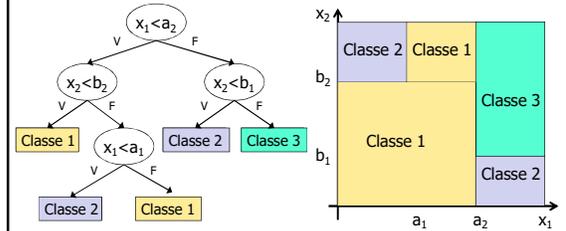
- Valores numéricos e simbólicos
- Grau de pureza é calculado para cada atributo
 - Atributo que gera menos impureza é escolhido
- Critério de parada precisa ser definido



© André de Carvalho - ICMC/USP

55

Árvore e partição do espaço de hipóteses



© André de Carvalho - ICMC/USP

56

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 um hiper-retângulo no espaço de soluções
 - Cada classe é representada por um conjunto de hiper-retângulos
 - Interseção de hiper-retângulos é um conjunto vazio
 - União de hiper-retângulos cobre todo o espaço

© André de Carvalho - ICMC/USP

57

Aspectos positivos das ADs

- Baixo custo de indução e dedução
- Fácil interpretação da hipótese induzida
 - Para árvores pequenas
- Acurácia comparável a de outros classificadores
 - Para conjuntos de dados de baixa complexidade
- Indica atributos preditivos mais relevantes
- Atributos preditivos podem ser numéricos ou simbólicos

© André de Carvalho - ICMC/USP

58

Aspectos negativos das ADs

- Dificuldade para predição de valores contínuos
 - Árvores de regressão
- Baixo desempenho em problemas com muitas classes e poucos dados
- Abordagem gulosa
- Limitação de hipóteses a hiper-retângulos

© André de Carvalho - ICMC/USP

59

Overfitting

- Partição recursiva pode gerar árvores perfeitamente ajustadas aos dados
- Decisões são baseadas em conjuntos cada vez menores de dados
 - Níveis mais profundos podem ter muito poucos dados
 - Presença de ruído nos dados afeta bastante escolha de atributos para esses nós
 - Reduz capacidade de generalização
 - Poda

© André de Carvalho - ICMC/USP

60

Poda de árvores

- Elimina parte da árvore
- Pode ser realizada em duas etapas
 - Durante indução (pré-poda)
 - Parar o crescimento da árvore mais cedo
 - Após indução (pós-poda)
 - Crescer a árvore completa e depois podá-la
 - Mais lento, porém mais confiável

© André de Carvalho - ICMC/USP

61

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

© André de Carvalho - ICMC/USP

62

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)

© André de Carvalho - ICMC/USP

63

Conclusão

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

© André de Carvalho - ICMC/USP

64

Perguntas



© André de Carvalho - ICMC/USP

65