

# ARVORES DE DECISÃO

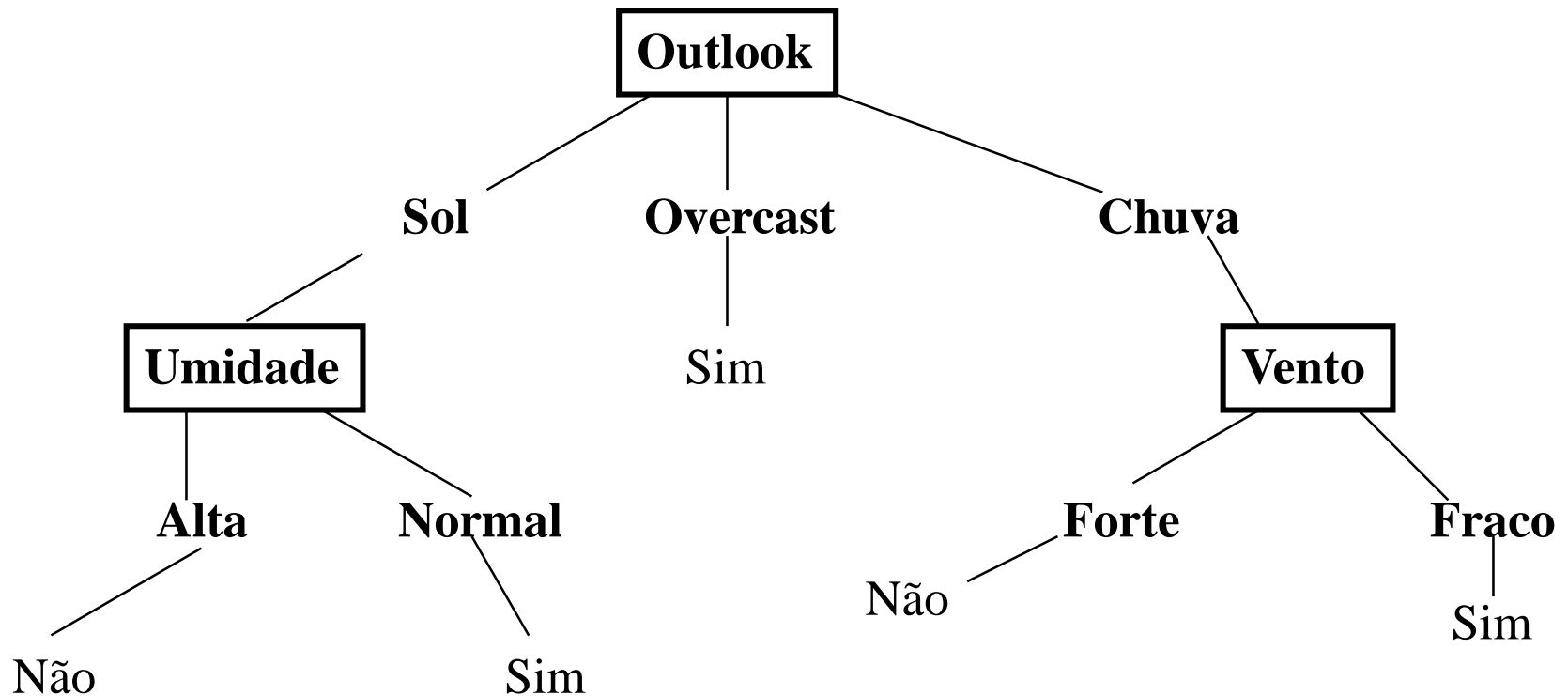
SCC0276/SCC5871/MAI5025

APRENDIZADO DE MÁQUINA  
Profa. Roseli Ap. Francelin Romero

# Árvore de Decisão

- Representação de Árvores de Decisão
- Algoritmo ID3
- Conceito de Entropia e Ganho de Informação
- Overfitting

# Árvore de Decisão

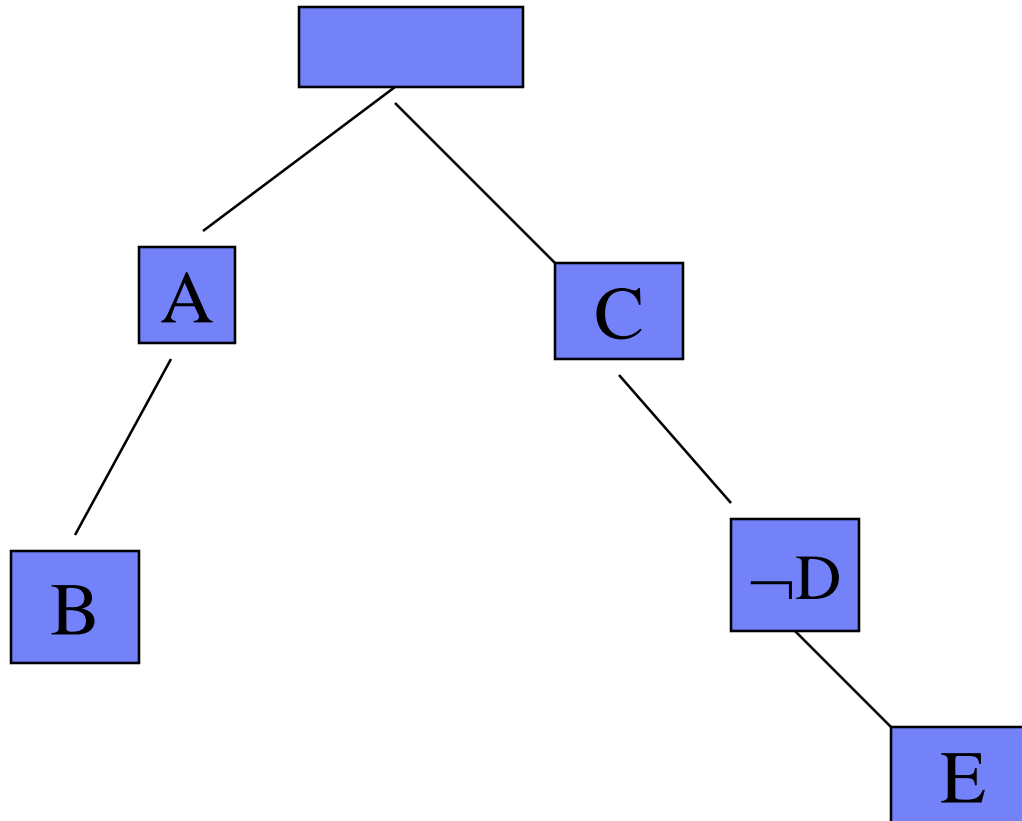


# Árvore de Decisão

- PlayTennis
- Cada nó interno testa um ATRIBUTO
- Cada ramo corresponde a um valor do atributo
- Cada nó terminal designa uma classificação.

# Árvore de Decisão

■  $(A \wedge B) \vee (C \wedge \neg D \wedge E)$



# Árvore de Decisão

## ■ Quando utilizar?

- Problemas descritos por pares de atributo/valor
- Função objetivo é discreta
- Hipóteses disjuntivas são requeridas
- ruídos nos dados

Exemplos: diagnóstico médicos e de equipamentos, análise de crédito.

# Árvore de Decisão

## Indução Top-Down

### Main Loop

1.  $A$  o melhor atributo de decisão para o próximo nó.
2. Designar  $A$  como o atributo de decisão p/ o nó.
3. Para cada valor de  $A$ , criar um novo descendente.
4. Escolher exemplos de treinamento para os nós folha
5. Se exemplos de treinamento forem perfeitamente classificados, então PARE, senão iterar sobre novos nós folha.

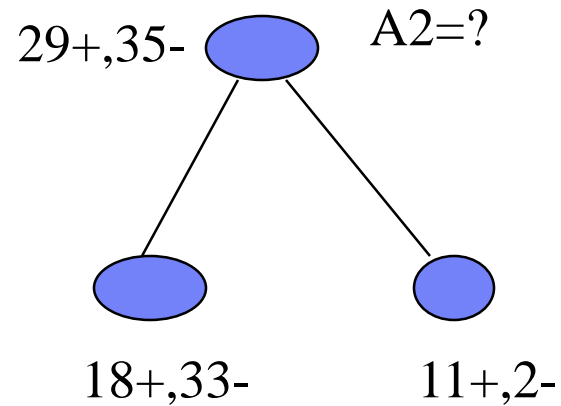
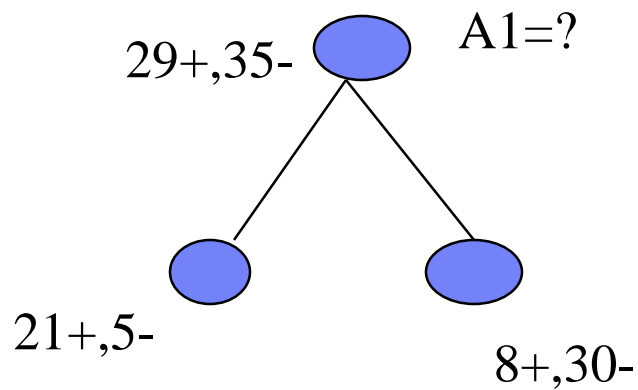
# Exemplos de Treinamento

<i><b>DAY</b></i>	<i><b>OUTLOOK</b></i>	<i><b>TEMPERATURA</b></i>	<i><b>UMIDADE</b></i>	<i><b>VENTO</b></i>	<i><b>PLAYTENN</b></i>
<b>D1</b>	SOL	QUENTE	ALTA	FRACO	NÃO
<b>D2</b>	SOL	QUENTE	ALTA	FORTE	NÃO
<b>D3</b>	NUBLADO	QUENTE	ALTA	FRACO	SIM
<b>D4</b>	CHUVA	AMENO	ALTA	FRACO	SIM
<b>D5</b>	CHUVA	FRIO	NORMAL	FRACO	SIM
<b>D6</b>	CHUVA	FRIO	NORMAL	FORTE	NÃO
<b>D7</b>	NUBLADO	FRIO	NORMAL	FORTE	SIM
<b>D8</b>	SOL	AMENO	ALTA	FRACO	NÃO
<b>D9</b>	SOL	FRIO	NORMAL	FRACO	SIM
<b>D10</b>	CHUVA	AMENO	NORMAL	FRACO	SIM
<b>D11</b>	SOL	AMENO	NORMAL	FORTE	SIM
<b>D12</b>	NUBLADO	AMENO	ALTA	FORTE	SIM
<b>D13</b>	NUBLADO	QUENTE	NORMAL	FRACO	SIM
<b>D14</b>	CHUVA	AMENO	ALTA	FORTE	NÃO



# Árvore de Decisão

- Qual atributo é o MELHOR?



# Entropia

Entropia S

1.

.5

0.

0.5

1.0

$S$  = conjunto de exemplos treinamento

$P_+$  = proporção de exemplos positivos

$p_-$  = proporção de exemplos negativos

Entropia mede a IMPURIDADE de S

$$E(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

# Entropia

- Da teoria de Informação:

Entropia (  $S$  ) = número esperado de bits necessários para representar uma classe (+ or -) dos membros de  $S$  (**sob código de menor comprimento e ótimo**).

- PORQUE? (se  $p_+ = 1$  ou  $p_+ = 0.5$ )

Um código de compr. ótimo designa  $-\log_2 p$  bits com probabilidade  $p$ . Então, o número esperado de bits representar + ou - membros de  $S$  é:

$$p_+ (-\log_2 p_+) + p_- (-\log_2 p_-)$$

$$\text{Entropia ( } S \text{ )} \equiv - p_+ \log_2 p_+ - p_- \log_2 p_-$$

# Entropia

## ■ EXEMPLO:

$$S = [ 9+ , 5- ]$$

$$\text{ENTROPIA} ( [ 9+ , 5- ] ) = -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) = 0.940$$

# Ganho de Informação

- $\text{Gain} ( S, A )$  = redução esperada na entropia devido a escolha do atributo  $A$ .

$$\text{Gain} ( S, A ) \equiv \text{Entropia} ( S ) - \sum_{v \in \text{Values} (A)} \frac{|S_v|}{|S|} \text{Entropia} (S_v)$$

Valor (Wind) = Fraco, Forte

$S = [9+, 5-]$        $S_{\text{fraco}} = [6+, 2-]$        $S_{\text{forte}} = [3+, 3-]$

$$\text{Gain} (S, \text{Wind}) = \text{Entropia}(S) - \sum_{v \in \{\text{Fraco}, \text{Forte}\}} \frac{|S_v|}{|S|} \text{Entropia}(S_v)$$

$$\begin{aligned}
 &= \text{Entropia}(S) - (8/14)\text{Entropia}(\text{Fraco}) - \\
 &\quad (6/14)\text{Entropia}(\text{Forte}) = \\
 &\quad 0.94 - (8/14) 0.811 - (6/14) 1.00 = \\
 &\quad 0.048
 \end{aligned}$$

# Exemplos de Treinamento

<i><b>DAY</b></i>	<i><b>OUTLOOK</b></i>	<i><b>TEMPERATURA</b></i>	<i><b>UMIDADE</b></i>	<i><b>VENTO</b></i>	<i><b>PLAYTENN</b></i>
<b>D1</b>	SOL	QUENTE	ALTA	FRACO	NÃO
<b>D2</b>	SOL	QUENTE	ALTA	FORTE	NÃO
<b>D3</b>	NUBLADO	QUENTE	ALTA	FRACO	SIM
<b>D4</b>	CHUVA	AMENO	ALTA	FRACO	SIM
<b>D5</b>	CHUVA	FRIO	NORMAL	FRACO	SIM
<b>D6</b>	CHUVA	FRIO	NORMAL	FORTE	NÃO
<b>D7</b>	NUBLADO	FRIO	NORMAL	FORTE	SIM
<b>D8</b>	SOL	AMENO	ALTA	FRACO	NÃO
<b>D9</b>	SOL	FRIO	NORMAL	FRACO	SIM
<b>D10</b>	CHUVA	AMENO	NORMAL	FRACO	SIM
<b>D11</b>	SOL	AMENO	NORMAL	FORTE	SIM
<b>D12</b>	NUBLADO	AMENO	ALTA	FORTE	SIM
<b>D13</b>	NUBLADO	QUENTE	NORMAL	FRACO	SIM
<b>D14</b>	CHUVA	AMENO	ALTA	FORTE	NÃO

# Selecionando o Próximo Atributo

Qual atributo é o melhor classificador?

$S = [9+, 5-]$   
 $E = 0.940$

**Umidade**

ALTA

NORMAL

$[3+, 4-]$

$E = 0.985$

$[6+, 1-]$

$E = 0.592$

$\text{GAIN}(S, \text{Umidade}) = 0.94 - (7/14) 0.985 +$   
 $- (7/14) 0.592 = 0.151$

$S = [9+, 5-]$   
 $E = 0.940$

**Vento**

FRACO

FORTE

$[6+, 2-]$

$E = 0.811$

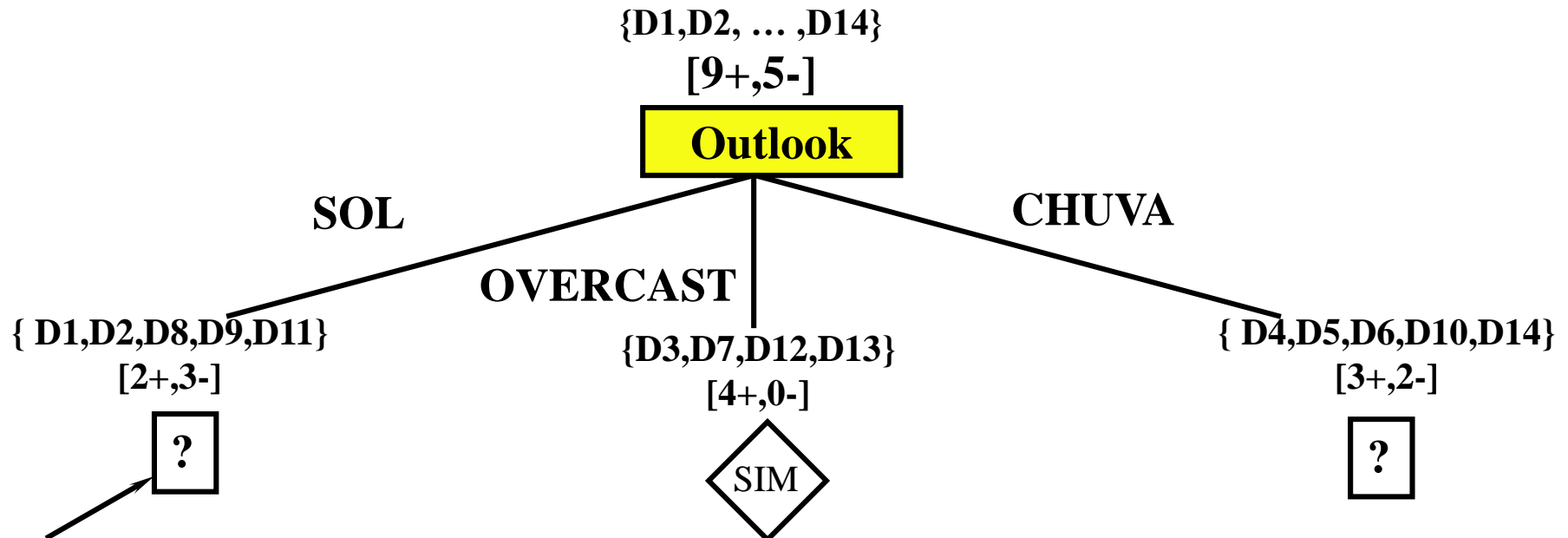
$[3+, 3-]$

$E = 1.00$

$\text{GAIN}(S, \text{Vento}) = 0.94 - (8/14) 0.811 +$   
 $- (6/14) 1.00 = 0.048$



# Selecionando o Próximo Atributo



■ Qual atributo deveria ser testado aqui?

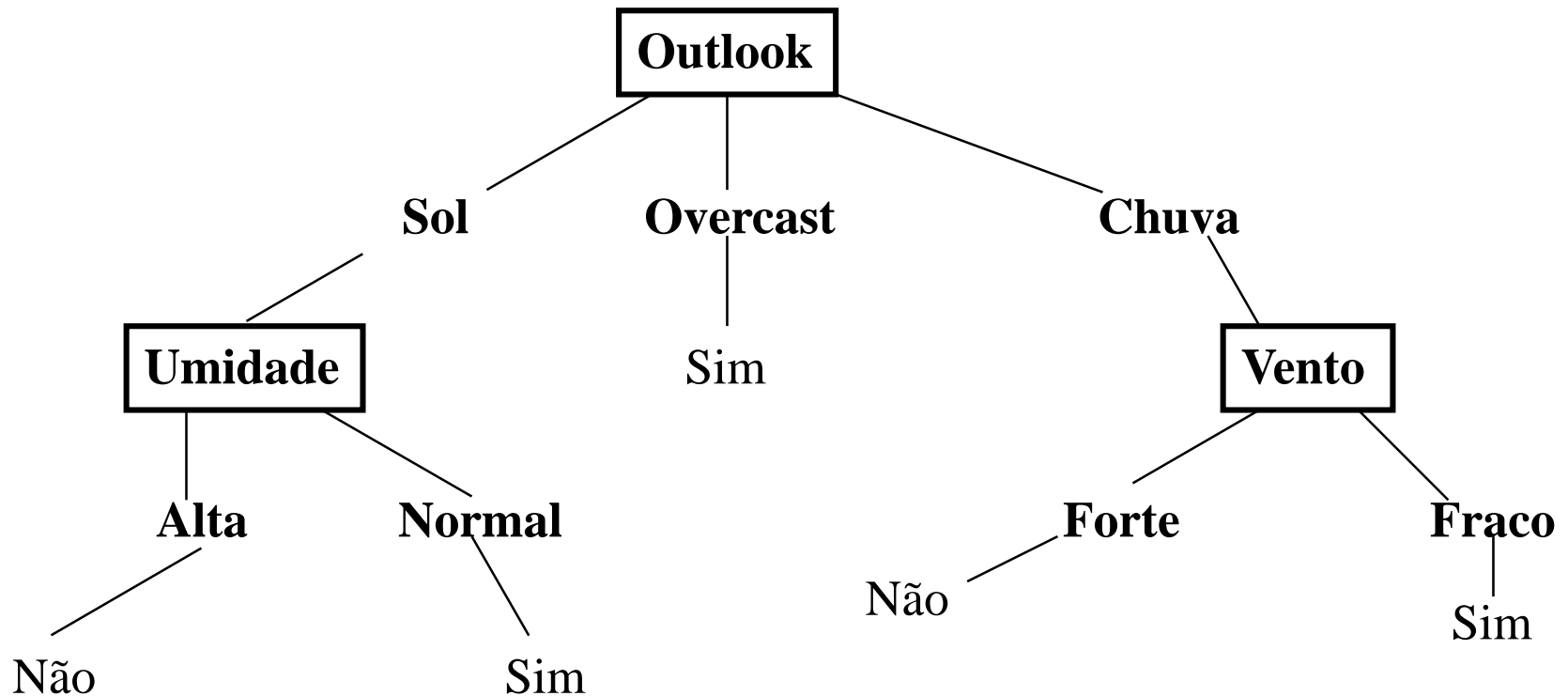
$$S_{\text{sol}} = \{ D1, D2, D8, D9, D11 \}$$

$$\text{Gain}(S_{\text{sol}}, \text{Umidade}) = 0.97 - (3/5) 0.0 - (2/5) 0.0 = 0.97$$

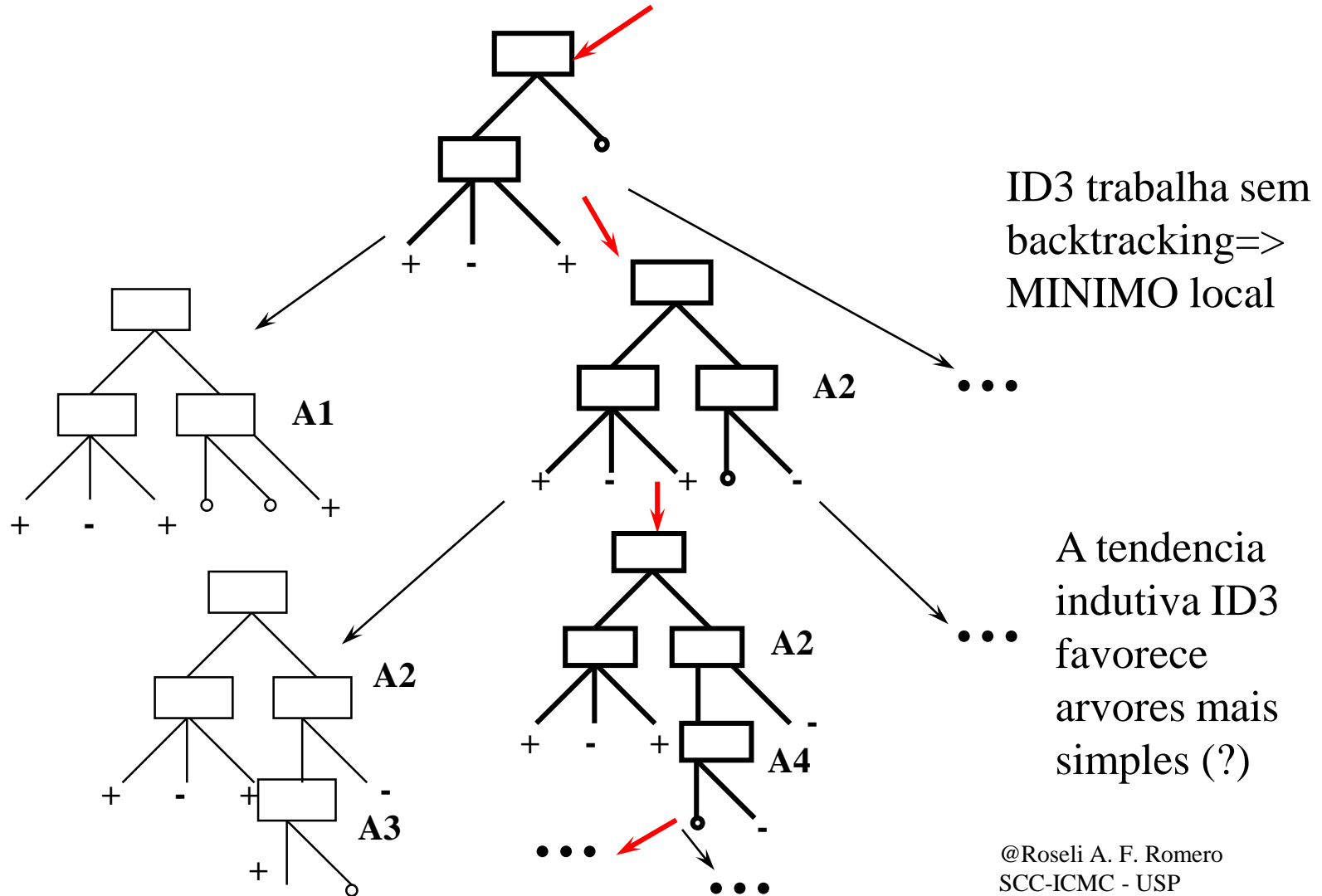
$$\text{Gain}(S_{\text{sol}}, \text{Temperatura}) = 0.97 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = 0.57$$

$$\text{Gain}(S_{\text{sol}}, \text{Vento}) = 0.97 - (2/5) 1.0 - (3/5) 0.918 = 0.19$$

# Árvore de Decisão



# Espaço de Hipóteses pesquisado por ID3



# Espaço de Hipóteses pesquisado por ID3

- Espaço de Hipótese é completo!
  - ◆ Com certeza a função alvo se encontra no espaço de busca
- ID3 fornece uma única saída
  - ◆ Não se pode considerar quantas DT alternativas são consistente com o conjunto de treinamento
- Nenhum back-tracking
  - ◆ Mínimo Local ...
- Usa propr.estatística dos dados (ganho de inf.)
  - ◆ Robusto a dados com ruídos ...

# Inductive Bias

- Como ID3 consegue generalizar a partir dos exemplos aprendidos?
- Inductive bias de ID3:
  - seleciona a favor de “árvores menores preferidas”
  - seleciona árvores que colocam os atributos com maior “ganho de informação” próximo ao nó raiz.

# Overfitting

- Considere o erro da hipótese  $h$  sobre
  - ◆ conjunto de treinamento:  $error_{train}(h)$
  - ◆ Distribuição inteira  $D$  dos dados:  $error_D(h)$
- Dada uma Hipótese  $h \in H$  para a qual ocorre um **overfitting** se existir uma hipótese alternativa  $h' \in H$  tal que

$$error_{train}(h) < error_{train}(h')$$

mas

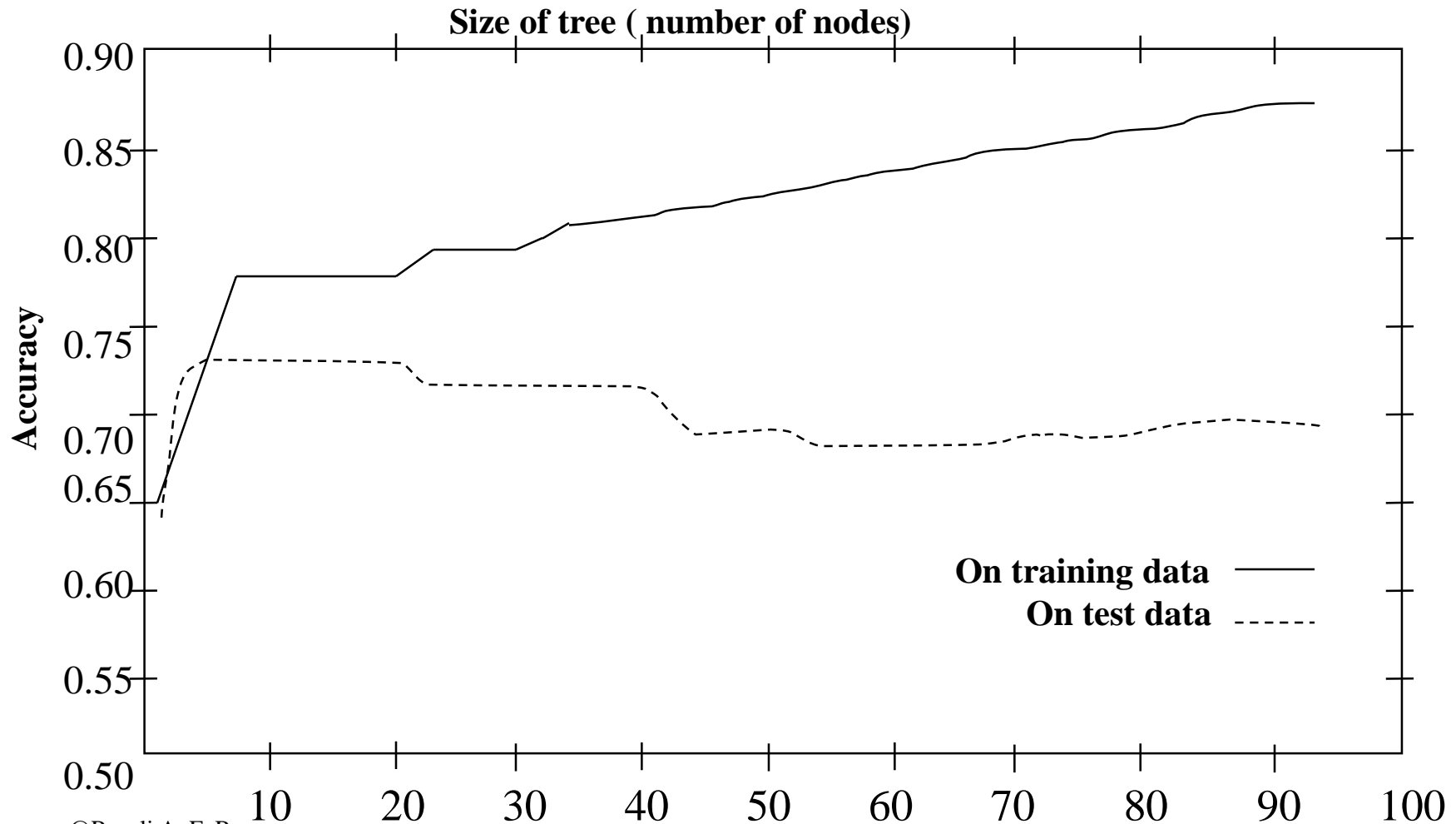
$$error_D(h) > error_D(h')$$

Overfitting é possível de ocorrer quando:

- ◆ existem ruídos nos dados ou erros aleatórios.
- ◆ Existem poucos exemplos associados com os nós folhas.

Ocorre muito na prática em métodos de aprendizado. Um estudo experimental com ID3 envolvendo 5 tarefas com dados não determinísticos e ruídos, overfitting diminuiu a precisão da DT em 10 - 25% dos problemas.

# Overfitting in Decision Trees Learning (Mitchell, 1998)



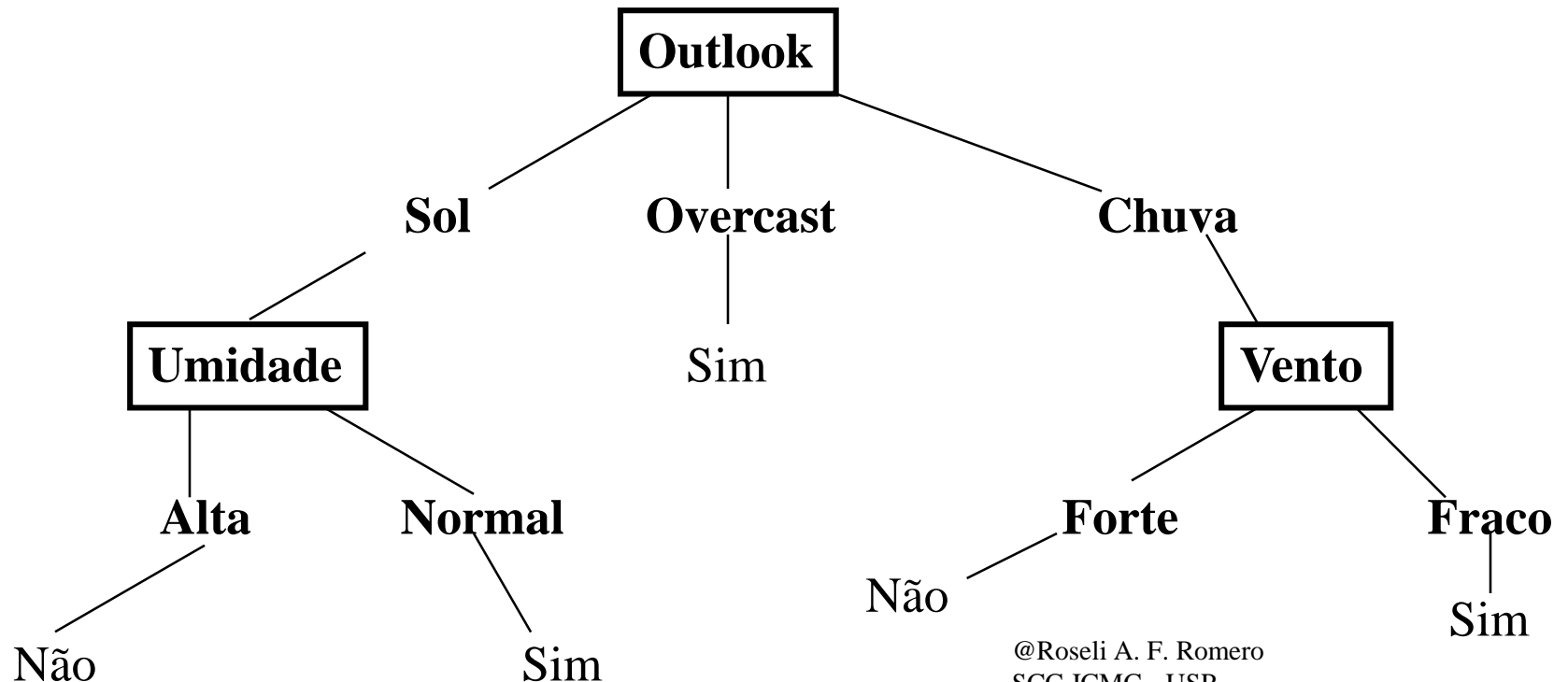


# Overfitting in Decision Trees

Considere a adição de ruído no Exemplo anterior no. 15:

(Sol, Quente, Normal, Forte, Playtennis = Não)

Qual o efeito na árvore de decisão anterior?



# Overfitting

## Evitar Overfitting

- ◆ Abordagens que cessam o crescimento da árvore antes que ela atinja o ponto onde ela perfeitamente classifica os dados de treinamento.
- ◆ Abordagens que permitem ocorrer um overfitting e então depois através de um pruning diminuir a árvore.

# Overfitting

A segunda abordagem é mais usada porque na primeira não se sabe exatamente qual é o ponto onde se deve parar.

## **Pruning de Erro Reduzido (Quinlan, 1987)**

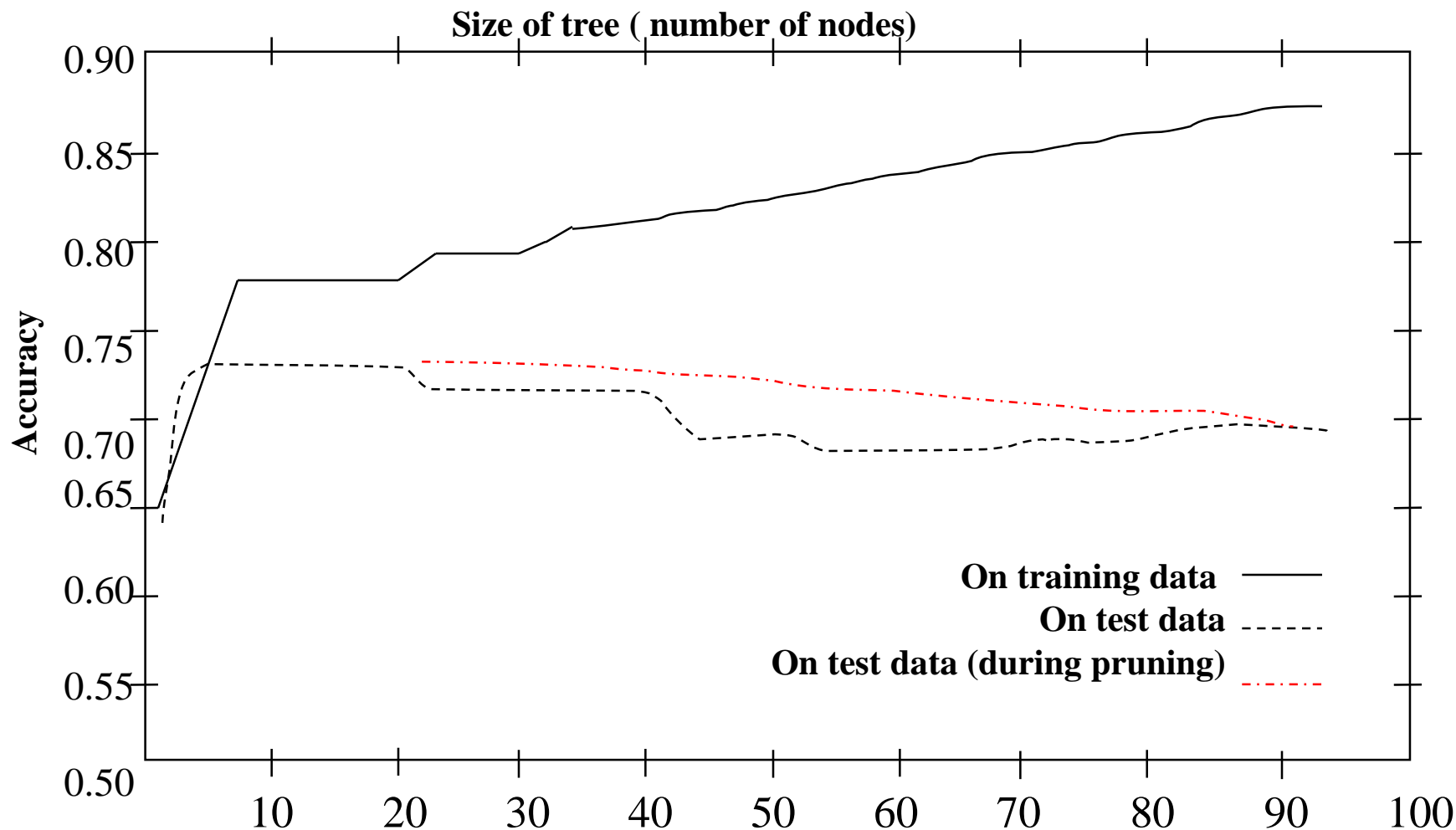
Considerar cada nó da árvore como um nó candidato a corte. Cortar um nó significa remover a sub- árvore a partir daquele nó, tornando-o um nó terminal (ou folha) e designando-o a mais comum classificação dos exemplos de treinamento afiliados com aquele nó.

# Pruning de Erro Reduzido

Nós são removidos apenas se a árvore resultante não desempenha pior que aquela original.

Nós são cortados iterativamente sempre escolhendo o nó cuja remoção aumenta a precisão da árvore de decisão sobre o conjunto de Validação.

# Efeito do Pruning de Erro Reduzido



# Regra Post-Pruning

Usado por C4.5 (Quinlan, 1993)

- ◆ Deixar ocorrer overfitting crescendo a árvore para treinar os dados.
- ◆ Converter a árvore aprendida num conjunto de regras.
- ◆ Cortar (generalizar) cada regra removendo algumas pré-condições que resultam em melhorar sua precisão estimada.
- ◆ Escolher as regras finais por sua precisão estimada e considerá-las nesta sequência quando classificando instâncias subsequentes.

# Regra Post-Pruning

- IF (Outlook = Sunny)  $\wedge$  (Humidity = High)  
THEN PlayTennis = No
- cada regra é podada removendo-se algum antecedente: ou (Outlook = Sunny) ou (Humidity = High)
- escolhe-se o antecedente que melhora a precisão estimada da regra. Nada é feito se a precisão piorar.

# Regra Post-Pruning

- Uma forma de avaliar a precisão estimada da regra é considerar um Conjunto de Validação disjunto do conjunto de treinamento.
- Usar um conjunto de validação para guiar a poda é uma indicação boa quando se tem grande volume de dados. Do contrário, ele acaba reduzindo ainda mais o conjunto de treinamento.



# Atributos de valores contínuos

- ID3 é restrito a assumir apenas valores discretos:
  - atributo alvo predito pela árvore é discreto
  - os atributos testado nos nós de decisão da árvore deve também ser discretos.

Mas, a segunda restrição pode ser relaxada para valores contínuos

Para um atributo  $A$ , que é um atributo de valor contínuo, o algoritmo cria um novo

- Um novo atributo booleano  $Ac$  que
  - se  $A < c$  então  $Ac = \text{true}$   
caso contrário  $Ac = \text{false}$

Exemplo:

Temperatura: 40 48 60 72 80 90

PlayTennis: No No Yes Yes Yes No

Qual **valor de  $c$**  escolher?

- O **valor de c** deveria ser escolhido de modo a produzir o maior ganho de informação.
  - Fayyad (1991) mostrou que o **valor de c** que maximiza o ganho de informação fica entre os limites de mudança do atributo.

Exemplo:

PlayTennys muda :  $(48+60)/2$  --- Temp  $>_{54}$

$(80+90)/2$  --- Temp  $>_{85}$

- Atributos candidatos:  $Temp_{>54}$   $Temp_{>85}$
- Calculado o ganho de informação para cada atributo é selecionamos o melhor:  
 $Temp_{>54}$ .

Este atributo booleano criado pode então competir com outros atributos candidatos discretos para o crescimento da árvores

# Exercícios

I - Construa árvores de decisão para as seguintes funções booleanas:

(a)  $A \wedge \neg B$

(b)  $A \vee [B \wedge C]$

(c)  $A \text{ XOR } B$

(d)  $[A \wedge B] \vee [C \wedge D]$

II - Seja os exemplos de treinamento:

		a1	a2
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T

- (a) Qual é a entropia desta coleção de exemplos de treinamento com a função objetivo de classificação.
- (b) Qual é o ganho e informação de  $a_2$  relativo aos exemplos de treinamento? II - Seja os exemplos de treinamento: