

Árvores de Decisão

Sistemas Inteligentes

Exemplos de situações do dia a dia em que a aprendizagem de máquina é importante

- ✉ A partir de informações sobre pacientes relativas a gravidez aprender a prever classes de futuros pacientes de alto risco que devem fazer cesárea
- ✉ Análise de risco de crédito: prever clientes mal pagadores
- ✉ Prever comportamento de compra de clientes
- ✉ Recomendar filmes para clientes
- ✉ etc

Multidisciplinaridade da Aprendizagem de Máquina

- ✉ Inteligência Artificial
- ✉ Estatística
- ✉ Teoria da Informação
- ✉ Teoria de Controle
- ✉ Filosofia
- ✉ Psicologia
- ✉ Neurobiologia
- ✉ ...

Problemas de Aprendizagem

- ✉ Melhorar a realização de uma tarefa a partir da experiência
 - Melhorar a realização da tarefa T
 - Em relação a uma medida de desempenho P
 - Baseada na experiência E

Problemas de Aprendizagem

✉ Exemplo: reconhecimento de caracteres manuscritos

- Tarefa T: reconhecer e classificar caracteres manuscritos
- Medida de desempenho P: percentagem de caracteres classificados corretamente
- Experiência a partir de treinamento E: base de dados de caracteres manuscritos com a respectiva classificação

Problemas de Aprendizagem

- ✉ O que é experiência adquirida a partir de treinamento?
 - É fornecida direta ou indiretamente?
 - É ensinada ou não por um professor
 - Problema: a experiência adquirida é suficiente para alcançar o desempenho desejado?
- ✉ Exatamente o que deve ser aprendido?
 - Aproximação de funções
 - Tipo de funções alvo: booleana, real, ...

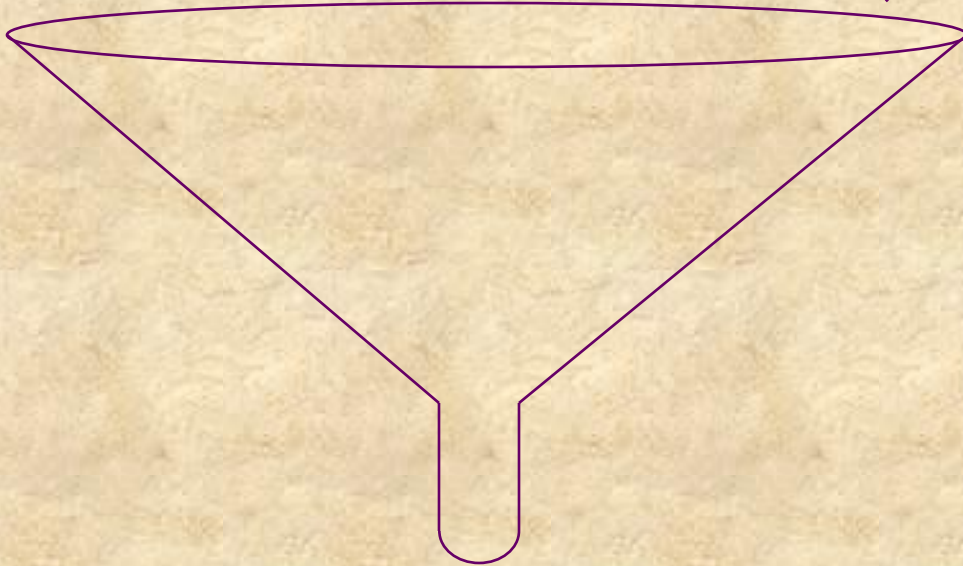
Problemas de Aprendizagem

- ✉ Como o que é aprendido deve ser representado?
 - Coleção de regras? Função polinomial? árvore de decisão?
- ✉ Qual o mecanismo de aprendizagem?
 - Qual o algoritmo de aprendizagem que deve ser usado?

Objetivo da aprendizagem

Conhecimento em extensão

(exemplos percepção-ação,
características-conceitos, etc.)



Conhecimento em intenção

(regras definições.)

Exemplos

dia 29, a Caxangá estava
engarrafada

dia 30, a Caxangá estava
engarrafada

dia 01, a Caxangá estava
engarrafada

dia 03, a Caxangá estava
engarrafada

Hipótese indutiva

Todo dia, a Caxangá está
engarrafada

Aprendizagem indutiva

- ✉ Inferência de uma regra geral (hipótese) a partir de exemplos particulares
 - ex. trânsito na caxangá
- ✉ Precisão diretamente proporcional à quantidade de exemplos

Aprendizagem indutiva

✉ Pode ser

- **incremental**: atualiza hipótese a cada novo exemplo
 - ◆ mais flexível, situada... Porém a ordem de apresentação é importante (backtracking)
- **não incremental**: gerada a partir de todo conjunto de exemplos
 - ◆ mais eficiente e prática

Uma Abordagem típicas em aprendizagem simbólica

✉ Árvores de decisão: inductive decision trees (ID3)

- Lógica de ordem 0+ (Instâncias (exemplos) são representadas por pares atributo-valor)
- Fáceis de serem implementadas e utilizadas
- aprendizagem não incremental
- estatística (admite exceções)

Árvores de Decisão

✉ Uma árvore de decisão utiliza uma estratégia de *dividir-para-conquistar*:

- Um problema complexo é decomposto em sub-problemas mais simples.
- Recursivamente a mesma estratégia é aplicada a cada sub-problema.

✉ A capacidade de discriminação de uma árvore vem da:

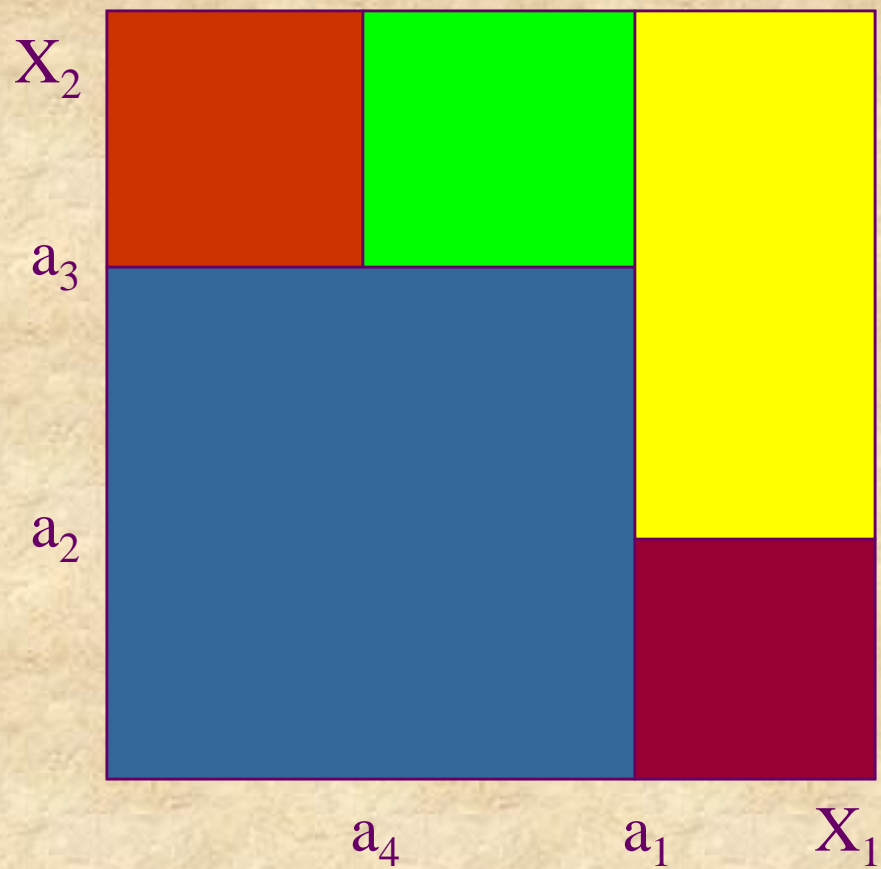
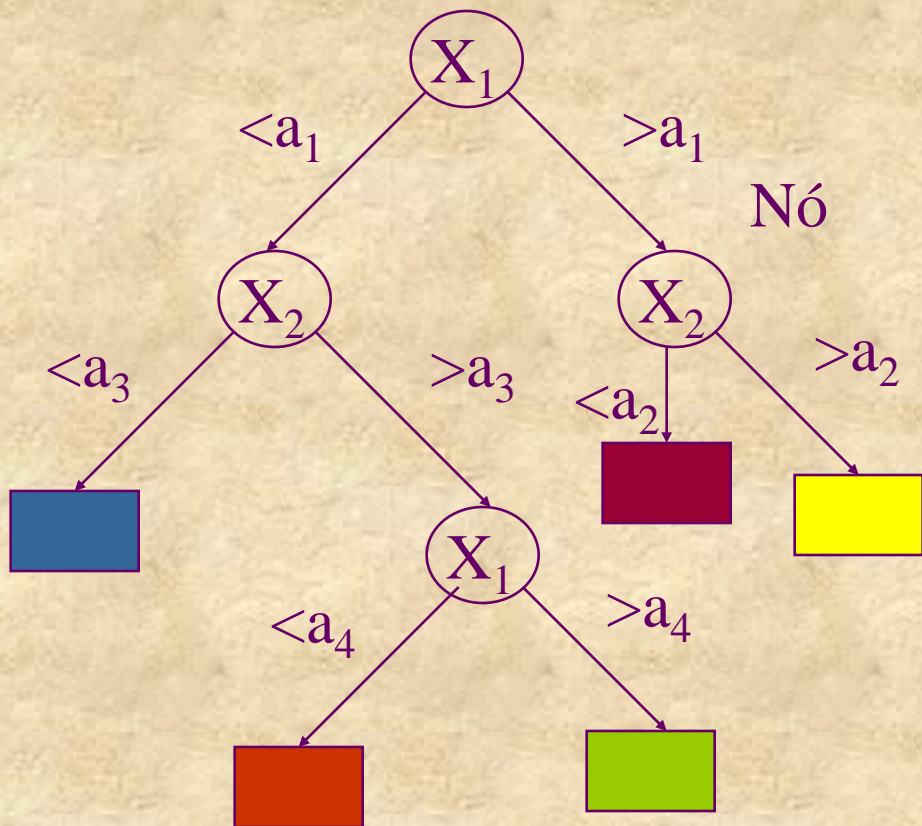
- Divisão do espaço definido pelos atributos em sub-espacos.
- A cada sub-espaco é associada uma classe.

Árvores de Decisão

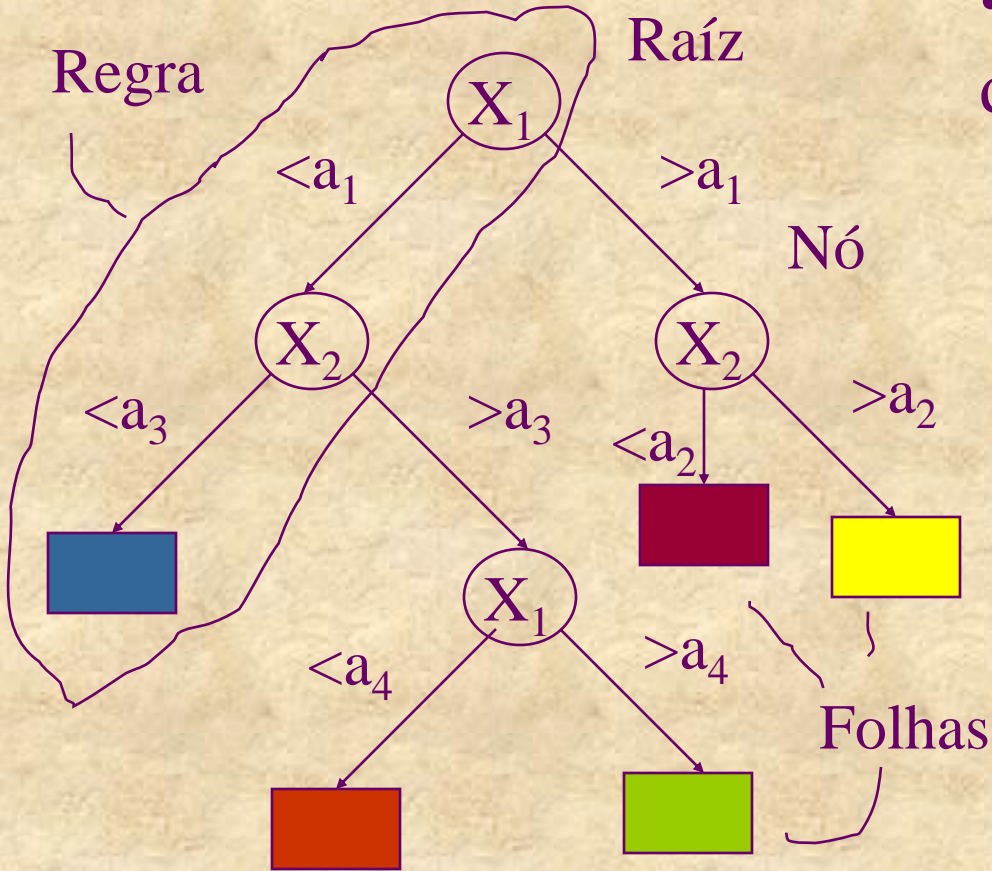
✉ Crescente interesse

- CART (Breiman, Friedman, et.al.)
- C4.5 (Quinlan)
- S plus , Statistica, SPSS, SAS

Árvores de Decisão



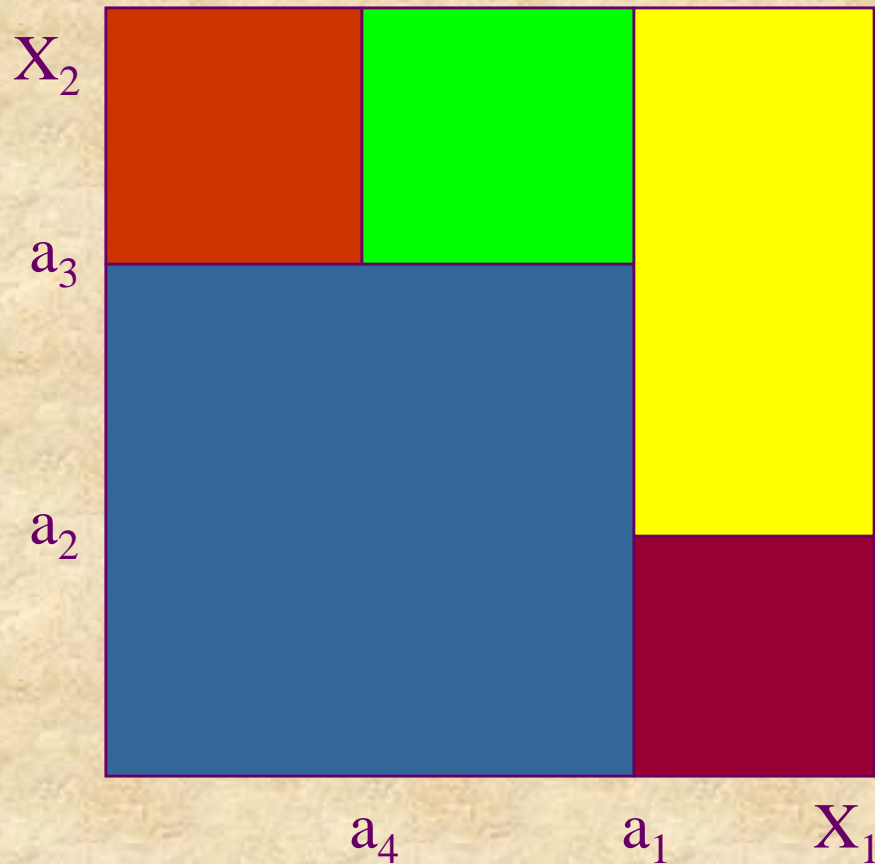
O que é uma Árvore de Decisão



- Representação por árvores de decisão:

- Cada nó de decisão contém um teste num atributo.
- Cada ramo descendente corresponde a um possível valor deste atributo.
- Cada Folha está associada a uma classe.
- Cada percurso na árvore (da raiz à folha) corresponde a uma regra de classificação.

Árvores de Decisão



- No espaço definido pelos atributos:

- Cada folha corresponde a uma região: Hiper-retângulo

- A intersecção dos hiper-retângulos é vazia

- A união dos hiper-retângulos é o espaço completo

Quando usar árvores de decisão?

- Instâncias (exemplos) são representadas por pares atributo-valor
- Função objetivo assume apenas valores discretos
- Hipóteses disjuntivas podem ser necessárias
- Conjunto de treinamento possivelmente corrompido por ruído
- Exemplos:
 - Diagnóstico médico, diagnóstico de equipamentos, análise de crédito

Construção de uma Árvore de Decisão

- A idéia *base*:
 1. Escolher um atributo.
 2. Estender a árvore adicionando um ramo para cada valor do atributo.
 3. Passar os exemplos para as folhas (tendo em conta o valor do atributo escolhido)

Construção de uma Árvore de Decisão

4. Para cada folha

1. Se todos os exemplos são da mesma classe, associar essa classe à folha
2. Senão repetir os passos 1 a 4

Exemplo

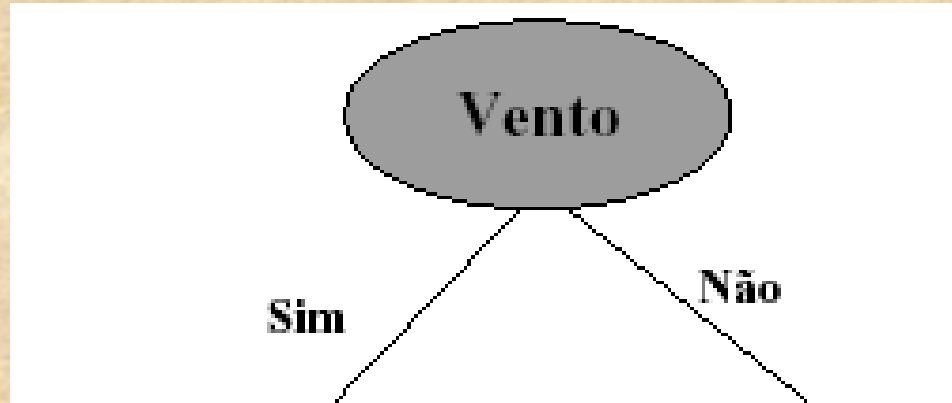
O conjunto de dados original

Tempo	Temperatura	Humidade	de vento	Joga
Sol	85	85	Não	Não
Sol	80	90	Sim	Não
Nublado	83	86	Não	Sim
Chuva	70	96	Não	Sim
Chuva	68	80	Não	Sim
Chuva	65	70	Sim	Não
Nublado	64	65	Sim	Sim
Sol	72	95	Não	Não
Sol	69	70	Não	Sim
Chuva	75	80	Não	Sim
Sol	75	70	Sim	Sim
Nublado	72	90	Sim	Sim
Nublado	81	75	Não	Sim
Chuva	71	91	Sim	Não

Exemplo

Seleciona um atributo

Qual o melhor atributo?



Tempo	Temperatura	Humidade	vento	Agua
Sol	80	90	Sim	Não
Chuva	80	70	Sim	Não
Nublado	80	80	Sim	Sim
Sol	70	70	Sim	Sim
Nublado	70	90	Sim	Sim
Chuva	70	90	Sim	Não

Tempo	Temperatura	Humidade	vento	Agua
Sol	80	80	Não	Não
Nublado	80	80	Não	Sim
Chuva	70	90	Não	Sim
Chuva	80	80	Não	Sim
Sol	70	80	Não	Não
Sol	80	70	Não	Sim
Chuva	70	80	Não	Sim
Nublado	80	70	Não	Sim

Critérios para Escolha do Atributo

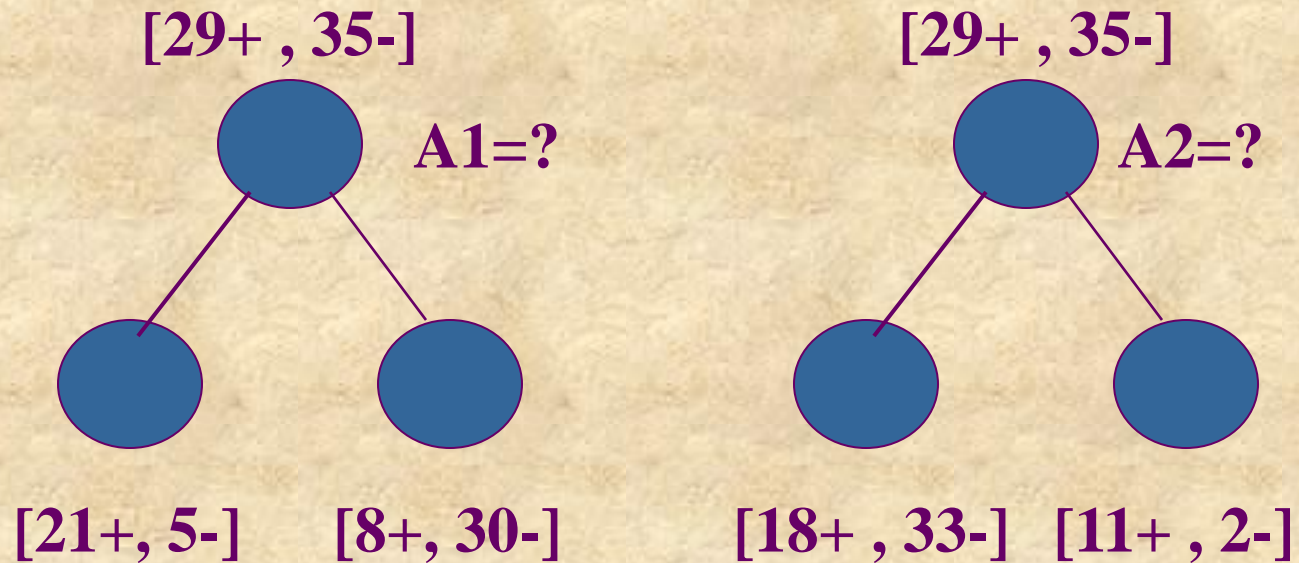
- Como medir a *habilidade* de um dado atributo discriminar as classes?
- Existem muitas medidas.

Todas concordam em dois pontos:

- Uma divisão que mantém as proporções de classes em todas as partições é inútil.
- Uma divisão onde em cada partição todos os exemplos são da mesma classe tem utilidade máxima.

Critérios para Escolha do Atributo

✉ Qual é o melhor atributo?



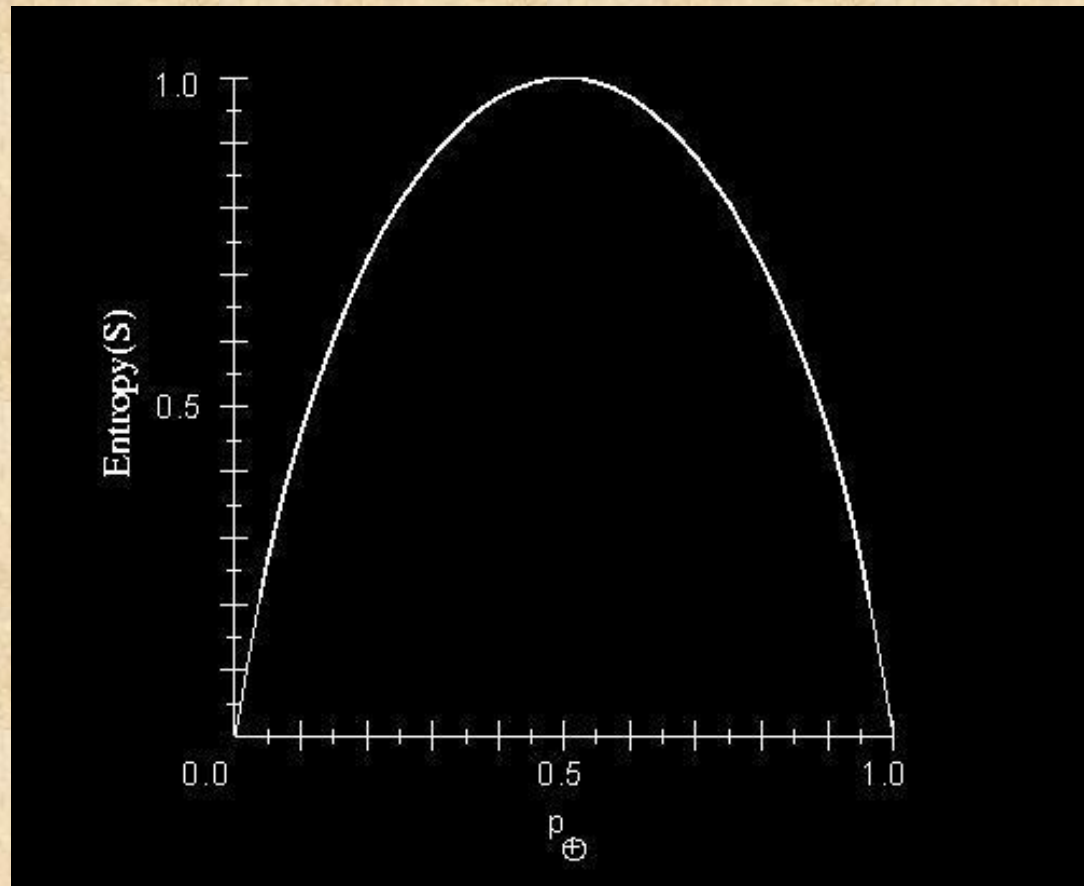
Entropia

- ✉ S é uma amostra dos exemplos de treinamento
- ✉ p_{\oplus} é a proporção de exemplos positivos em S
- ✉ p_{\ominus} é a proporção de exemplos negativos em S
- ✉ Entropia mede a "impureza" de S :
 - $Entropia(S) = -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$

Entropia - Exemplo I

- ✉ Se p_{\oplus} é 1, o destinatário sabe que o exemplo selecionado será positivo
 - Nenhuma mensagem precisa ser enviada
 - Entropia é 0 (mínima)
- ✉ Se p_{\oplus} é 0.5, um bit é necessário para indicar se o exemplo selecionado é \oplus ou \ominus
 - Entropia é 1 (máxima)

Entropia - Gráfico



Entropia

- Entropia é uma medida da aleatoriedade (impureza) de uma variável.

- A entropia de uma variável nominal X que pode tomar i valores:

$$\text{entropia}(X) = -\sum_i p_i \log_2 p_i$$

- A entropia tem máximo ($\log_2 i$) se $p_i = p_j$ para qualquer $i \neq j$

- A entropia(x) = 0 se existe um i tal que $p_i = 1$

- É assumido que $0 * \log_2 0 = 0$

Entropia - Exemplo II

- ✉ Suponha que S é uma coleção de 14 exemplos, incluindo 9 positivos e 5 negativos
 - Notação: [9+,5-]
- ✉ A entropia de S em relação a esta classificação booleana é dada por:

$$\begin{aligned} \text{Entropy}([9+,5-]) &= -(9/14)\log_2(9/14) - (5/14)\log_2(5/14) \\ &= 0.940 \end{aligned}$$

Ganho de Informação

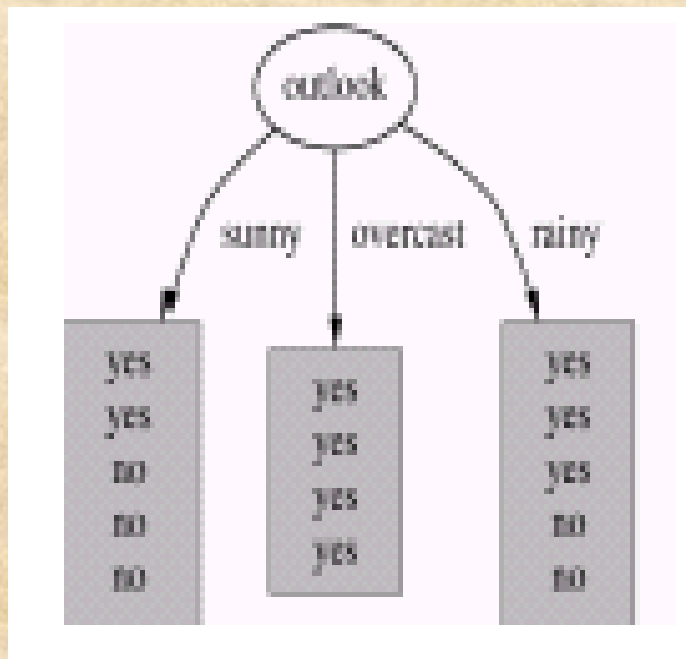
- No contexto das árvores de decisão a entropia é usada para estimar a aleatoriedade da variável a prever (classe).
- Dado um conjunto de exemplos, que atributo escolher para teste?
 - Os valores de um atributo definem partições do conjunto de exemplos.
 - O ganho de informação mede a redução da entropia causada pela partição dos exemplos de acordo com os valores do atributo.

Ganho de Informação

$$\text{ganho}(Exs, Atri) = \text{entropia}(Exs) - \sum_v \frac{\# Exs_v}{\# Exs} \text{entropia}(Exs_v)$$

A construção de uma árvore de decisão é guiada pelo objetivo de diminuir a entropia ou seja a aleatoriedade - dificuldade de previsão- da variável que define as classes.

Cálculo do Ganho de Informação de um atributo nominal



- ✉ • Informação da Classe:
 - $p(\text{sim}) = 9/14$
 - $p(\text{não}) = 5/14$
 - $\text{Ent}(\text{joga}) = - 9/14 \log_2 9/14 - 5/14 \log_2 5/14 = 0.940$
- ✉ • Informação nas partições:
 - $p(\text{sim} \mid \text{tempo}=\text{sol}) = 2/5$
 - $p(\text{não} \mid \text{tempo}=\text{sol}) = 3/5$

Cálculo do Ganho de Informação de um atributo nominal

✉ Informação nas partições:

- $\text{Ent}(\text{joga} \mid \text{tempo}=\text{sol})$
- $= -2/5 \log_2 2/5 - 3/5 \log_2 3/5 = 0.971$
- $\text{Ent}(\text{joga} \mid \text{tempo}=\text{nublado}) = 0.0$
- $\text{Ent}(\text{joga} \mid \text{tempo}=\text{chuva}) = 0.971$
- $\text{Info}(\text{tempo}) = 5/14 * 0.971 + 4/14 * 0 + 5/14 * 0.971 = 0.693$

	Sol	Nublado	Chuva
Sim	2	4	3
Não	3	0	2

✉ Ganho de Informação obtida neste atributo:

- $\text{Ganho}(\text{tempo}) = \text{Ent}(\text{joga}) - \text{Info}(\text{tempo})$
- $\text{Ganho}(\text{tempo}) = 0.940 - 0.693 = 0.247$

Ganho (vento)

Values (Wind) = Weak, Strong

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

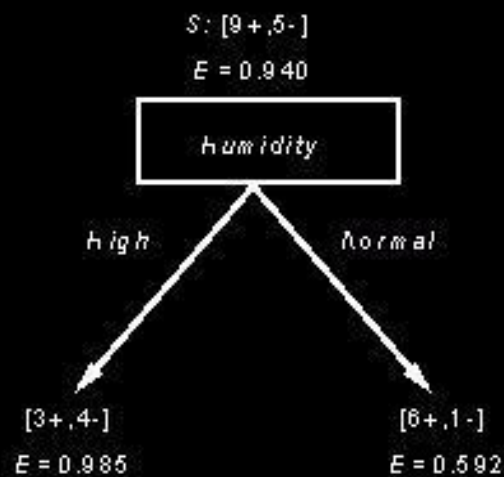
$$S_{Weak} = [6+, 2-]$$

$$S_{Strong} = [3+, 3-]$$

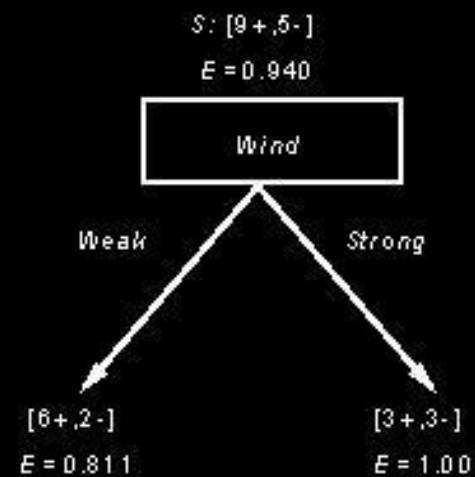
$$\begin{aligned} Gain(S, Wind) &= Entropy(S) - \sum_{v \in \{Weak, Strong\}} \frac{|S_v|}{|S|} Entropy(S_v) \\ &= Entropy(S) - (8/14) Entropy(S_{Weak}) - (6/14) Entropy(S_{Strong}) \\ &= 0.940 - (8/14) * 0.811 - (6/14) * 1.00 \\ &= 0.048 \end{aligned}$$

Critério de ganho

Which attribute is the best classifier?



$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= .940 - (7/14) \cdot 0.985 - (7/14) \cdot 0.592 \\ &= .151 \end{aligned}$$



$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= .940 - (8/14) \cdot 0.811 - (6/14) \cdot 1.0 \\ &= .048 \end{aligned}$$

Exemplos de treinamento

Considere a tarefa de aprendizagem representada pelos exemplos de treinamento na tabela abaixo, onde o objetivo é prever o atributo *PlayTennis* baseando-se nos outros atributos

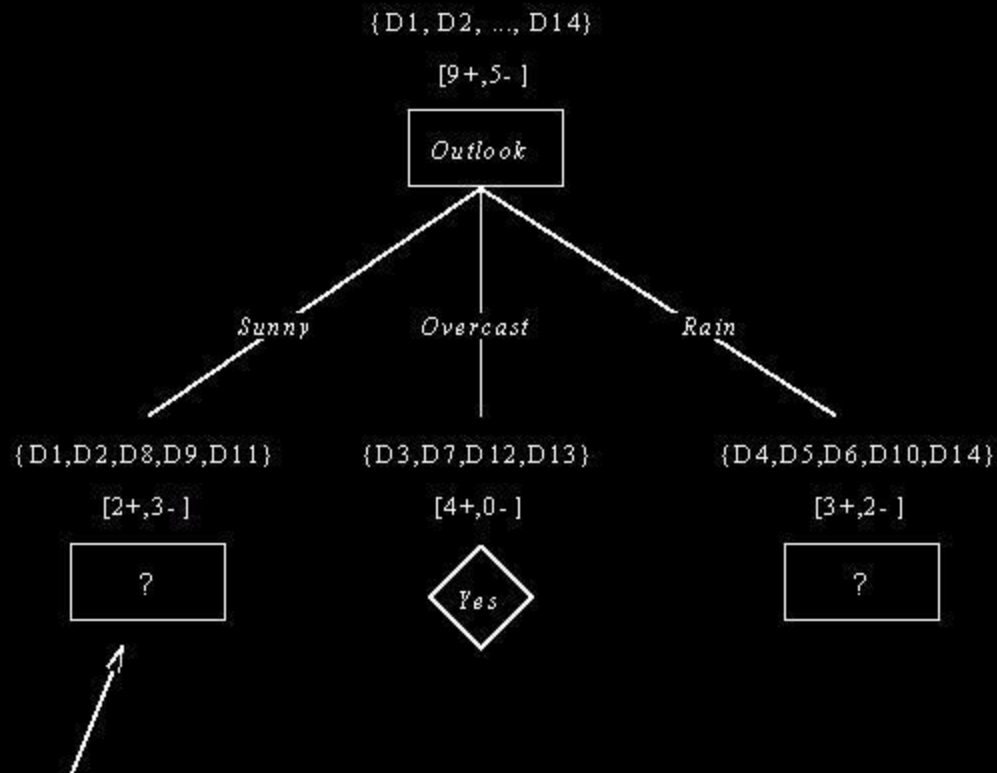
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Exemplos de treinamento

✉ Que atributo deve ser selecionado para ser a raiz da árvore?

- $Gain(S, Outlook) = 0.247$
- $Gain(S, Humidity) = 0.151$
- $Gain(S, Wind) = 0.048$
- $Gain(S, Temperature) = 0.029$

✉ onde S denota a coleção de exemplos na tabela anterior



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

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temperature}) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$$

$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = .970 - (2/5) 1.0 - (3/5) .918 = .019$$

Cálculo do Ganho para Atributos Numéricos

- ✉ Um teste num atributo numérico produz uma partição binária do conjunto de exemplos:
 - Exemplos onde $\text{valor_do_atributo} < \text{ponto_referência}$
 - Exemplos onde $\text{valor_do_atributo} > \text{ponto_referência}$
- ✉ Escolha do ponto de referência:
 - Ordenar os exemplos por ordem crescente dos valores do atributo numérico.
 - Qualquer ponto intermediário entre dois valores diferentes e consecutivos dos valores observados no conjunto de treinamento pode ser utilizado como possível ponto de referência.

Cálculo do Ganho para Atributos Numéricos

- É usual considerar o valor médio entre dois valores diferentes e consecutivos.
- Fayyad e Irani (1993) mostram que de todos os possíveis pontos de referência aqueles que maximizam o ganho de informação separam dois exemplos de classes diferentes.

Cálculo do Ganho para Atributos Numéricos

Temperatu.	Joga
64	Sim
65	Não
68	Sim
69	Sim
70	Sim
71	Não
72	Não
72	Sim
75	Sim
75	Sim
80	Não
81	Sim
83	Sim
85	Não

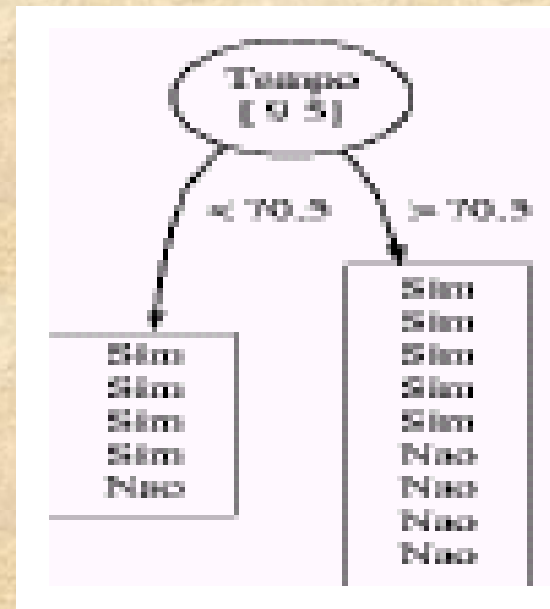
- ✉ Considere o ponto de referência temperatura = 70.5
- ✉ Um teste usando este ponto de referência divide os exemplos em duas classes:
 - Exemplos onde temperatura < 70.5
 - Exemplos onde temperatura > 70.5
- ✉ Como medir o ganho de informação desta partição?

Cálculo do Ganho para Atributos Numéricos

✉ Como medir o ganho de informação desta partição?

✉ Informação nas partições

- $p(\text{sim} \mid \text{temperatura} < 70.5) = 4/5$
- $p(\text{não} \mid \text{temperatura} < 70.5) = 1/5$
- $p(\text{sim} \mid \text{temperatura} > 70.5) = 5/9$
- $p(\text{não} \mid \text{temperatura} > 70.5) = 4/9$



Cálculo do Ganho para Atributos Numéricos

- $\text{Info}(\text{joga} \mid \text{temperatura} < 70.5) = -4/5 \log_2 4/5 - 1/5 \log_2 1/5 = 0.721$
- $\text{Info}(\text{joga} \mid \text{temperatura} > 70.5) = -5/9 \log_2 5/9 - 4/9 \log_2 4/9 = 0.991$
- $\text{Info}(\text{temperatura}) = 5/14 * 0.721 + 9/14 * 0.991 = 0.895$
- $\text{Ganho}(\text{temperatura}) = 0.940 - 0.895 = 0.045 \text{ bits}$

Critérios de Parada

- ✉ Quando parar a divisão dos exemplos?
- Todos os exemplos pertencem a mesma classe.
 - Todos os exemplos têm os mesmos valores dos atributos (mas diferentes classes).
 - O número de exemplos é inferior a um certo limite.
 - O mérito de todos os possíveis testes de partição dos exemplos é muito baixo.

Construção de uma Árvore de Decisão

✉ Input: Um conjunto de exemplos

✉ Output: Uma árvore de decisão

✉ Função Geraárvore(Exs)

- Se $\text{criterio_parada}(\text{Exs}) = \text{TRUE}$: retorna Folha
- Escolhe o atributo que maximiza o $\text{critério_divisão}(\text{Exs})$
- Para cada partição i dos exemplos baseada no atributo escolhido: $\text{árvore}_i = \text{Geraárvore}(\text{Exs}_i)$
- Retorna um nó de decisão baseado no atributo escolhido e com descendentes árvore_i .
- Fim

Construção de uma Árvore de Decisão

✉ O problema de construir uma árvore de decisão:

- Consistente com um conjunto de exemplos
- Com o menor número de nós
- É um problema *NP* completo.

✉ Dois problemas:

- Que atributo selecionar para teste num nó?
- Quando parar a divisão dos exemplos ?

Construção de uma Árvore de Decisão

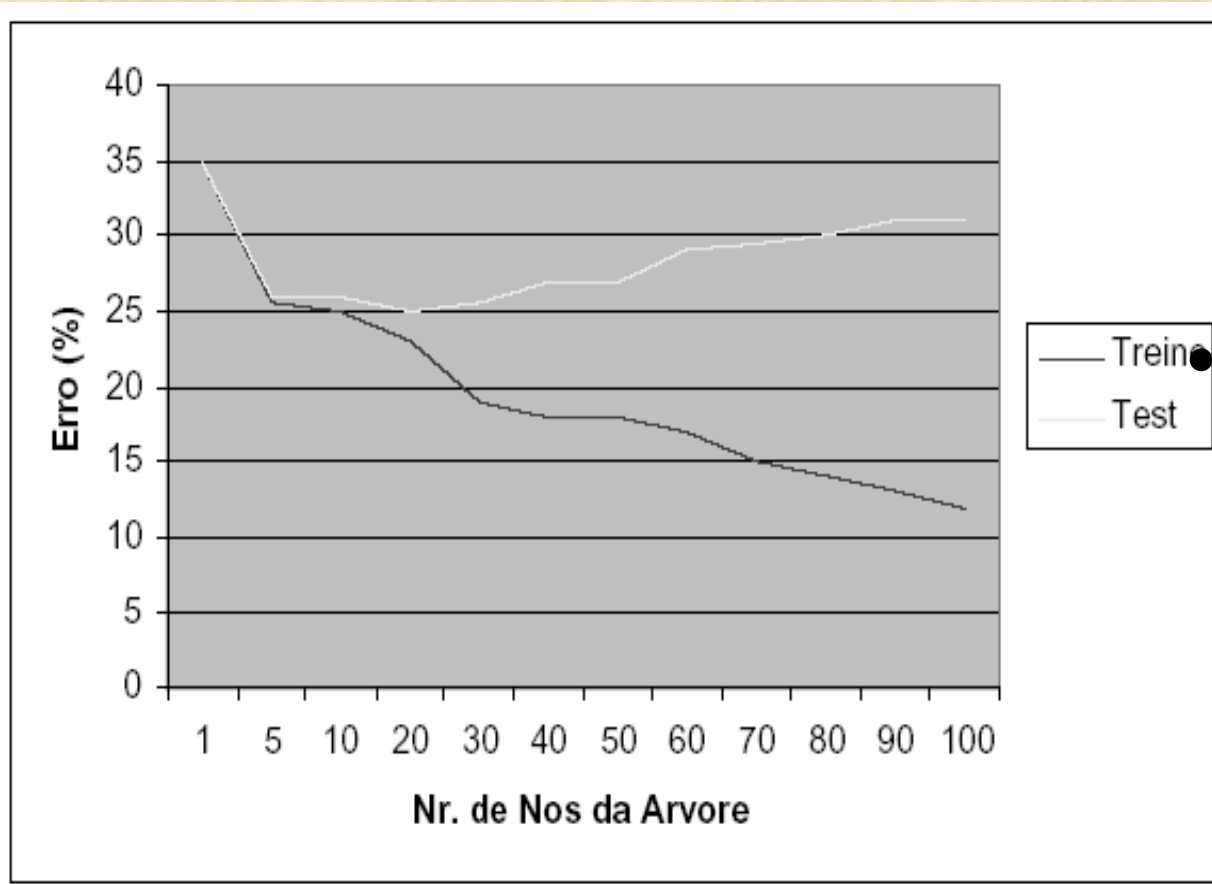
✉ Os algoritmos mais populares:

- Utilizam heurísticas que tomam decisões olhando para a frente um passo.
- Não reconsideram as opções tomadas
 - ◆ Não há backtracking
 - ◆ Mínimo local

Sobre-ajustamento (Overfitting)

- ✉ O algoritmo de partição recursiva do conjunto de dados gera estruturas que podem obter um ajuste aos exemplos de treinamento perfeito.
 - Em domínios sem ruído o nr. de erros no conjunto de treinamento pode ser 0.
- ✉ Em problemas com *ruído* esta capacidade é problemática:
 - A partir de uma certa profundidade as decisões tomadas são baseadas em pequenos conjuntos de exemplos.
 - A capacidade de generalização para exemplos não utilizados no crescimento da árvore diminui.

Variação do erro com o nr. de nós



Sobre-ajustamento (“*overfitting*”)

✉ Definição:

- Uma árvore de decisão d faz sobre-ajustamento aos dados se existir uma árvore d' tal que:
 d tem menor erro que d' no conjunto de treinamento mas d' tem menor erro na população.

✉ Como pode acontecer:

- Ruído nos dados;

✉ O número de parâmetros de uma árvore de decisão cresce linearmente com o número de exemplos.

- Uma árvore de decisão pode obter um ajuste perfeito aos dados de treinamento.

Sobre-ajustamento (“*overfitting*”)

✉ Occam’s razor: preferência pela hipótese mais simples.

- Existem menos hipóteses simples do que complexas.
- Se uma hipótese simples explica os dados é pouco provável que seja uma coincidência.
- Uma hipótese complexa pode explicar os dados apenas por coincidência.

Simplificar a árvore

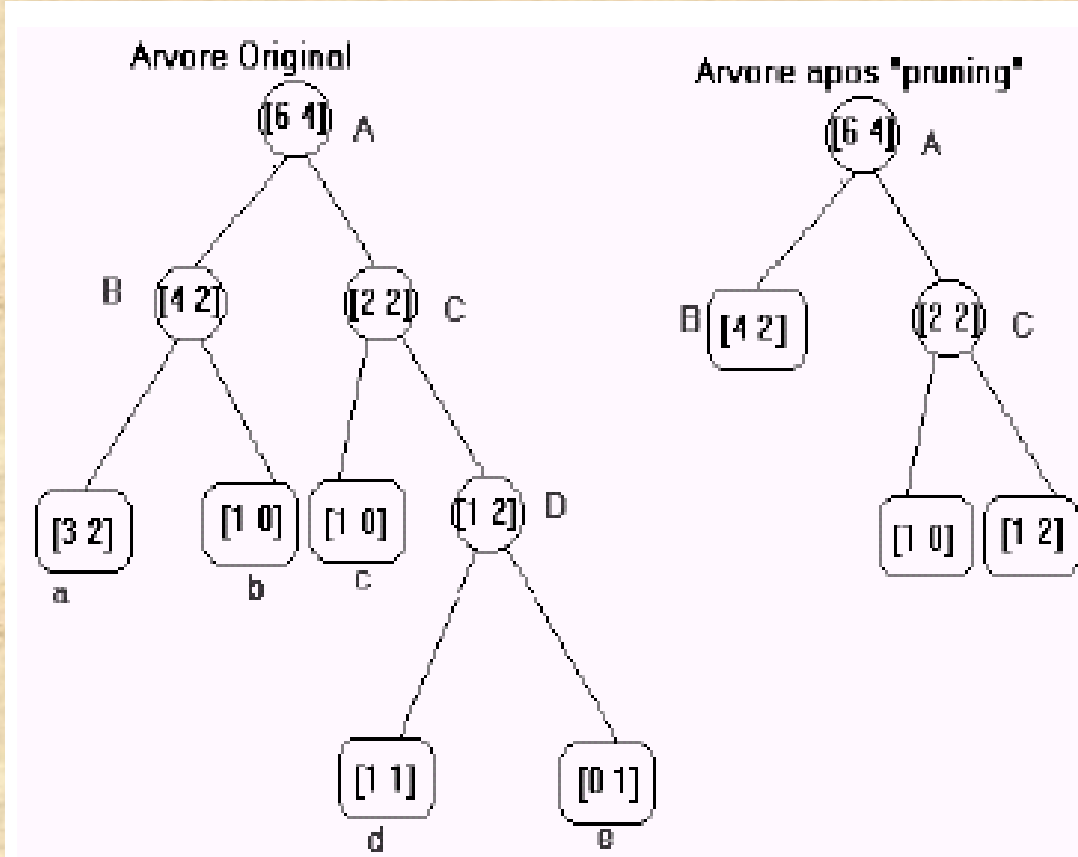
✉ Duas possibilidades:

- Parar o crescimento da árvore mais cedo (pre-pruning).
- Construir uma árvore completa e podar a árvore (pos-pruning).
- “*Growing and pruning is slower but more reliable*”
 - ◆ Quinlan, 1988

Um algoritmo básico de pruning

- ✉ Percorre a árvore em profundidade
- ✉ Para cada nó de decisão calcula:
 - Erro no nó
 - Soma dos erros nos nós descendentes
- ✉ Se o erro no nó é menor ou igual à soma dos erros dos nós descendentes, o nó é transformado em folha.

Um algoritmo básico de pruning



✉ Exemplo do nó B:

- Erro no nó = 2
- Soma dos erros nos nós descendentes:
 - ◆ $2 + 0$
- Transforma o nó em folha
 - ◆ Elimina os nós descendentes.

Critérios de como escolher a melhor árvore.

- ✉ Obter estimativas confiáveis do erro a partir do conjunto de treinamento.
- ✉ Otimizar o erro num conjunto de validação independente do utilizado para construir a árvore.
- ✉ Minimizar:
 - *erro no treinamento + dimensão da árvore*
 - ◆ *Cost Complexity pruning (Cart)*
 - *dimensão da árvore + quantidade de exemplos mal classificados*
 - ◆ *MDL pruning (Quinlan)*

Estimativas de Erro

- ✉ O problema fundamental do algoritmo de poda é a estimativa de erro num determinado nó.
 - O erro estimado a partir do conjunto de treino não é um estimador confiável.
- ✉ O “*reduced error pruning*”
 - consiste em obter estimativas de erro a partir de um conjunto de validação independente do conjunto de treino.
 - Reduz o volume de informação disponível para crescer a árvore.

Valores de atributo desconhecidos

- ✉ E se valores do atributo A estão faltando para alguns exemplos?
 - Substituir o valor desconhecido durante o pré-processamento pelo valor mais provável (ex. média)
- ✉ Mesmo assim use os exemplos de treinamento, e organize a árvore como segue:
 - Se um nó n testa A , atribua um valor para A que seja o mais comum entre os outros exemplos classificados no nó n
 - Atribua para A um valor que seja o mais comum entre os outros exemplos com o mesmo valor objetivo (*target value*)

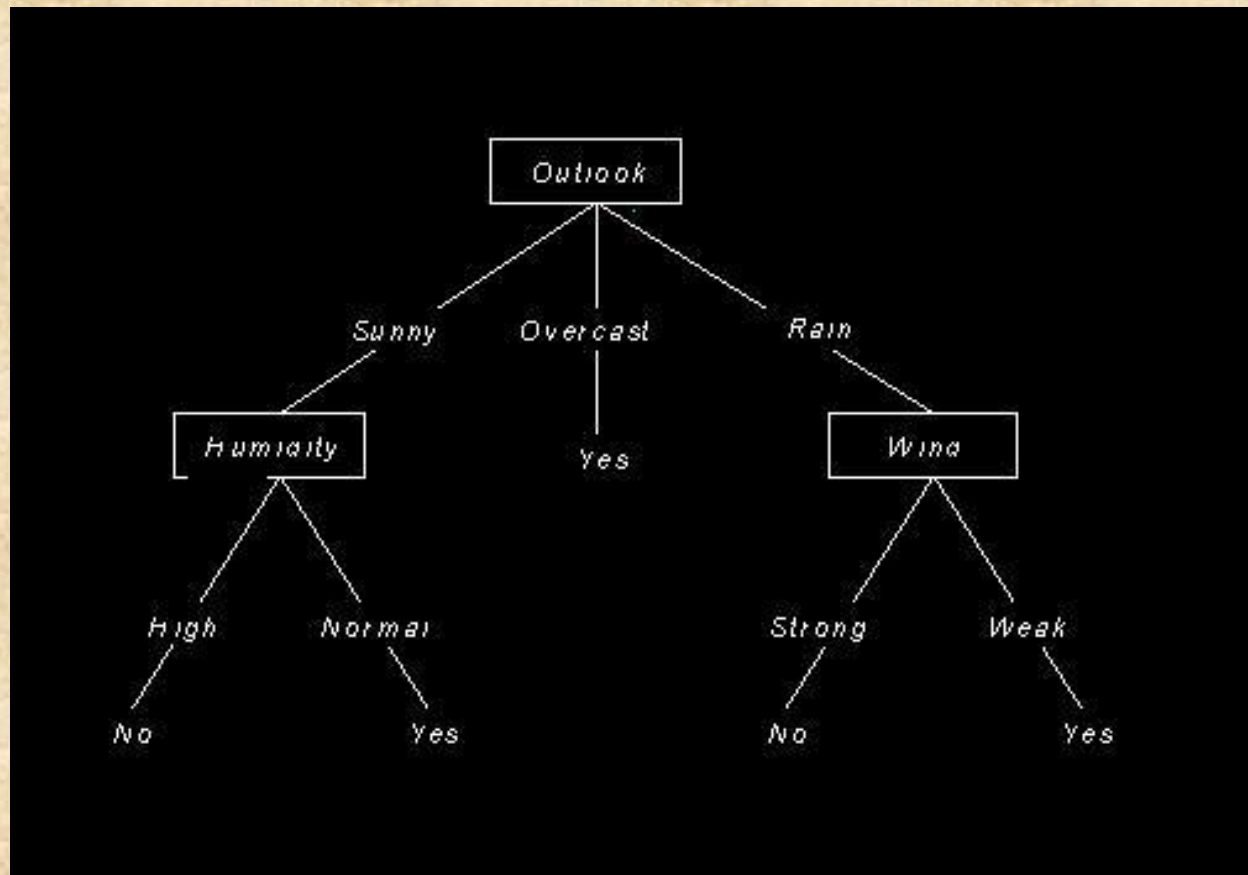
Valores de atributo desconhecidos

- Atribua uma probabilidade p_i para cada valor possível v_i de A
 - ◆ atribua uma fração p_i de exemplos para cada descendente da árvore

Transformação de árvores em regras de decisão

- ✉ Regras podem ser auto-interpretadas.
- ✉ Uma transformação:
 - Cada ramo dá origem a uma regra
 - ◆ A regra prediz a classe associada á folha
 - ◆ A parte condicional da regra é obtida pela conjunção das condições de cada nó.
- ✉ Em cada regra é testado a eliminação de condições. Uma condição é eliminada se:
 - O erro não aumenta
 - A estimativa de erro não aumenta

Convertendo uma árvore em regras



Convertendo uma árvore em regras

- ✉ IF (*Outlook = Sunny*) \wedge (*Humidity = High*) THEN
PlayTennis = No
- ✉ IF (*Outlook = Sunny*) \wedge (*Humidity = Normal*) THEN
PlayTennis = YES

.....

Porquê Regras ?

- ✉ Permite eliminar um teste numa regra, mas pode reter o teste em outra regra.
- ✉ Elimina a distinção entre testes perto da raiz e testes perto das folhas.
- ✉ Maior grau de interpretabilidade.

Referências

✉ Machine Learning. Tom Mitchell. McGraw-Hill.1997.