

Aprendizado de Máquina: Árvore de Decisão

Prof. Arnaldo Candido Junior
UTFPR – Medianeira

Visão geral

- **Árvores de decisão** (ADs)
 - Utilizam a estratégia de divisão e conquista
 - Recursivamente: dividir problemas difíceis em problemas mais simples
 - É uma das técnicas mais utilizadas

Visão geral ₍₂₎

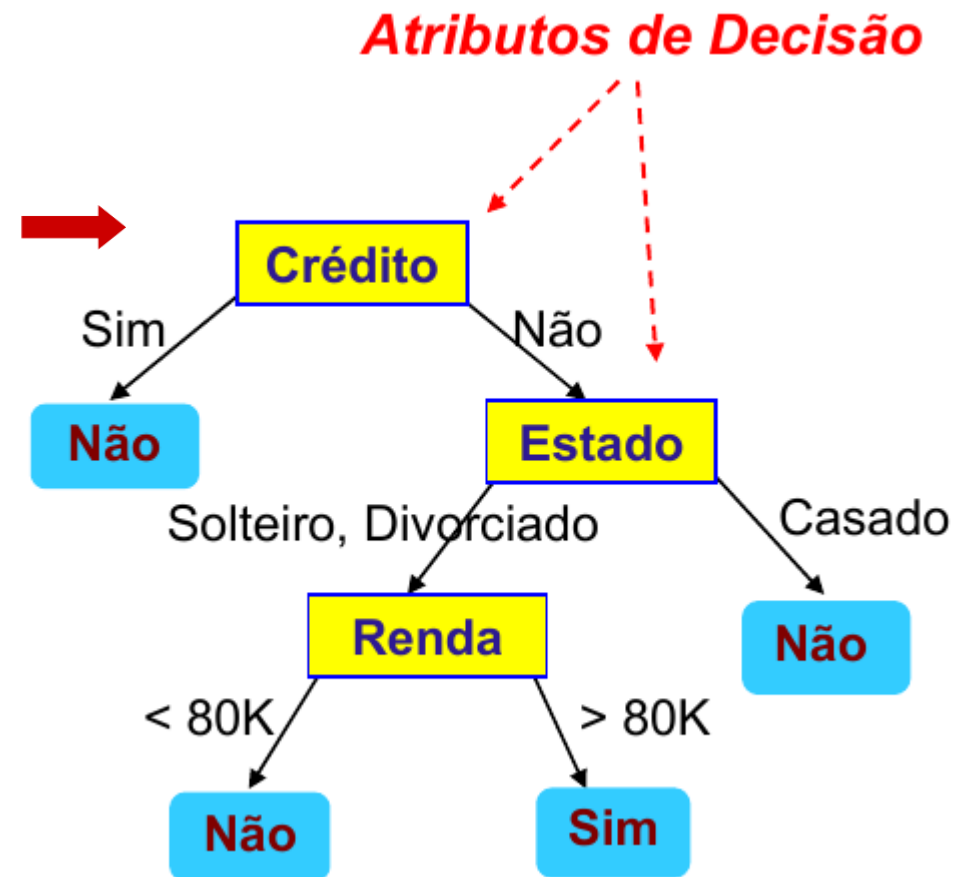
- Composta por:
 - Nó raiz: nenhum pai e com $n \geq 0$ filhos
 - Nós intermediários: um pai e $n \geq 2$ filhos
 - Nós folhas: um pai e nenhum filho

Exemplo: deve dinheiro?

Dados de treinamento

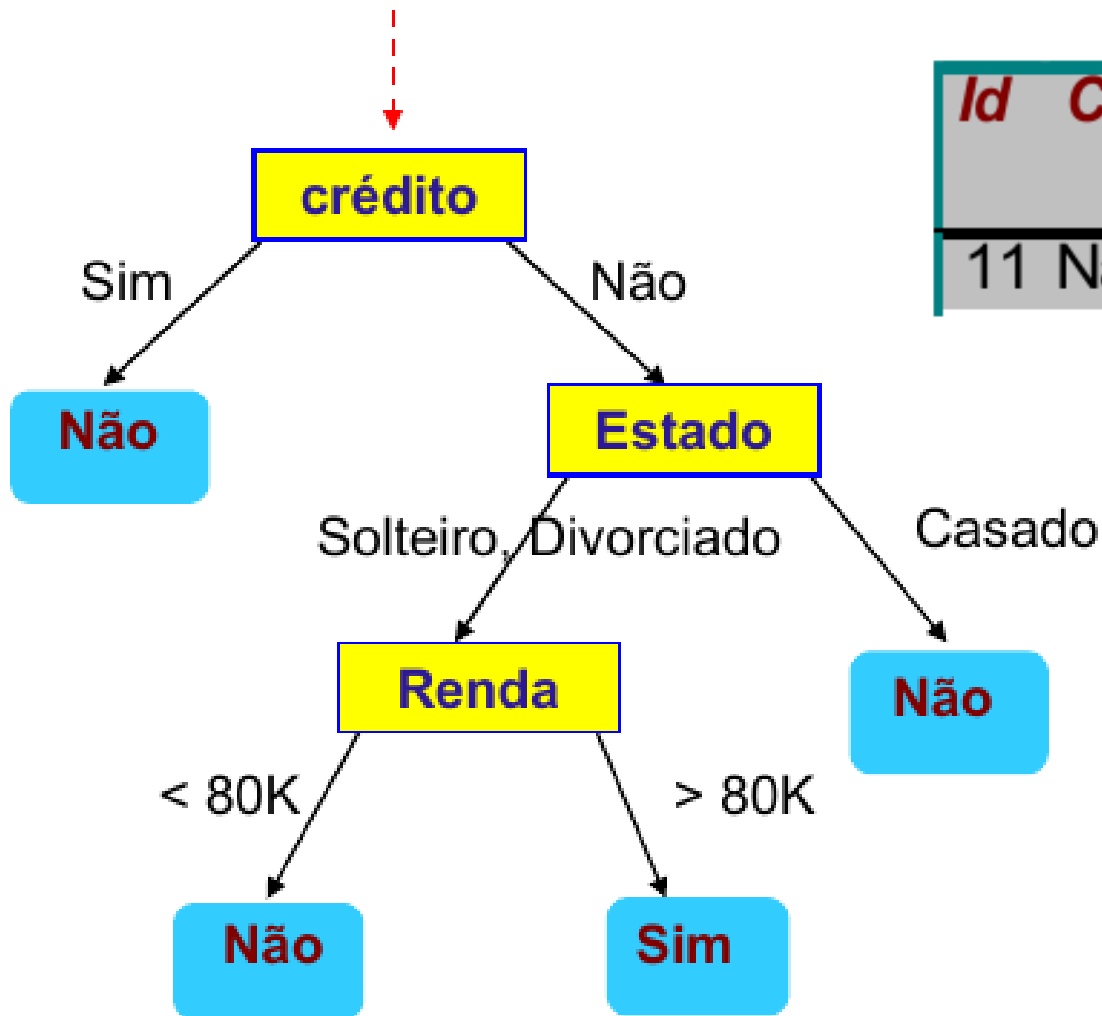
Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Modelo induzido



Exemplo: deve dinheiro ⁽²⁾

Começar da raiz

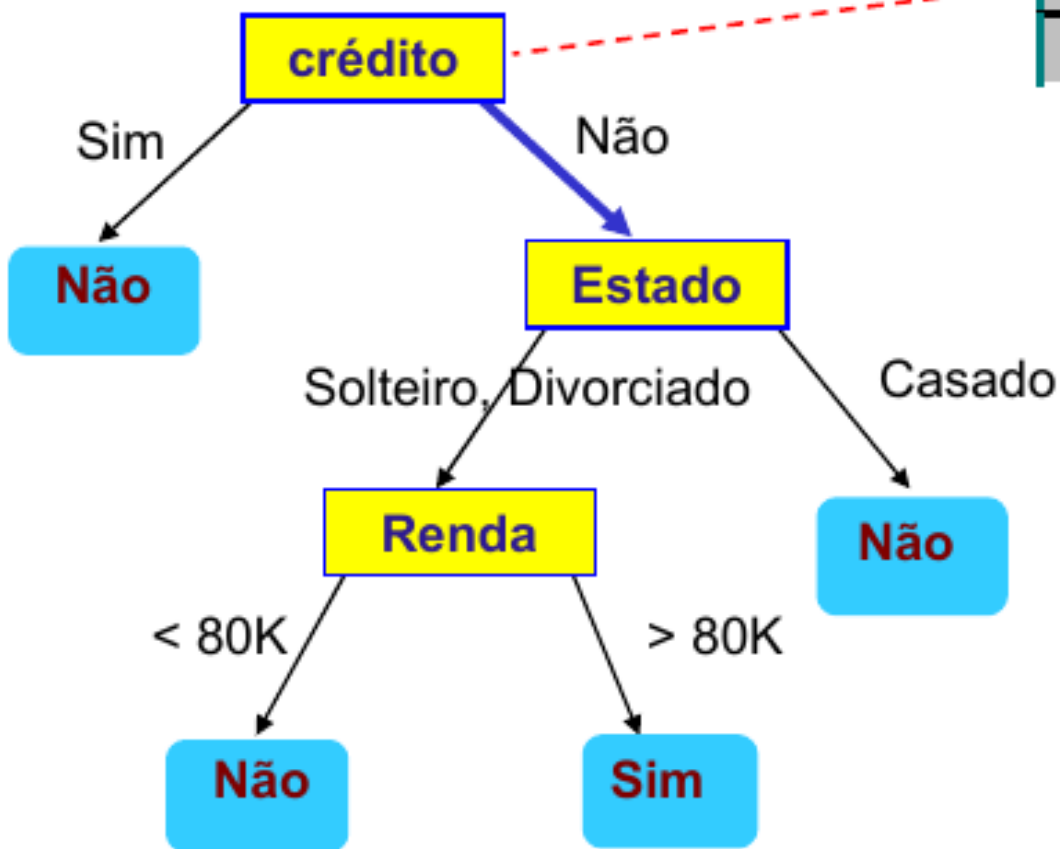


Dados de teste

<i>Id</i>	<i>Crédito</i>	<i>Estado Civil</i>	<i>Renda</i>	<i>Deve</i>
11	Não	Casado	80K	?

Exemplo: deve dinheiro ⁽³⁾

Dados de teste

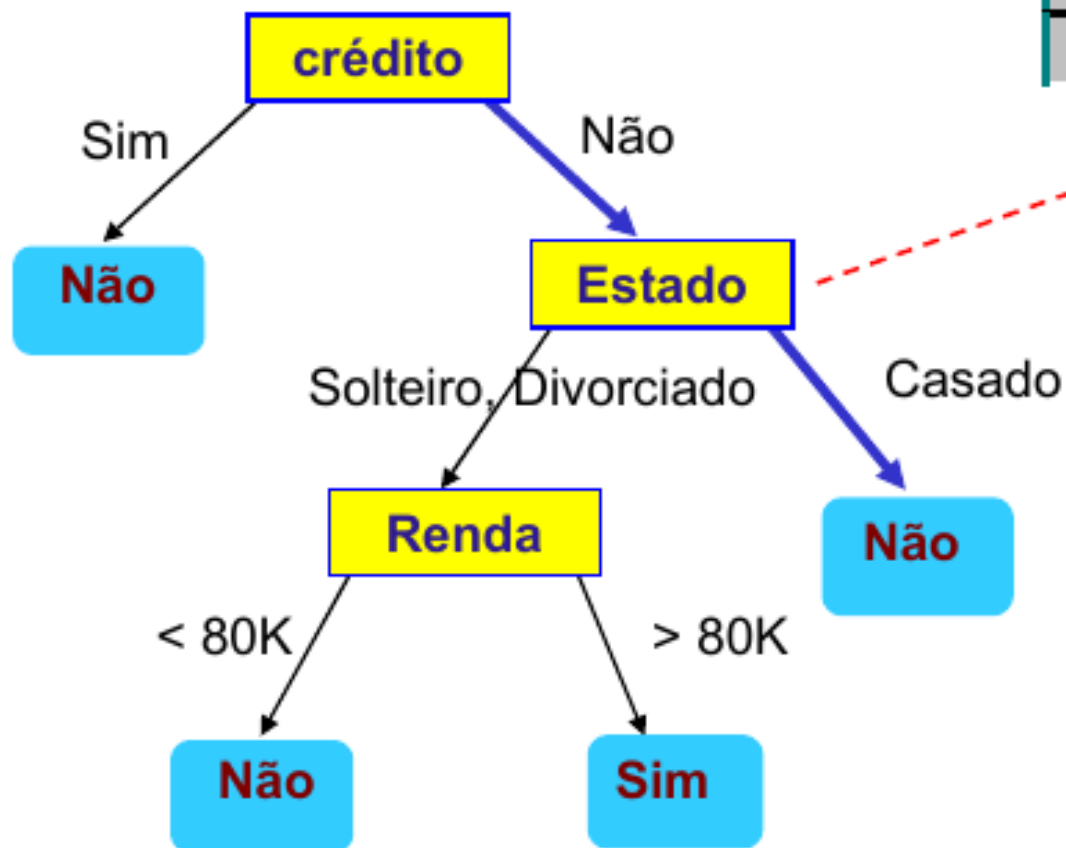


<i>Id</i>	<i>Crédito</i>	<i>Estado Civil</i>	<i>Renda</i>	<i>Deve</i>
11	Não	Casado	80K	?

Exemplo: deve dinheiro ₍₄₎

Dados de teste

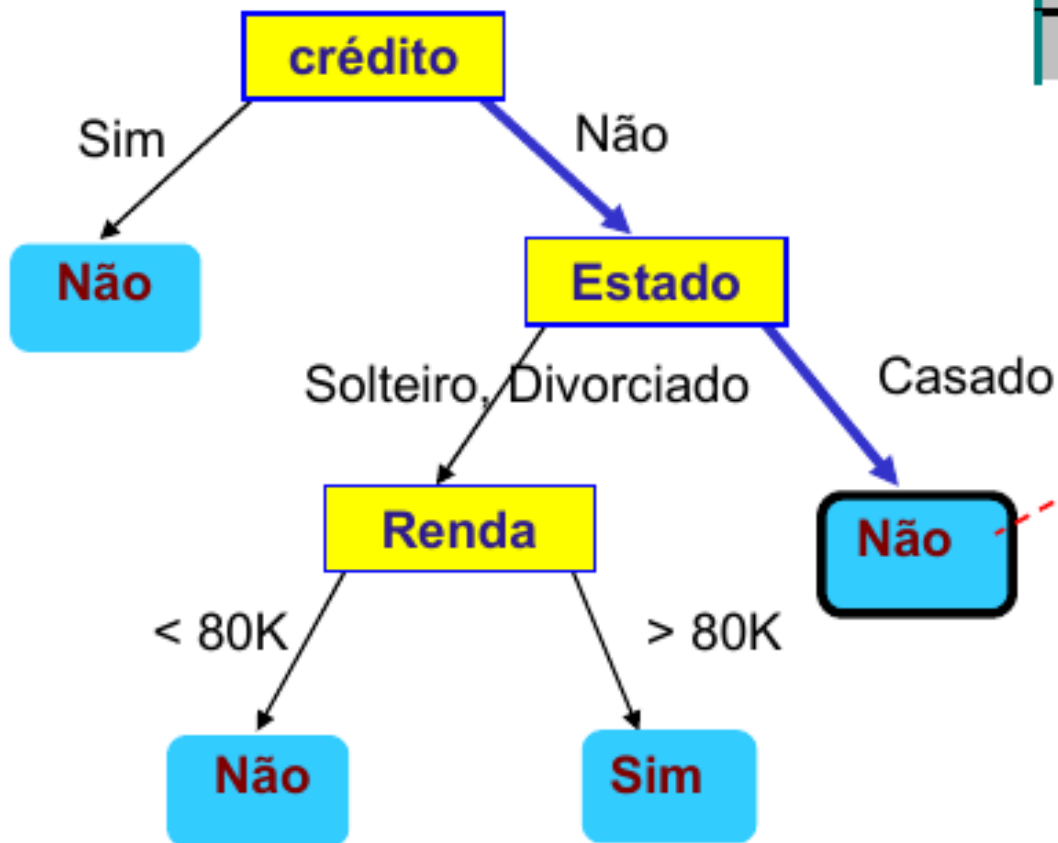
<i>Id</i>	<i>Crédito</i>	<i>Estado Civil</i>	<i>Renda</i>	<i>Deve</i>
11	Não	Casado	80K	?



Exemplo: deve dinheiro ⁽⁵⁾

Dados de teste

<i>Id</i>	<i>Crédito</i>	<i>Estado Civil</i>	<i>Renda</i>	<i>Deve</i>
11	Não	Casado	80K	?

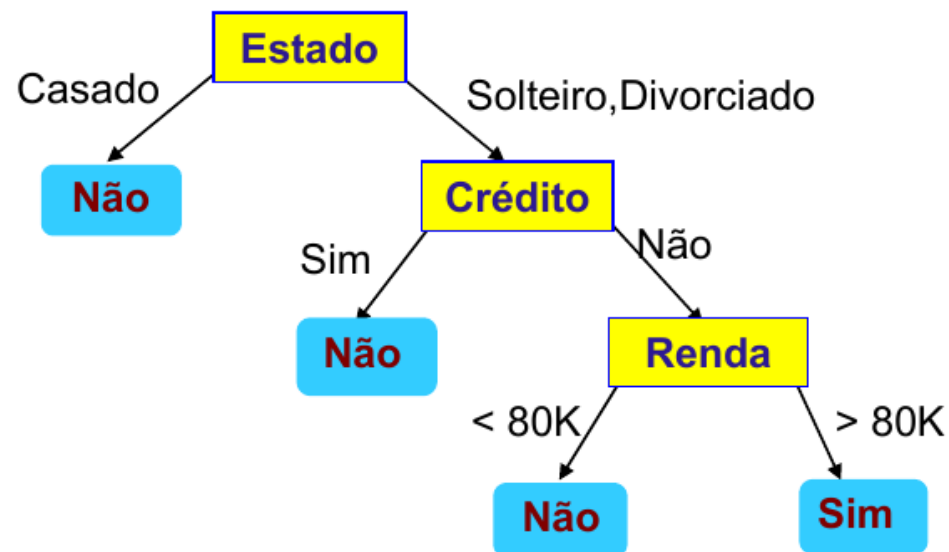


Atribui classe **não**

Exemplo: deve dinheiro ⁽⁶⁾

- Mais de uma árvore pode ser reajustada para os mesmos dados

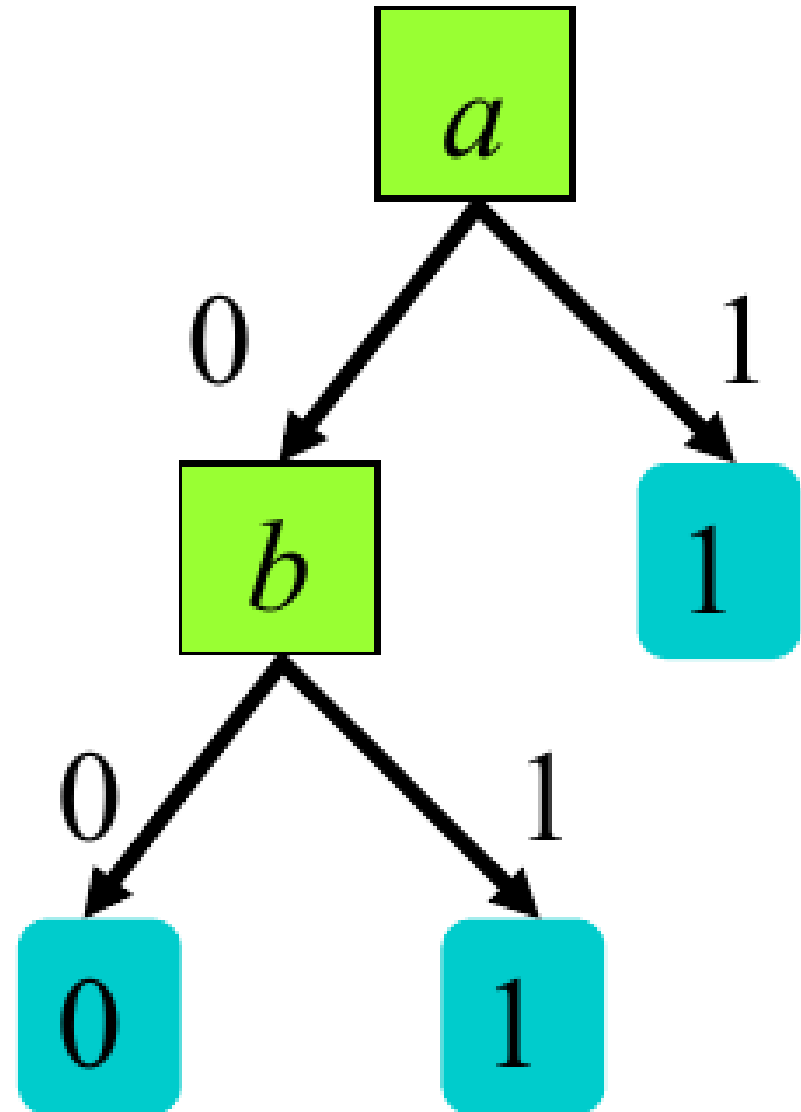
Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim



Exemplo: operações lógicas

$a \vee b$

a	b	$a \vee b$
0	0	0
0	1	1
1	0	1
1	1	1



Exercício 1

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

Indução de Árvores de Decisão ⁽³⁾

- Geralmente usa estratégia gulosa
- Divide progressivamente objetos baseado em um atributo de teste
- Escolha do atributo normalmente tenta otimizar algum critério

Indução de Árvores de Decisão ⁽²⁾

- Fase 1: **treinamento**
 - **Conj. treinamento**: permite induzir o modelo
 - **Conj. teste**: permite avaliar seu desempenho com dados não vistos durante indução
 - **Conj. validação**: semelhante a teste, mas o objetivo é podar a árvore
- Fase 2: **classificação**
 - Modelo induzido é utilizado em alguma aplicação do mundo real

Indução de Árvores de Decisão ⁽³⁾

- **Algoritmo de Hunt** (**foco**)
 - Um dos primeiros
 - Base de vários algoritmos atuais
- **J48**: no Weka
 - Hunt → ID3 → C4.5 (J48)

Algoritmo de Hunt

Se todas as instâncias \hat{x}_i do dataset X \hookleftarrow
são da mesma classe y_j então

t é um nó folha rotulado como y_j

Se $X = \emptyset$ então

t é um nó folha rotulado pela \hookleftarrow
classe majoritária, y_m

Se existem x_i de diferentes classes $\in X$ então

Dividir X em subconjuntos X_1 e X_2 com \hookleftarrow
um atributo de teste

Repetir procedimento a cada subconjunto

- Obs: X = conjunto de instâncias que atingem o nó t

Algoritmo de Hunt (2)

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Raiz:
classe default não (7/10)

Não

Algoritmo de Hunt ₍₃₎

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

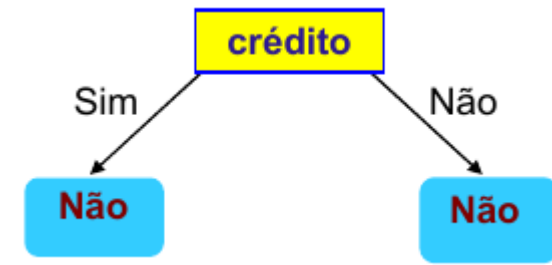
Nó crédito (nova raiz):
sim: classe não (3/3)

Não

Algoritmo de Hunt (4)

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

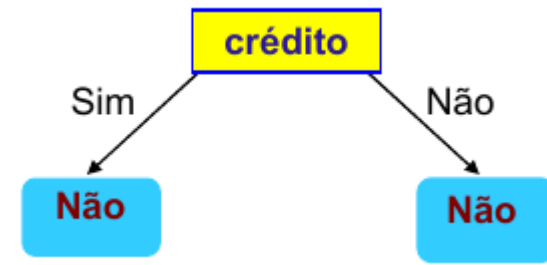
Nó crédito (nova raiz):
sim: classe não (3/3)
não: classe não (4/7)



Algoritmo de Hunt (5)

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

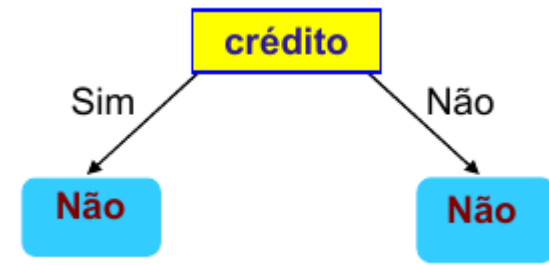
Acerto para crédito = sim já está ok.
Vamos seguir para crédito = não



Algoritmo de Hunt ⁽⁶⁾

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Nó estado
casado: classe não (3/3)



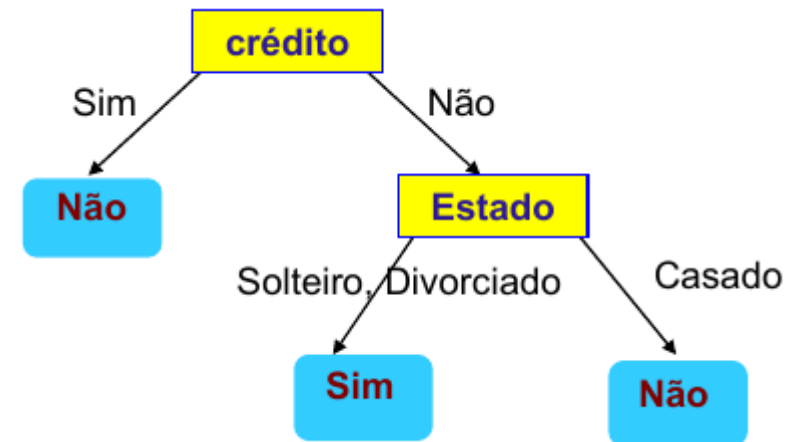
Algoritmo de Hunt ⁽⁷⁾

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Nó estado

casado: classe não (3/3)

solteiro/divorciado: classe sim (3/4)



Algoritmo de Hunt ₍₈₎

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

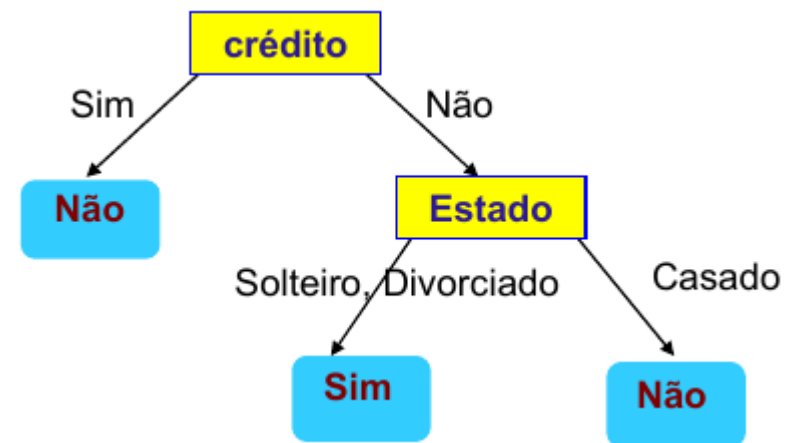
Análises ok:

crédito = sim

crédito = não && estado = casado

vamos seguir para

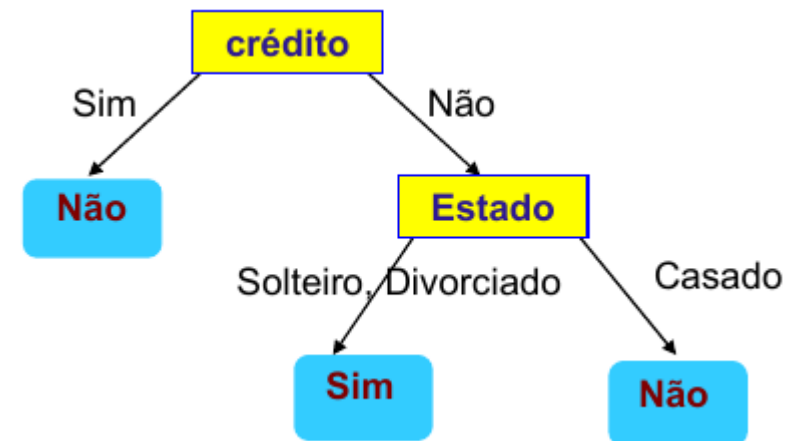
crédito = não && estado = sol/div.



Algoritmo de Hunt ₍₉₎

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Nó renda
<80k = classe não (1/1)



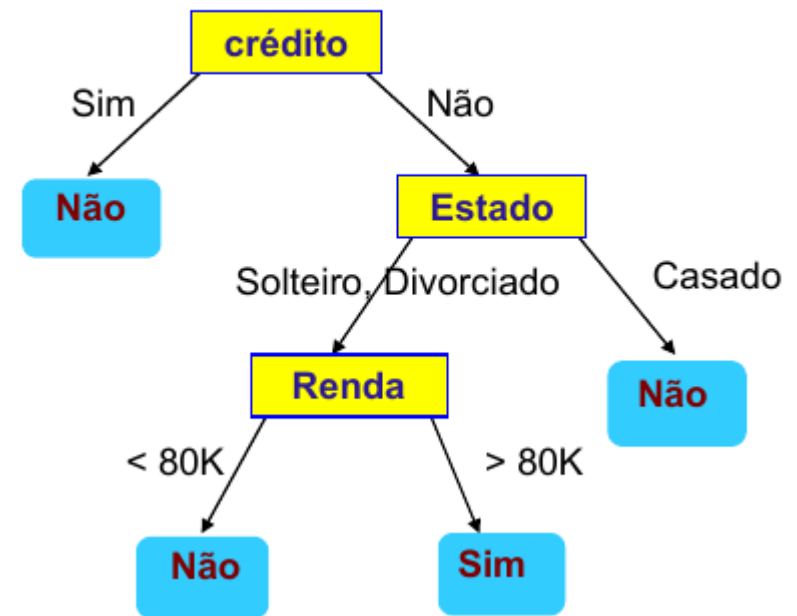
Algoritmo de Hunt ⁽¹⁰⁾

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim

Nó renda

<80k = classe não (1/1)

>80k = classe sim (3/3)



Algoritmo de Hunt ⁽¹¹⁾

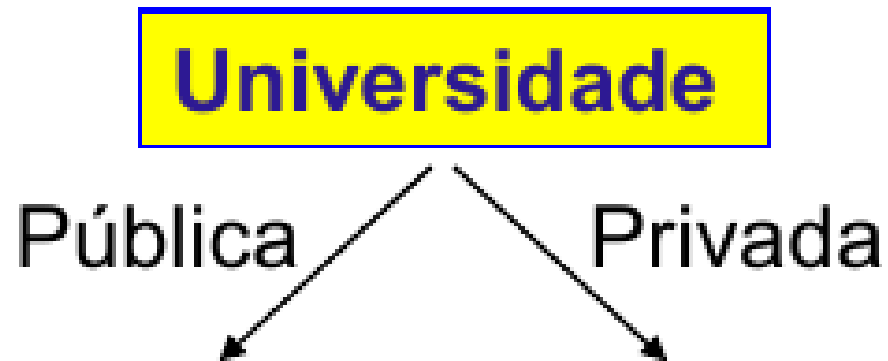
- Algoritmo de Hunt não especifica como:
 - Dividir atributos não binários
 - Escolher nó para dividir em cada nível
 - Escolher o número de divisões por nó: duas ou mais?
- Algoritmos derivados atacam essas questões

Dividindo atributos

- Divisão varia de acordo com as estratégias:
 - Binária: exatamente duas por nó (**foco**)
 - N-ária: um ramo para cada valor
 - Intermediária: alguns nós podem ficar agrupados

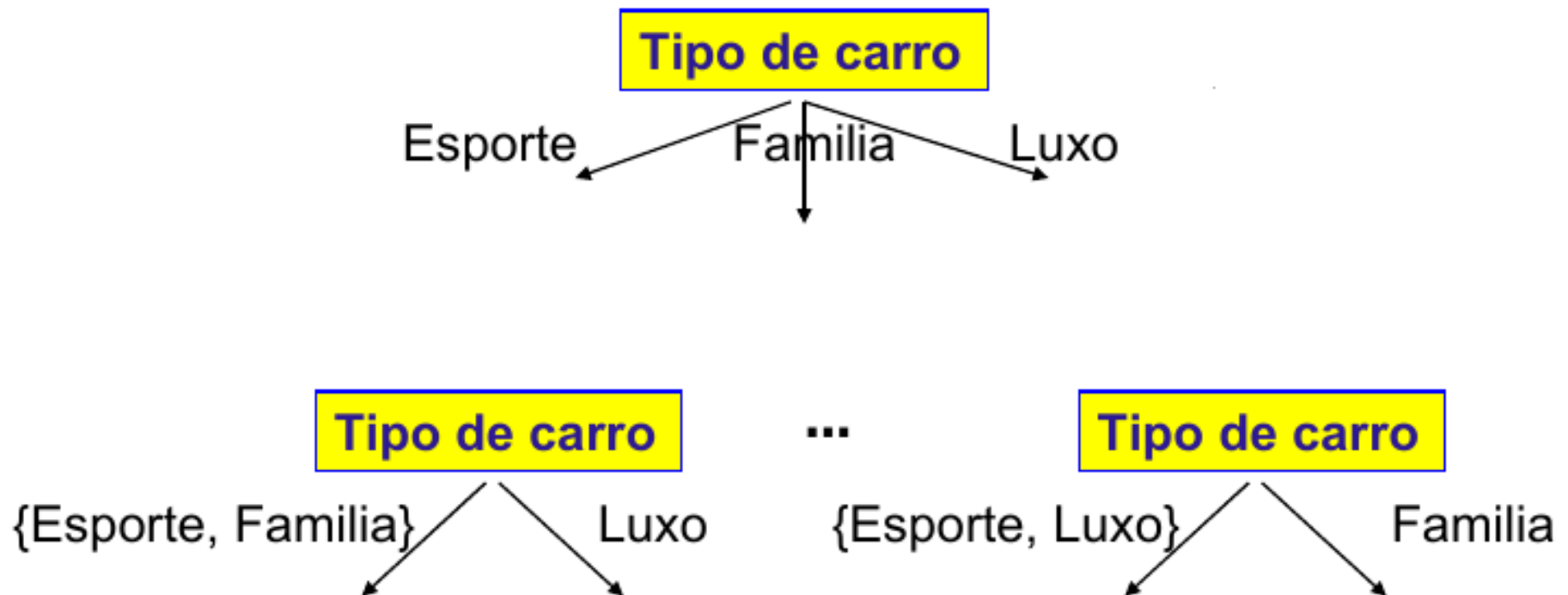
Dividindo atributos (2)

- **Atributo binário:** caso mais simples



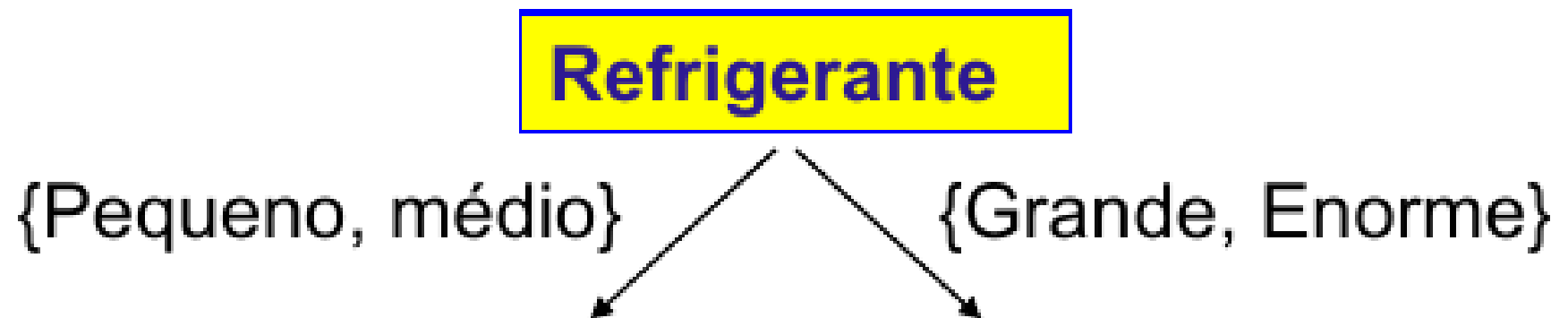
Dividindo atributos ⁽³⁾

- **Atributo nominal:** depende do número de divisões desejado: binária (foco) ou múltipla



Dividindo atributos (4)

- **Atributo ordinal**: semelhante binária, mas relação de ordem deve ser respeitada

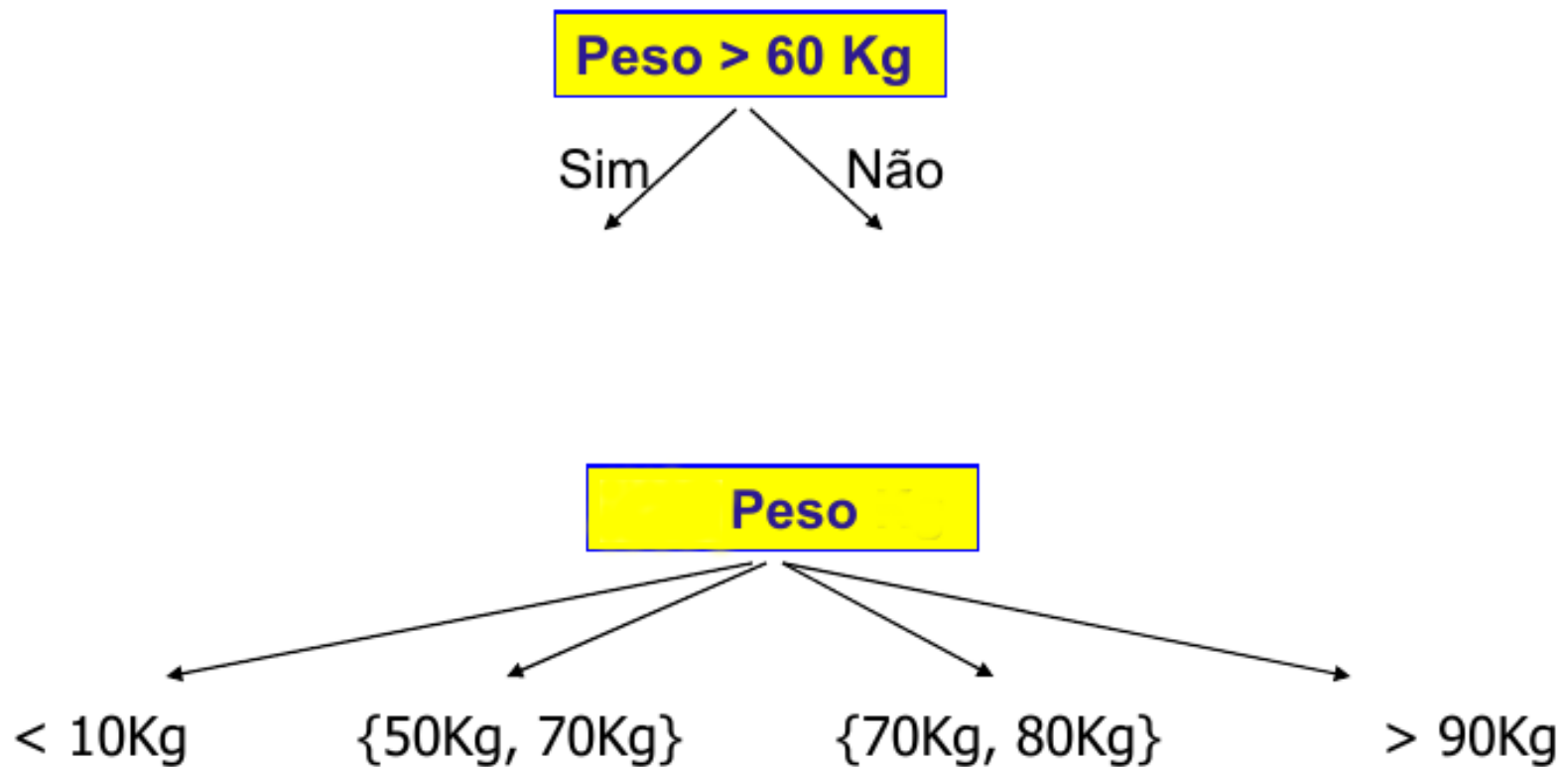


Dividindo atributos (5)

- **Atributo contínuo:**
 - Usa estratégia de discretização
 - Escolher pontos de corte que gera conjuntos mais puros
 - Mais sobre pureza vs cortes adiante...

Dividindo atributos (6)

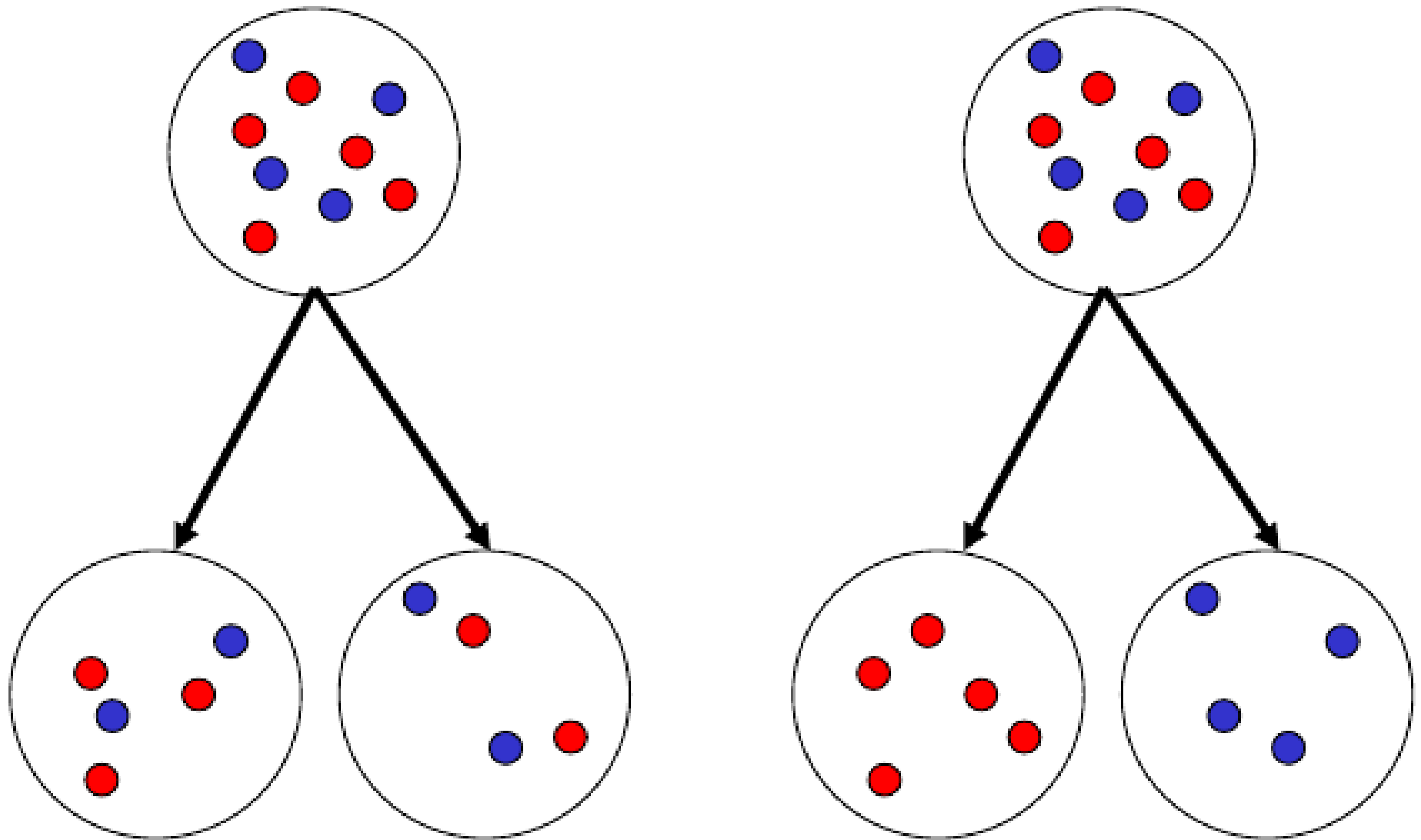
- Atributo contínuo:**



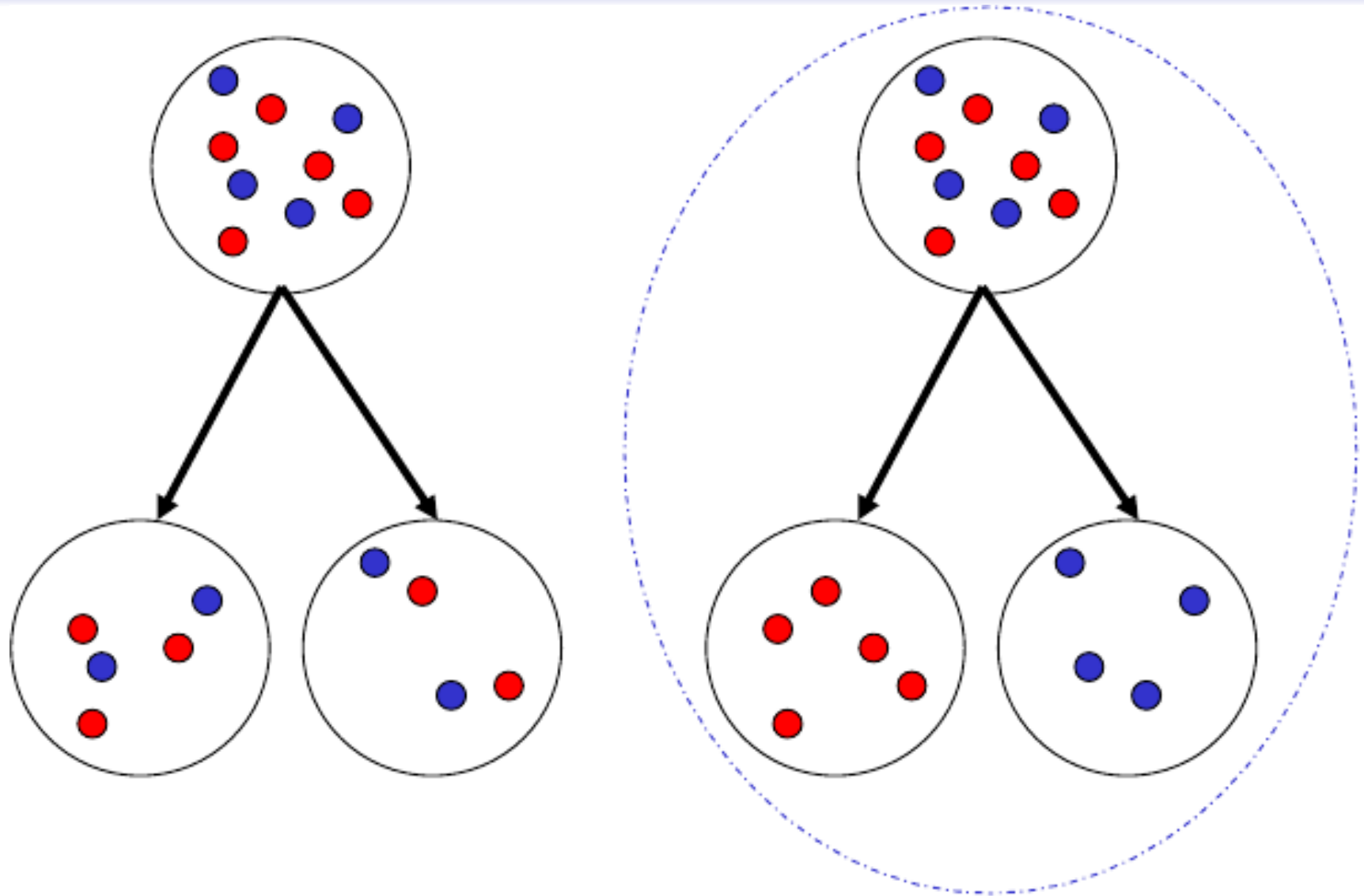
Escolhendo nós

- Idealmente: raiz deve trazer melhor divisão possível (estratégia gulosa)
- Precisamos de uma maneira de medir a qualidade:
 - (a) de um conjunto de dados
 - (b) da divisão de um conjunto de dados em conjuntos menores
- No slide a seguir, identifique qual a melhor divisão para um problema com a classe azul e a vermelha

Escolhendo nós ₍₂₎



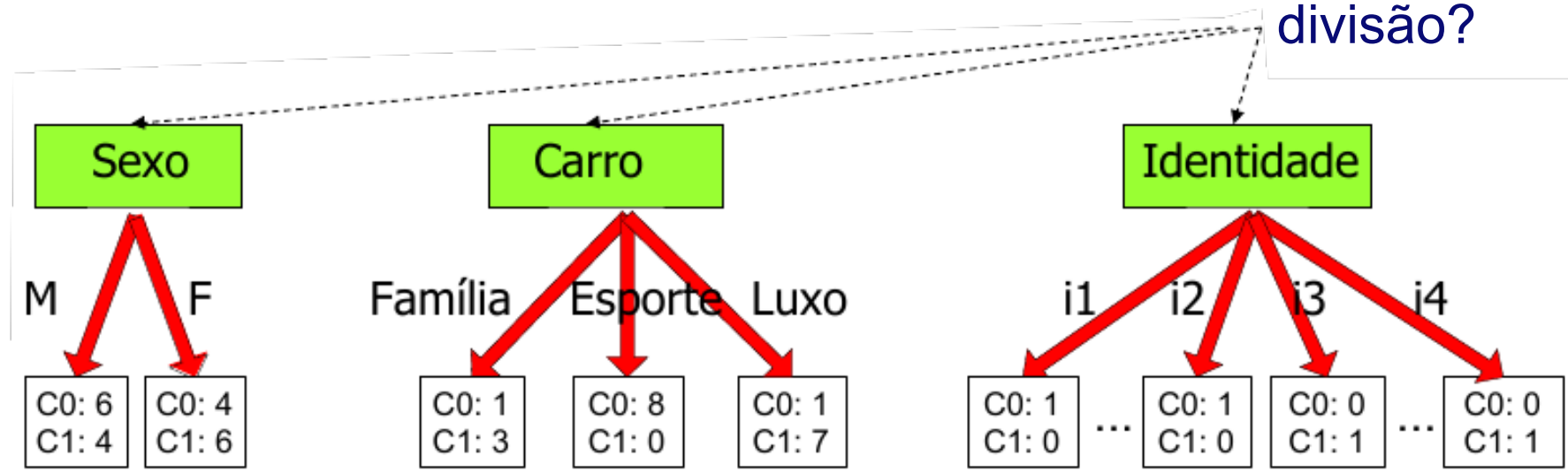
Escolhendo nós ₍₃₎



Escolhendo nós ₍₄₎

- Supor que D possui antes da divisão
 - 10 exemplos da classe 0 (C0: 10)
 - 10 exemplos da classe 1 (C1: 10)

Qual o melhor atributo para iniciar a divisão?



Medidas de impureza

- Desejamos que uma divisão gere subconjuntos mais puros o possível
- Para avaliar os subconjuntos, usamos **medidas de impureza**
 - Quanto menor seu valor, melhor
 - Desejamos minimizar a impureza

Medidas de impureza ⁽²⁾

- Divisão deve preferir nós com distribuição mais homogênea (pura) de classes

C0: 5
C1: 5

Não homogênea
Alto grau de impureza

C0: 9
C1: 1

Homogênea
Baixo grau de impureza

- Focaremos três medidas: entropia, gini, erro de classificação

Medidas de impureza ⁽³⁾

$$\text{Entropia (t)} = - \sum_{y_i \in Y} p(y_i) \log_2 p(y_i)$$

$$\text{Gini (t)} = 1 - \sum_{y_i \in Y} p(y_i)^2$$

$$\text{ErroClass (t)} = 1 - \max_{y_i \in Y} p(y_i)$$

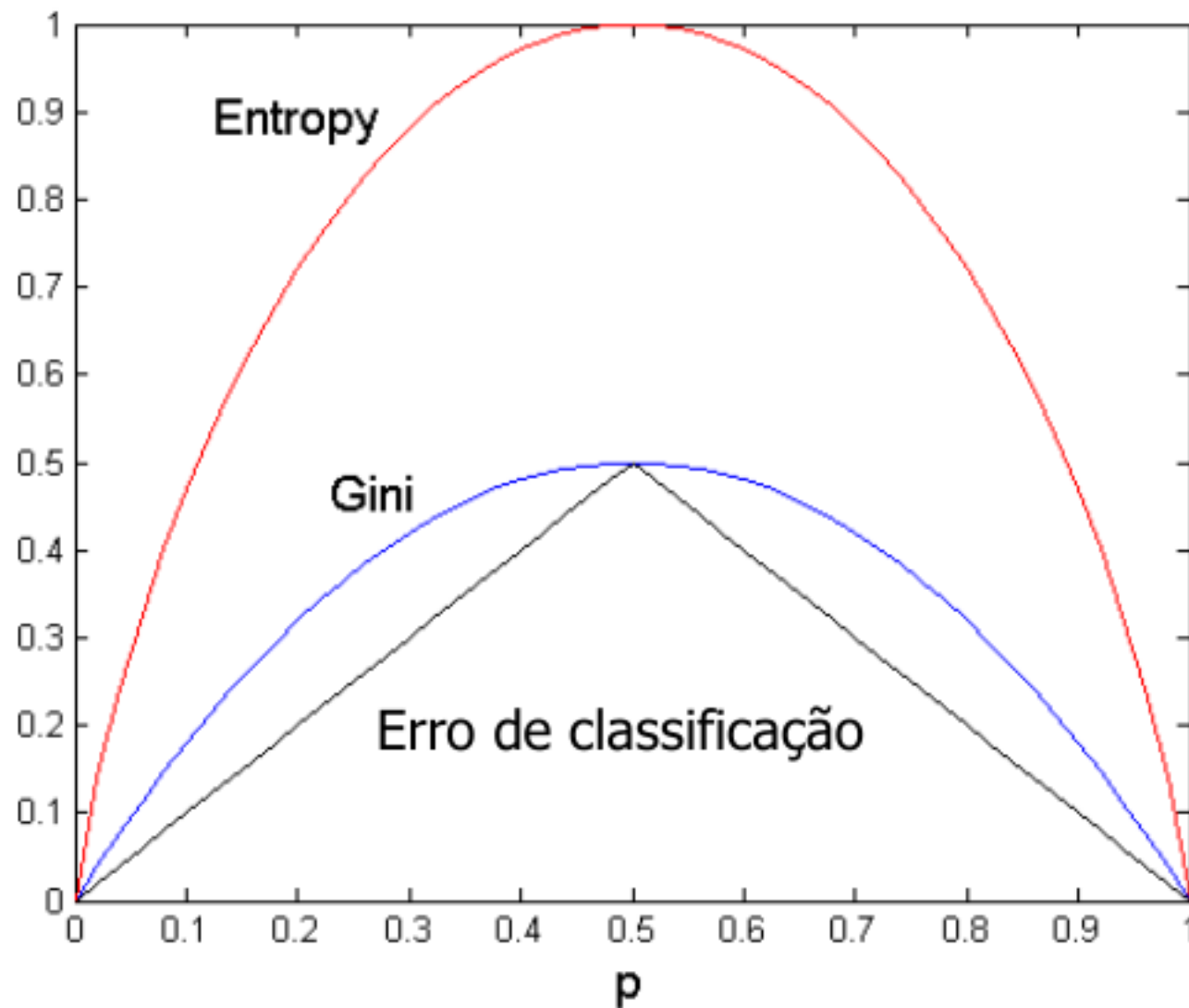
$$\text{percentual } p(y_i) = \frac{n(X, y_i)}{|X|}$$

Medidas de impureza ⁽⁴⁾

Onde:

- X : instâncias do subconjunto em questão
- Y : classes do subconjunto em questão
- y_i = a i -ésima classe
- $n(X, y_i)$ = número de instâncias da classe y_i no subconjunto em questão
- $0 \log_2 0 = 0$

Medidas de impureza ⁽⁵⁾



Medidas de impureza ⁽⁶⁾

- Máximos: quando os dados estão igualmente distribuídos entre todas as classes
 - Gini: $1 - 1/|Y|$
 - Erro de classificação: $1 - 1/|Y|$
 - Entropia: $\log_2 |Y|$
- Mínimos: quando todos os dados pertencem a uma classe (pureza máxima)
 - Todas: 0

Medida Gini

- Calcular a medida de impureza Gini para os dados abaixo

$$Gini(t) = 1 - \sum_{y_i \in Y} p(y_i)^2$$

C1	0
C2	6
Gini=?	

C1	1
C2	5
Gini=?	

C1	2
C2	4
Gini=?	

C1	3
C2	3
Gini=?	

Exemplo Medida Gini ₍₂₎

$$P(y_1) = 0/6 = 0 \quad P(y_2) = 6/6 = 1$$
$$\text{Gini} = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

$$P(y_1) = 1/6 \quad P(y_2) = 5/6$$
$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

$$P(y_1) = 2/6 \quad P(y_2) = 4/6$$
$$\text{Gini} = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

$$P(y_1) = 3/6 \quad P(y_2) = 3/6$$
$$\text{Gini} = 1 - (3/6)^2 - (3/6)^2 = 0.500$$

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	

Entropia

- Medida utilizada em Teoria da Informação proposta por Shannon
 - Permite determinar a capacidade de comunicação de um canal em termos de bits
 - Diferente da entropia definida na física (termodinâmica)
 - Relacionada a quantidade de energia em um sistema físico

Entropia ₍₂₎

- Mede quanta aleatoriedade existe em um sinal ou em um evento aleatório
 - Quanta informação é carregada por um sinal
- Shannon não sabia como chamar sua medida
 - Pediu opinião a Newton
 - “Você deveria chamá-la de Entropia”
 - “Ninguém sabe o que é entropia realmente é, assim, em um debate você estará em posição de vantagem”

Exercício

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

$$Entropia(t) = - \sum_{y_i \in Y} p(y_i) \log_2 p(y_i)$$

$$ErroClass(t) = 1 - \max_{y_i \in Y} p(y_i)$$

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=?	

Medida de Ganho

- As medidas ErroClass, Gini e Entropia dizem se um conjunto é puro ou não
- Uma divisão é boa quando a pureza dos filhos é maior que a pureza do pai
- A medida Ganho leva esse fator em conta
- Basicamente, escolhe-se o atributo que gere a divisão mais pura

Medida de Ganho ₍₂₎

$$\Delta = I(X_{pai}) - \sum_{i=1}^n \frac{|X_i|}{|X_{pai}|} I(X_i)$$

- X_{pai} : conjunto de instâncias do nó pai
- X_i : conjunto de instâncias em cada filho da divisão
- $|X|$: número de instâncias dentro do conjunto X
- $I(X)$: uma medida de impureza aplicada a um dado conjunto

Medida de Ganho ⁽³⁾

- A somatória calcula a impureza ponderada dos filhos
 - Conta a impureza de cada um, mas privilegia filhos com mais instâncias
- Δinfo : **ganho de informação**
 - Quando a medida de impureza é entropia

Exercício 2: calcular Gini's e Ganho



Medida GiniP

- Quando um nó pai é dividido em k filhos, a qualidade da divisão é definida por:

$$GiniP = \sum_{i=1}^n \frac{|X_i|}{|X_{pai}|} Gini(X_i)$$

- Muito similar à medida ganho
 - Mas note que, devido a troca de sinal, está medida deve ser minimizada

Medida GiniP₍₂₎

- GiniP é a abreviação de Gini da Média Ponderada
 - Usa a medida Gini
- Assim como o Ganho, também avalia a qualidade da divisão
 - Ganho deve ser maximizada; GiniP minimizada
 - Não leva nó pai em comparação: permitindo comparar divisões em pontos diferentes da hierarquia
- Usada pelos algoritmos CART, SLIQ, SPRINT

Exercício 3

- De acordo com GiniP, qual é a melhor divisão possível?

	Tipo de Carro		
	Família	Esporte	Luxo
C1	1	2	1
C2	4	1	1
GiniP	???		

	Tipo de Carro	
	{Esporte, Luxo}	{Família}
C1	3	1
C2	2	4
GiniP	???	

	Tipo de Carro	
	{Esporte}	{Family, Luxury}
C1	2	2
C2	1	5
GiniP	???	

Exercício 3 ₍₂₎

	Tipo de Carro		
	Família	Esporte	Luxo
C1	1	2	1
C2	4	1	1
GiniP	0.393		

	Tipo de Carro	
	{Esporte, Luxo}	{Família}
C1	3	1
C2	2	4
GiniP	0.400	

	Tipo de Carro	
	{Esporte}	{Família, Luxo}
C1	2	2
C2	1	5
GiniP	0.419	

Taxa de Ganho

- Ganho e GiniP acabam privilegiando divisões com o maior número de subgrupos
 - Como pode ser visto no slide anterior
- Conjuntos menores tendem a ser mais puros
 - Pense na pureza de um conjunto unitário
- Taxa de Ganho é uma melhoria na medida Ganho que penaliza divisões com muitos subgrupos
- A seguir, exemplo para entropia

Taxa de Ganho ₍₂₎

$$TG = \frac{\Delta_{info}}{IV}$$

$$IV = - \sum_{v_i \in X} p(v_i) \log_2(v_i)$$

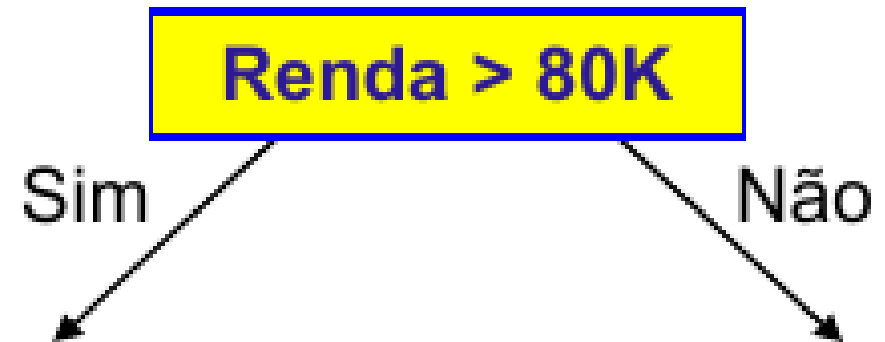
- Onde:
 - TG: taxa de ganho da divisão
 - Δ_{info} : ganho de informação da divisão
 - Continua no próximo slide...

Taxa de Ganho ₍₃₎

- ... continuação:
 - IV: valor intrínseco da divisão
 - v_i : valores do atributo x utilizados na divisão
 - $p(v_i)$ = percentagem das instâncias no nó pai que assumem o valor (v_i) para o atributo x
- Note que taxa de ganho lembra a entropia, mas a primeira é sobre valores de atributos e a segunda sobre as possíveis classes

Atributos contínuos

Crédito	Estado Civil	Renda	Deve
Sim	Solteiro	125K	Não
Não	Casado	100K	Não
Não	Solteiro	70K	Não
Sim	Casado	120K	Não
Não	Divorciado	95K	Sim
Não	Casado	60K	Não
Sim	Divorciado	220K	Não
Não	Solteiro	85K	Sim
Não	Casado	75K	Não
Não	Solteiro	90K	Sim



Atributos contínuos (2)

- Seleção de atributos para divisão
 - Binários: usar ganho ou giniP
 - Nominais/ordinais:
 - Divisão binária: ganho ou giniP
 - Divisão múltipla: taxa de ganho
 - **Contínuos**: mesma estratégia do caso anterior
 - Testando vários pontos de corte

Atributos contínuos ⁽³⁾

- Estratégias de escolha de pontos de corte
 - Força bruta: método mais simples consistindo em testar vários pontos de corte e calcular qualidade da divisão obtida
 - Ordenar valores do dataset em ordem crescente: avaliar qualidade do corte entre cada par de valores
 - Ordenação otimizada: idem anterior; testar corte só testar pares com classes diferentes

Atributos contínuos (4)

Valores
ordenados
→
Posições
de divisão
→

Deve	Não		Não		Não		Sim		Sim		Sim		Não		Não		Não		Não			
→ →	Renda																					
	60		70		75		85		90		95		100		120		125		220			
	55		65		72		80		87		92		97		110		122		172		230	
	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>
Sim	0	3	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0	3	0
Não	0	7	1	6	2	5	3	4	3	4	3	4	3	4	4	3	5	2	6	1	7	0
GiniP	0.420		0.400		0.375		0.343		0.417		0.400		<u>0.300</u>		0.343		0.375		0.400		0.420	

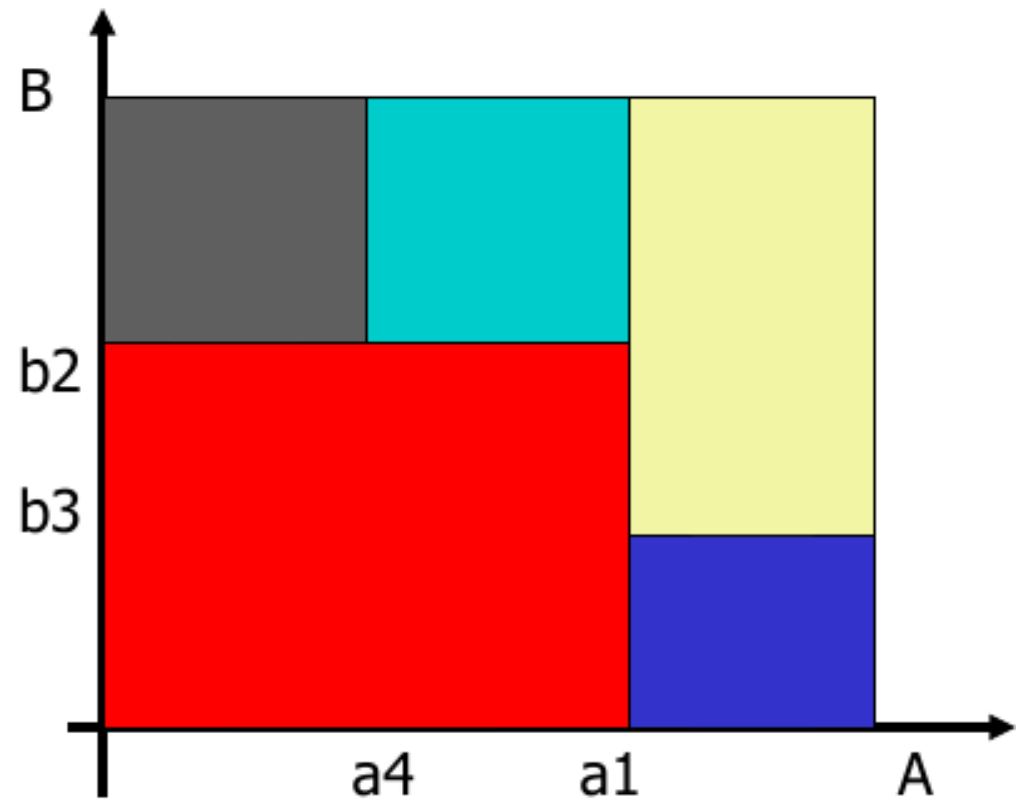
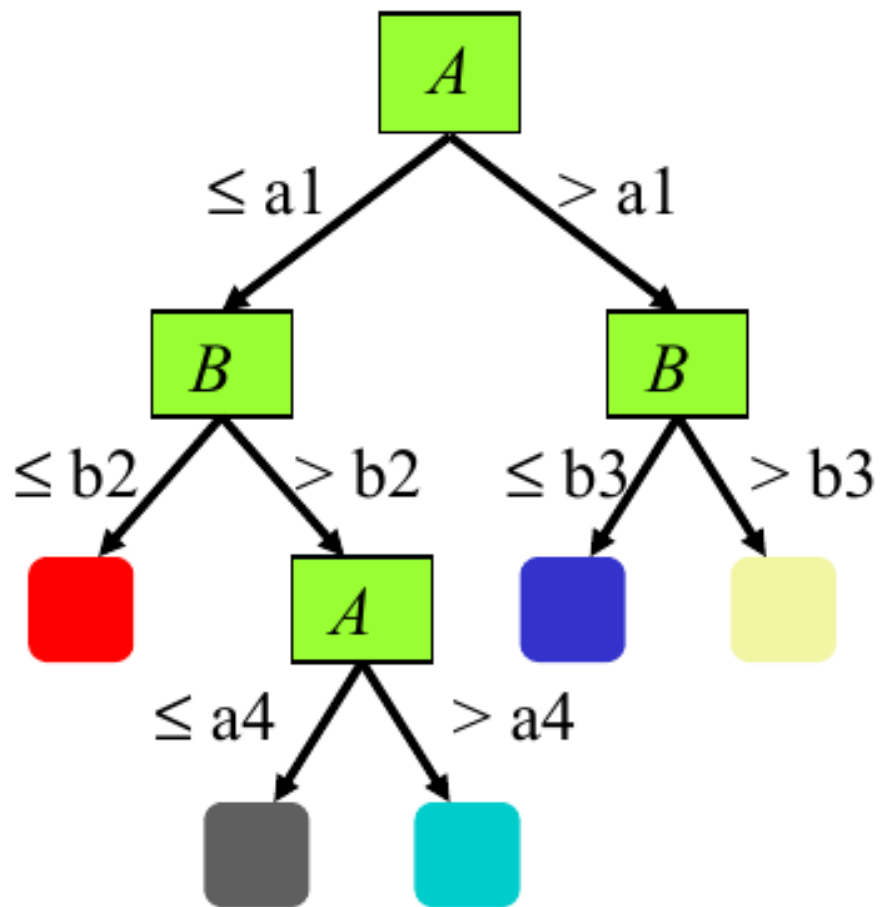
Treinamento

- Indução da árvore pode ser finalizada quando:
 - Quando os dados do nó atual têm o mesmo rótulo
 - Quando as instâncias no nó atual são duplicadas umas das outras (variando apenas a classe)
 - Quando o número de dados é menor que um dado valor
 - Todos os atributos já foram incluídos no caminho

Espaço de instâncias

- 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 instâncias
- Hiperretângulos:
 - Interseção entre eles é um conjunto vazio
 - União é o espaço total

Espaço de instâncias ⁽²⁾



Espaço de hipóteses ⁽³⁾

- Tem como saída uma única hipótese
 - Não há backtracking
 - Mínimo local
- Espaço de hipóteses completo (sujeito a overfitting)
 - Hipercubos podem aproximar hiperplanos não paralelos aos eixos

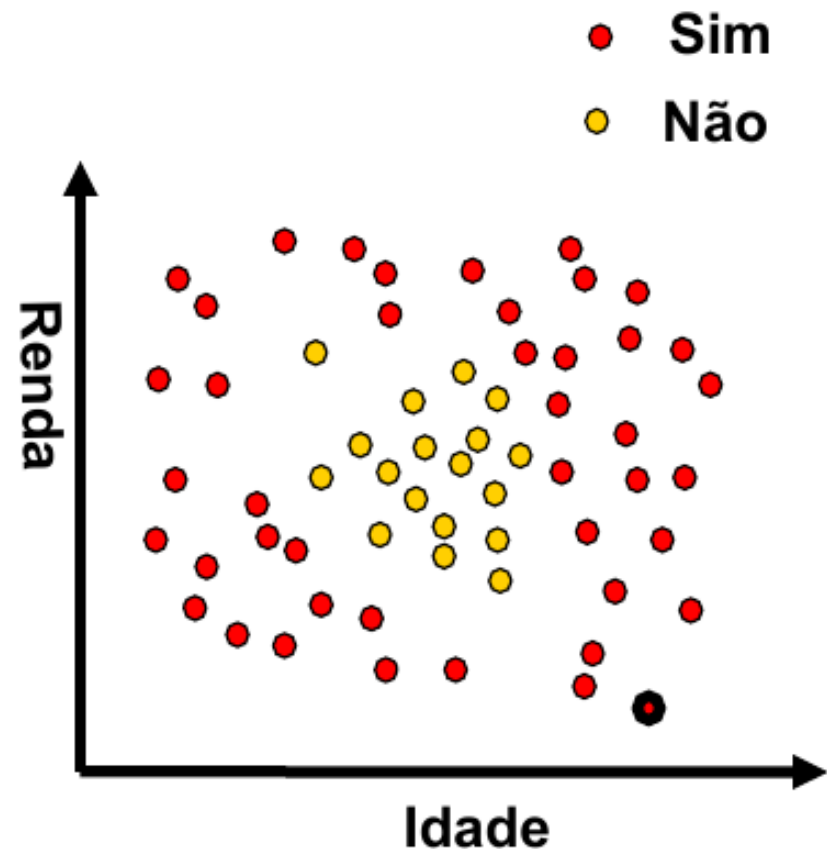
Exemplo 2

- Sejam os dados abaixo referentes a solicitações de crédito bancário
 - Construir uma árvore de decisão que classifica quem solicita crédito

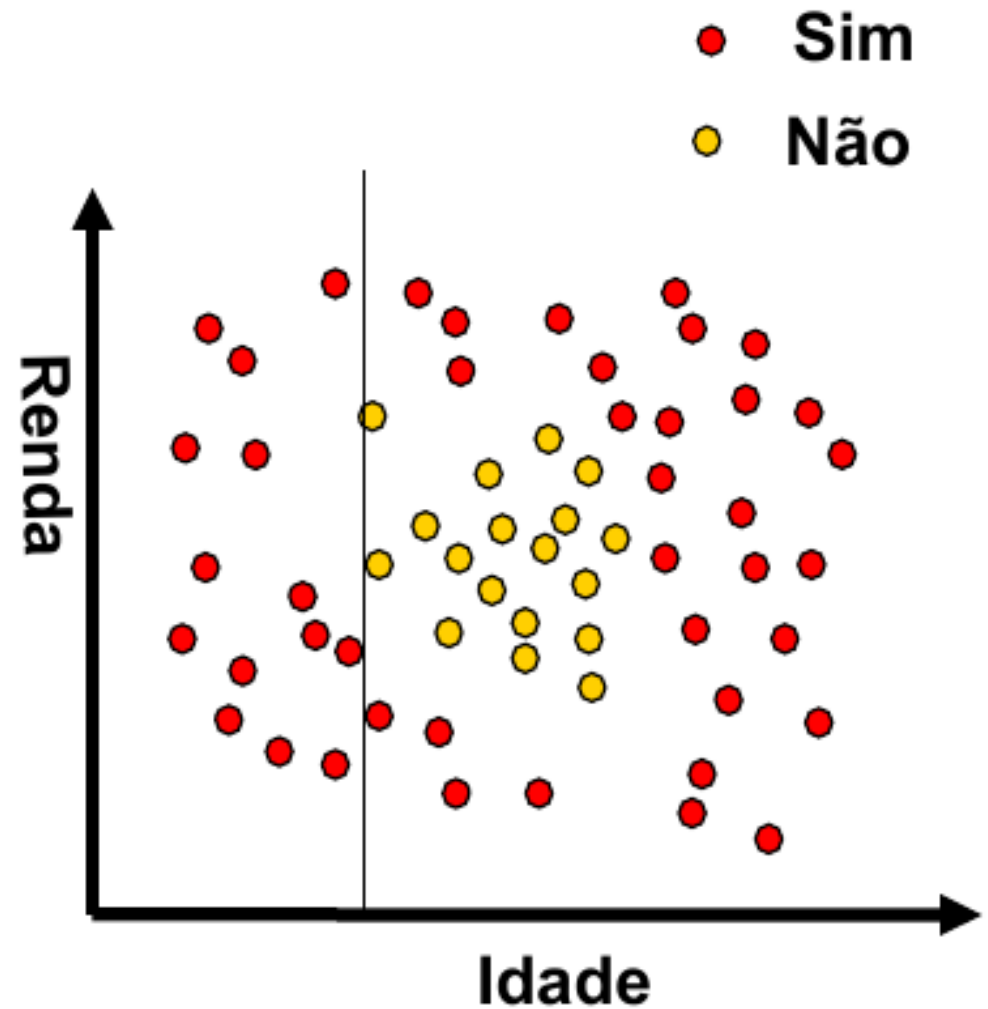
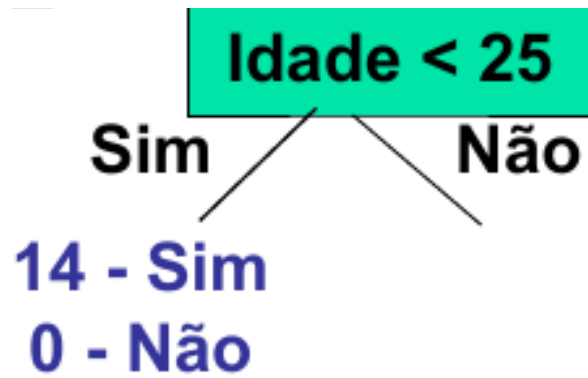
<i>Idade</i>	<i>Renda</i>	<i>Classe</i>
20	2000	Sim
30	5100	Não
60	5000	Sim
40	6000	Não
...

Indução

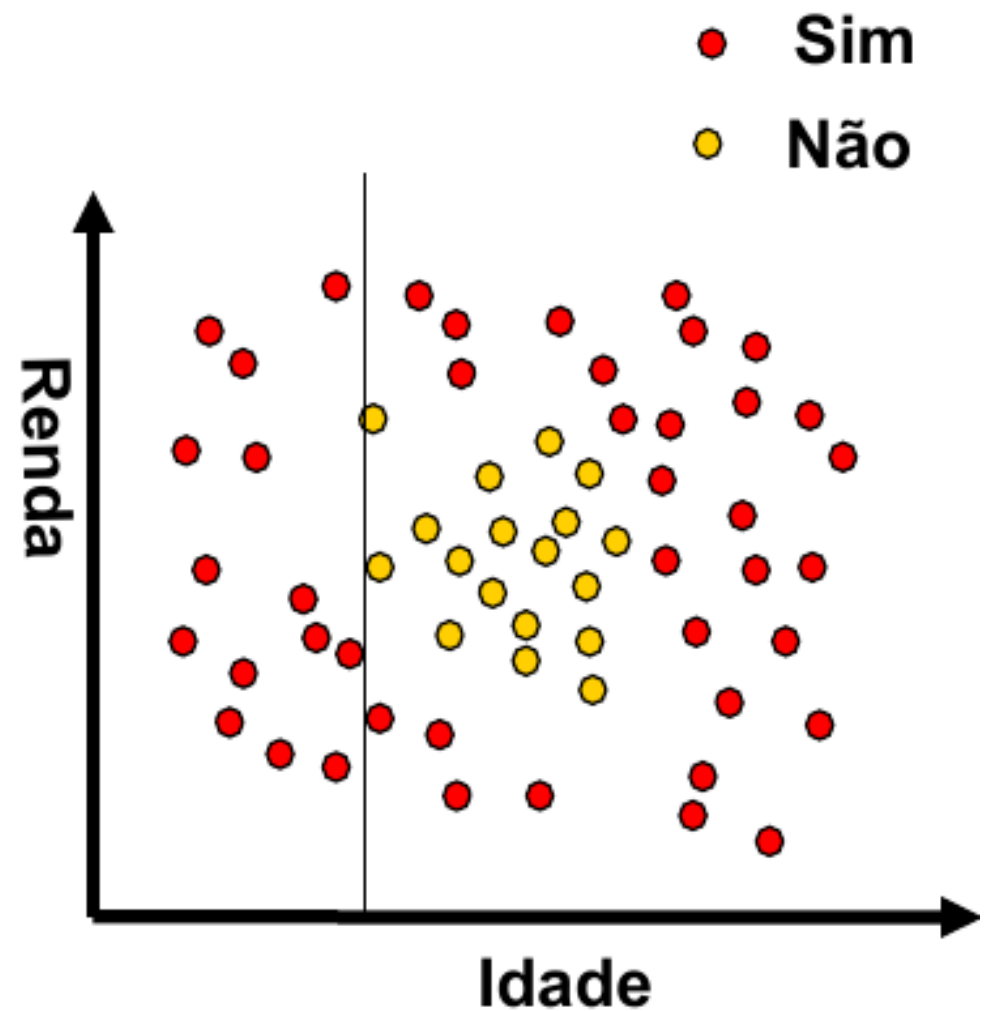
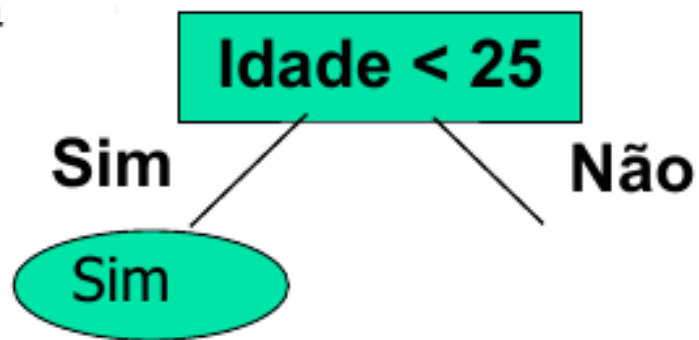
- Construir uma árvore de decisão que classifica quem solicitar crédito
 - Deve (Sim)
 - Não deve (Não)



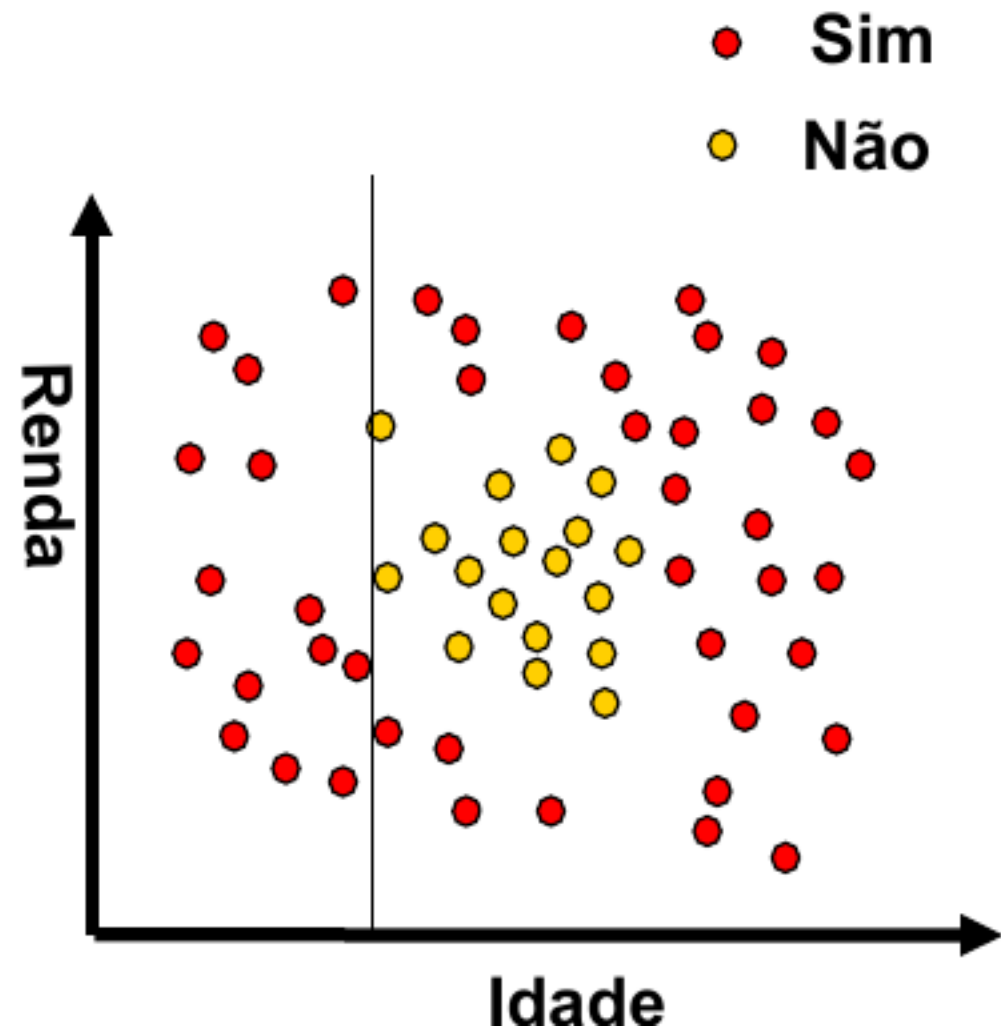
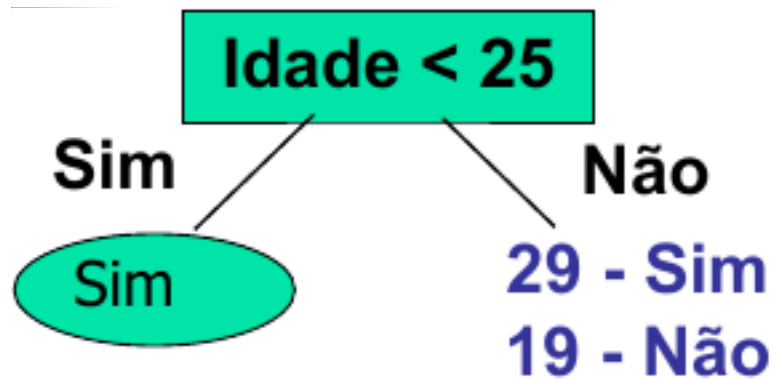
Indução ₍₂₎



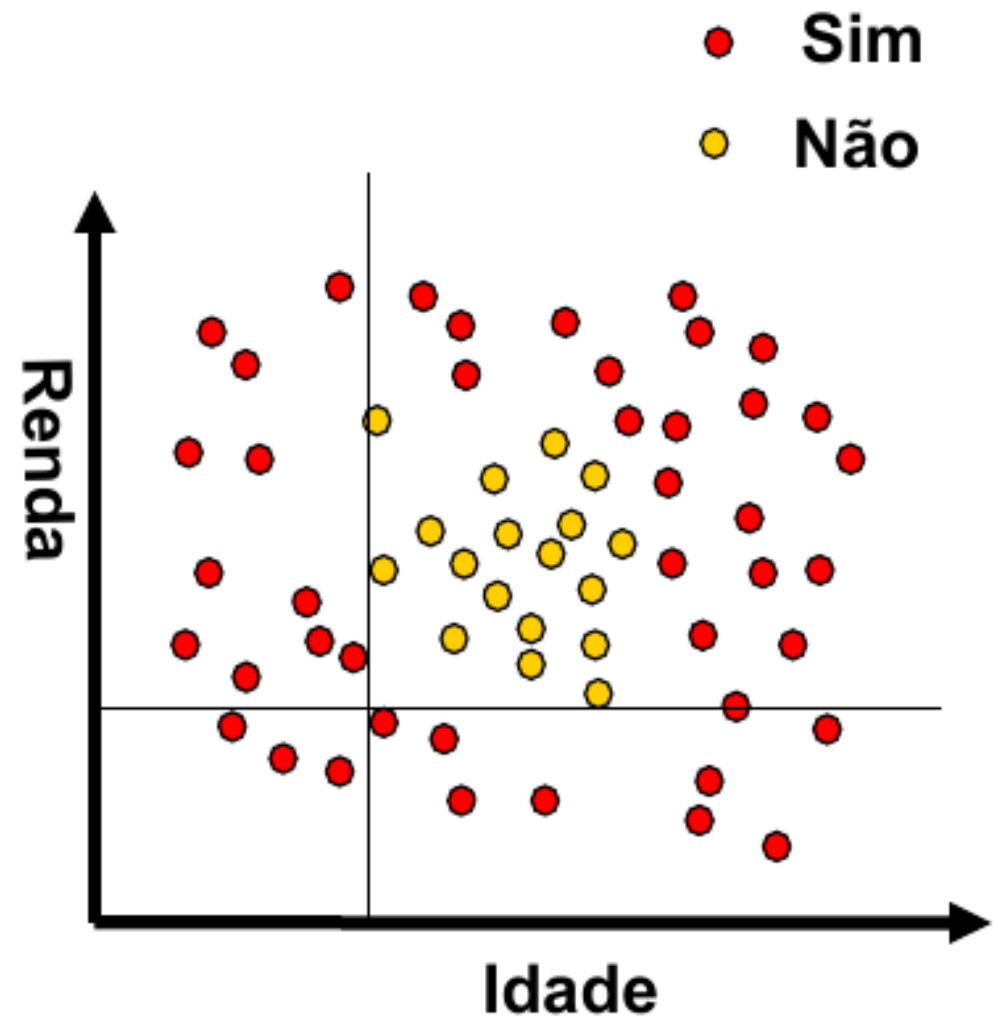
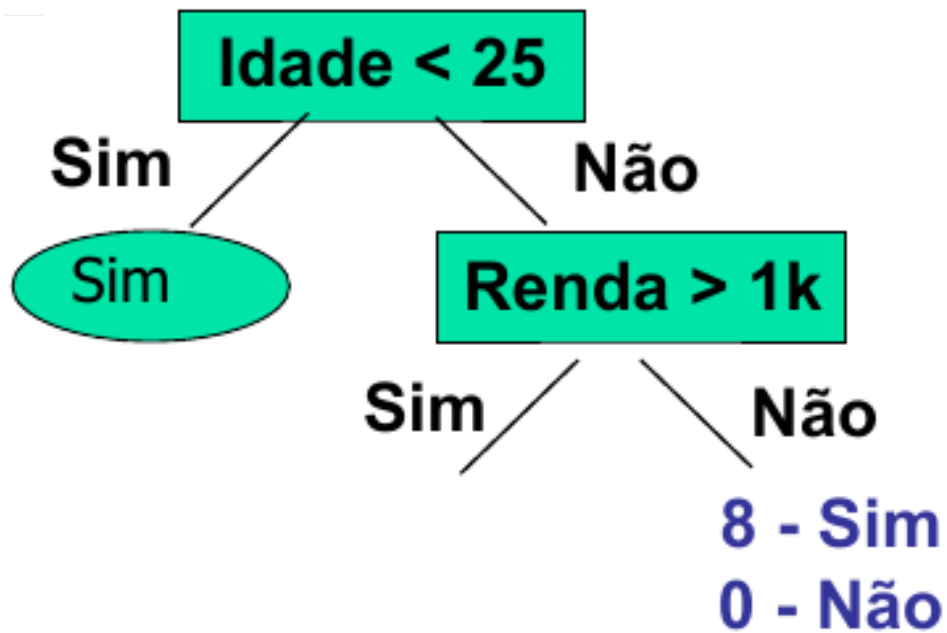
Indução₍₃₎



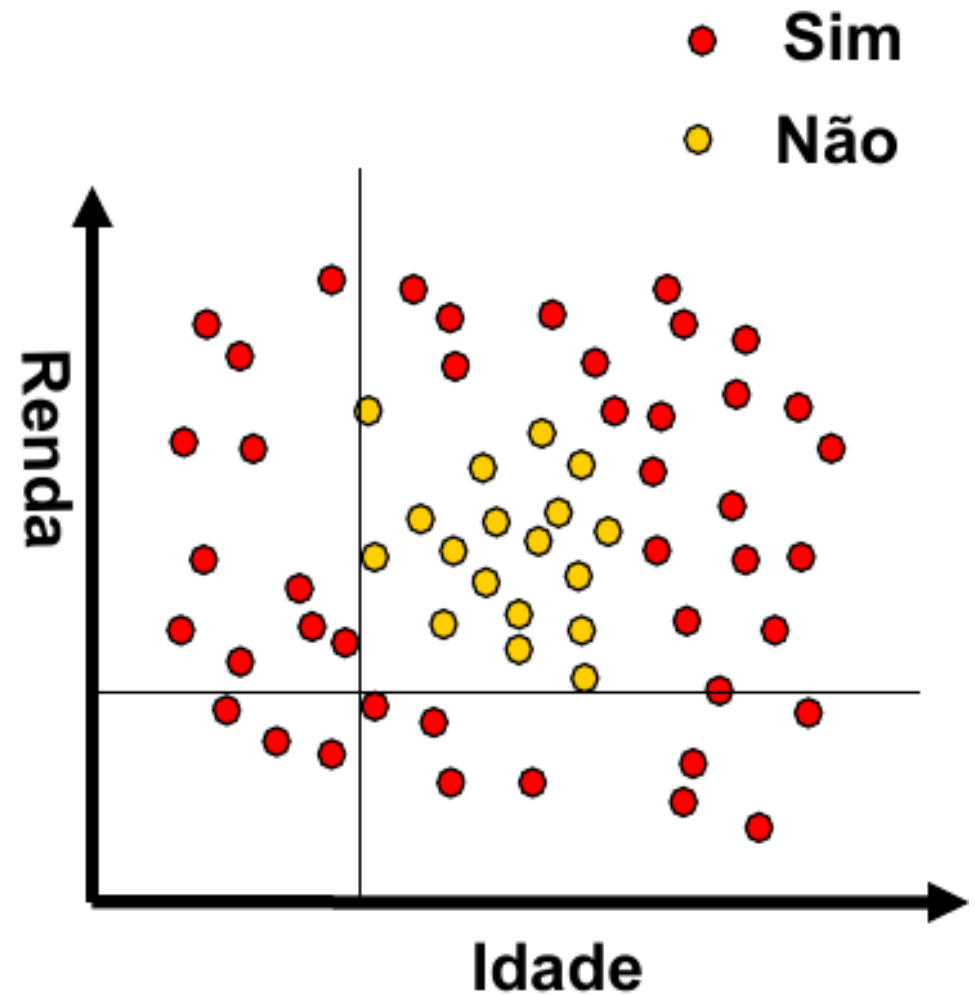
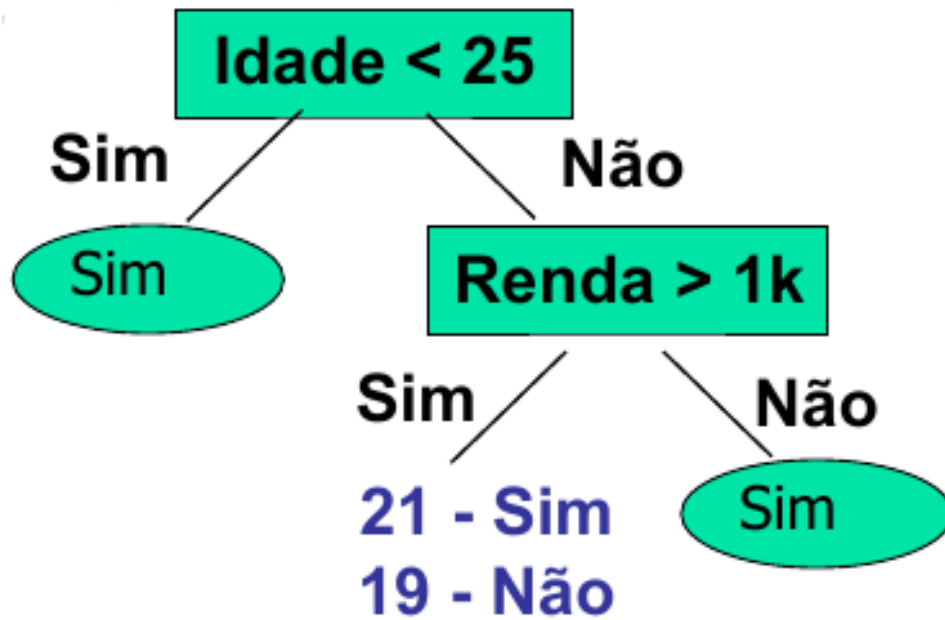
Indução ₍₄₎



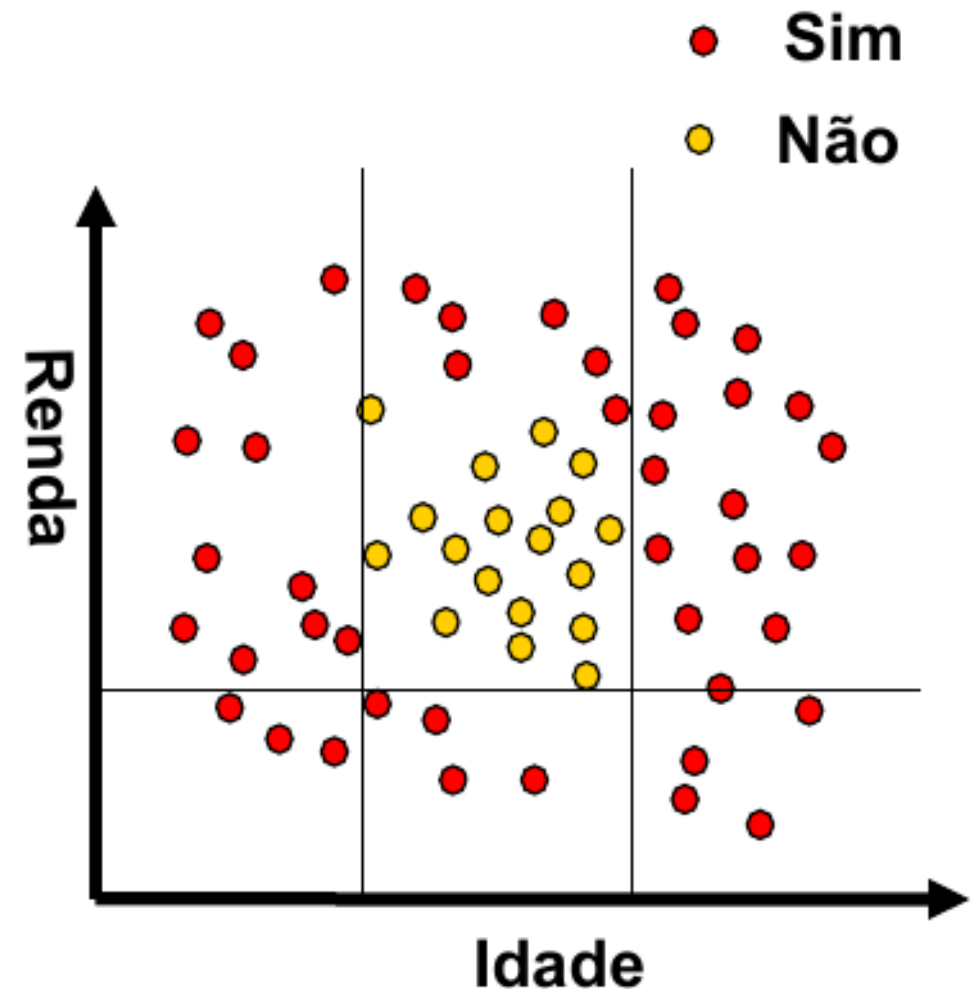
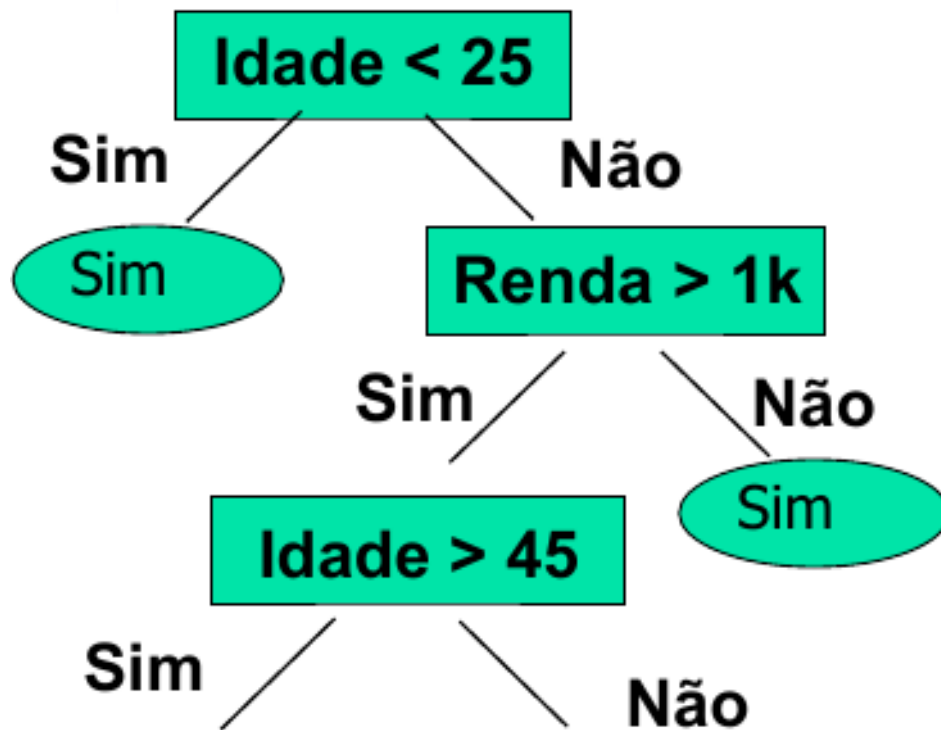
Indução ₍₅₎



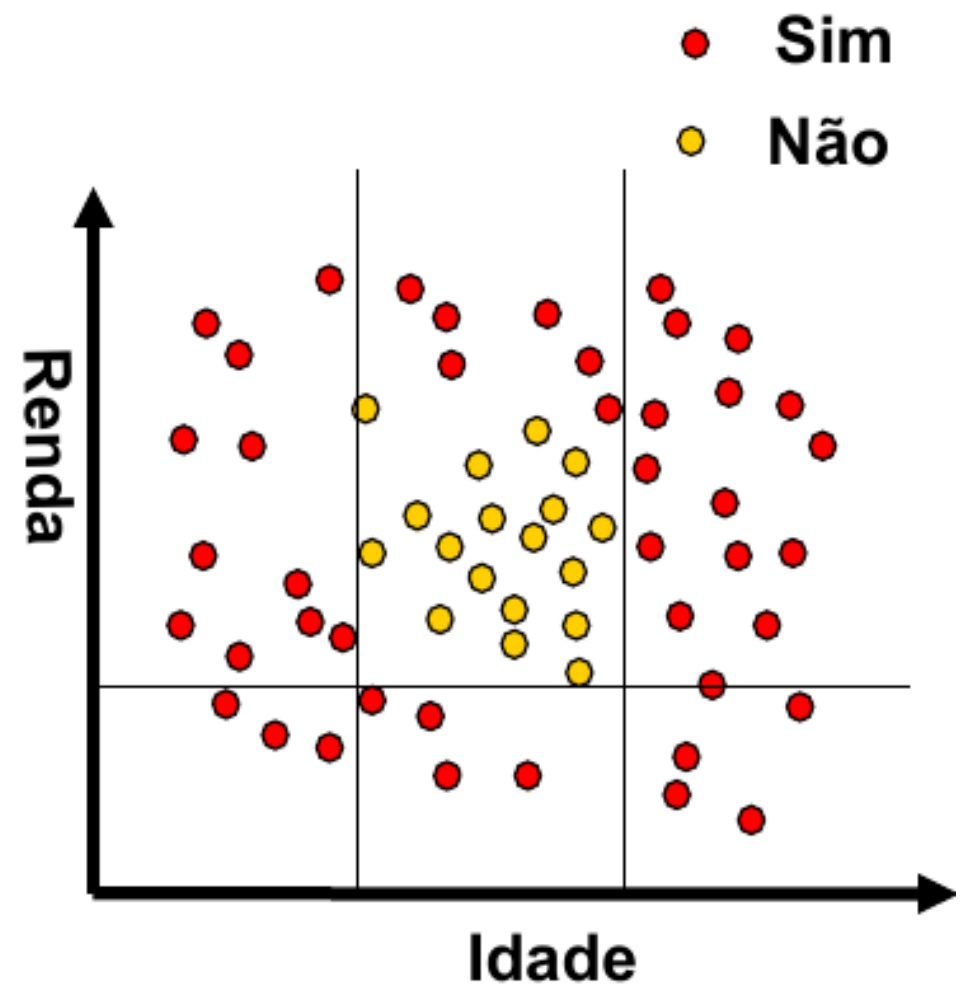
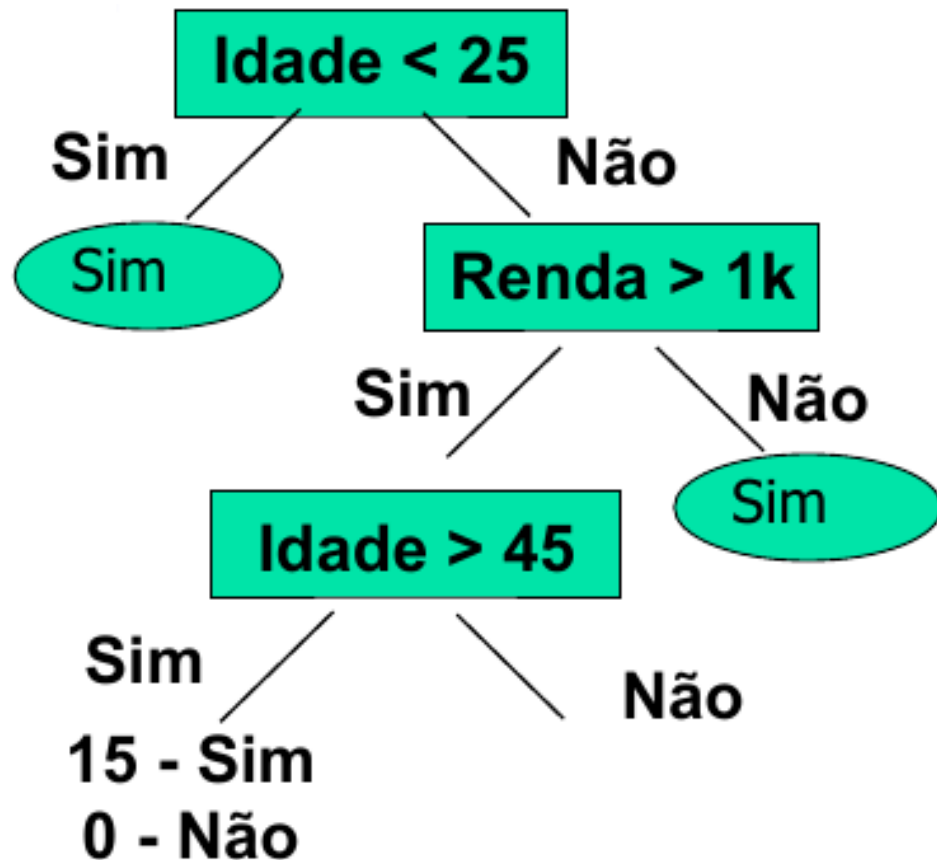
Indução ₍₆₎



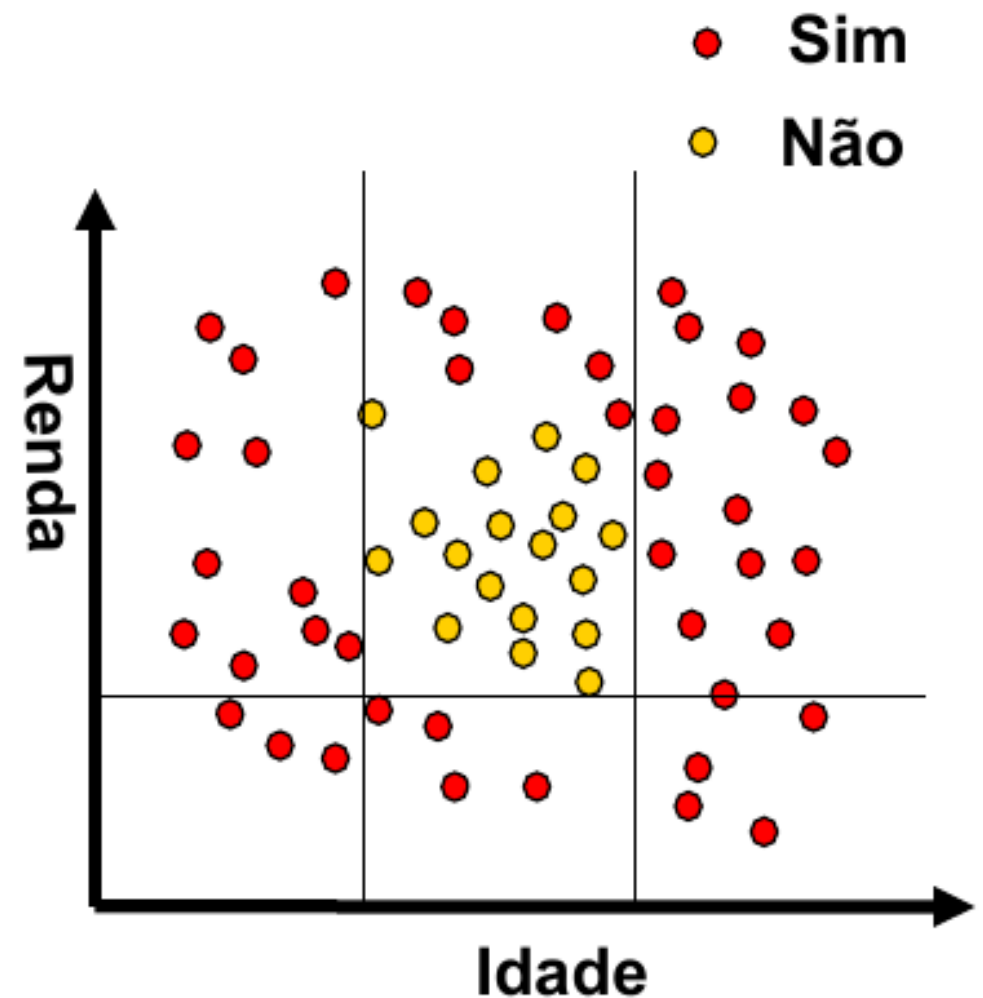
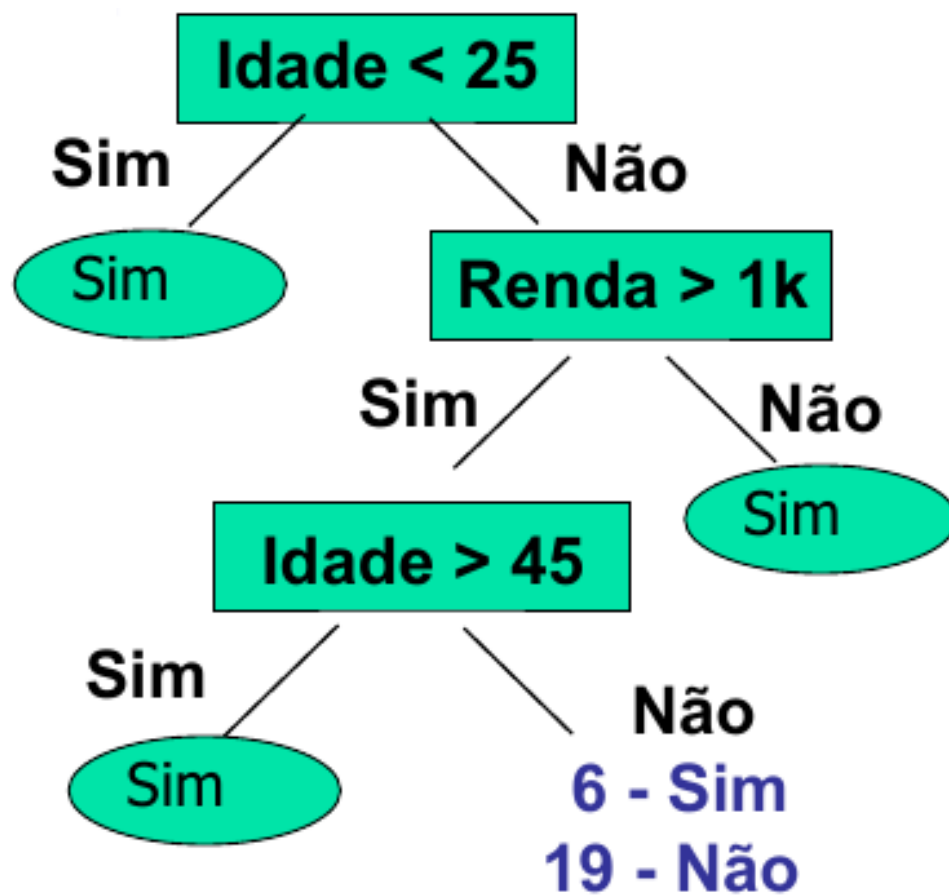
Indução ₍₇₎



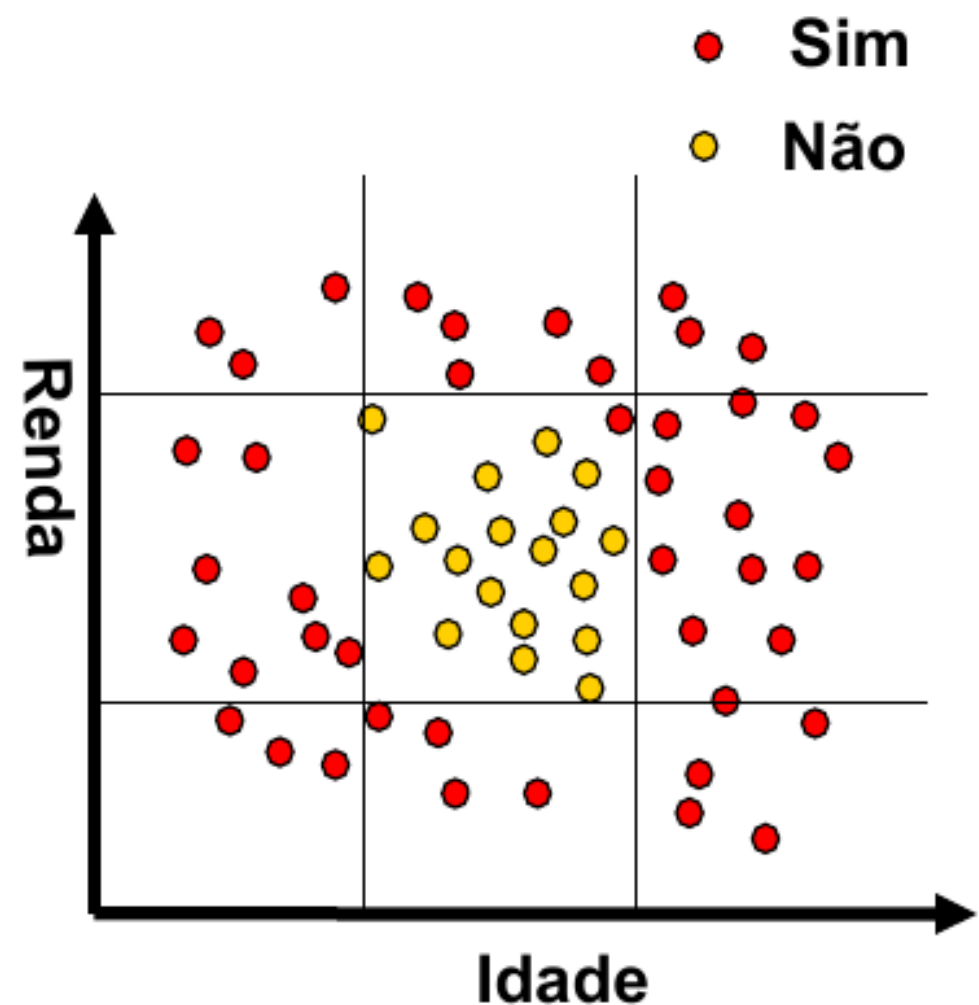
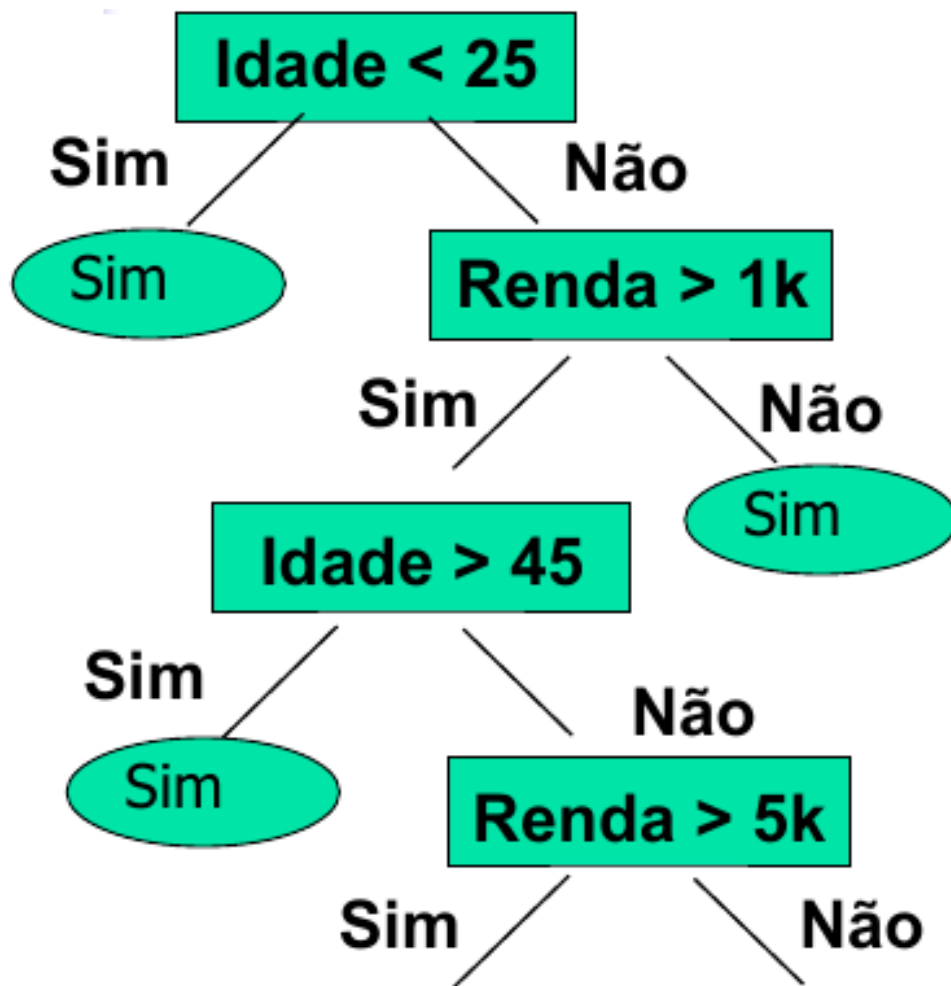
Indução ₍₈₎



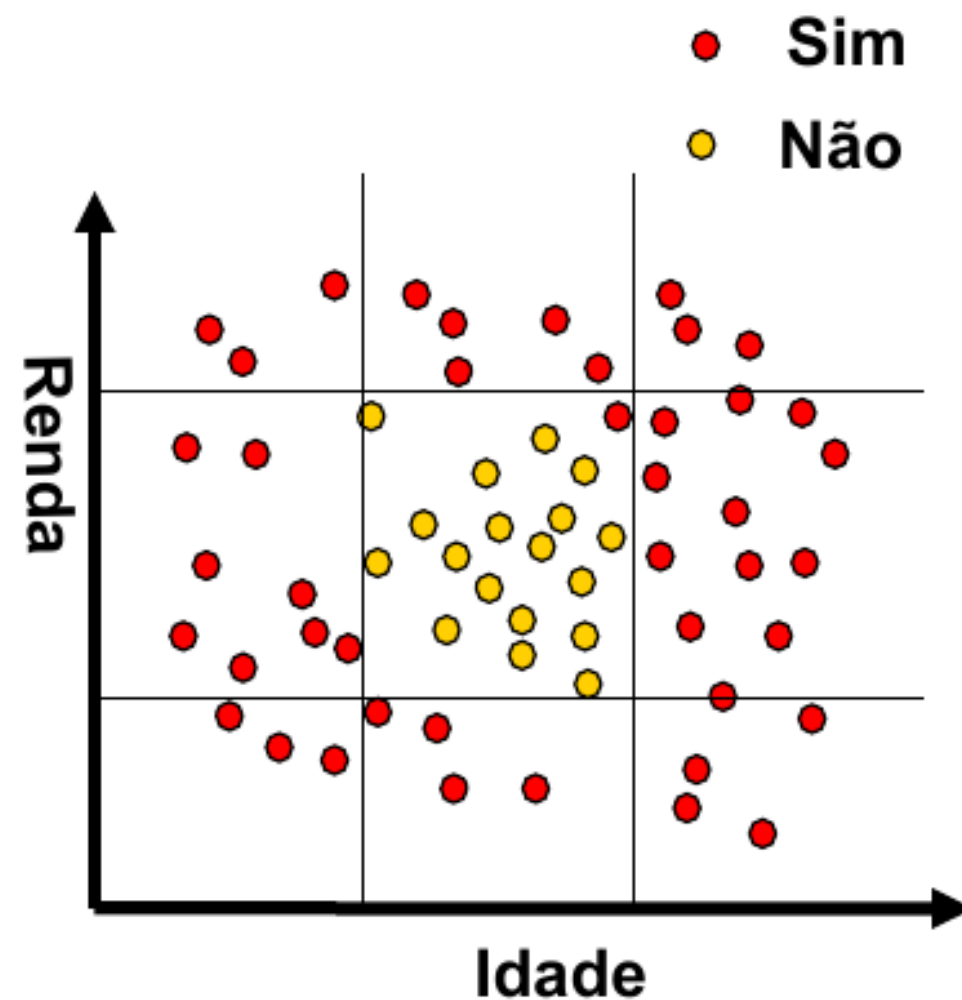
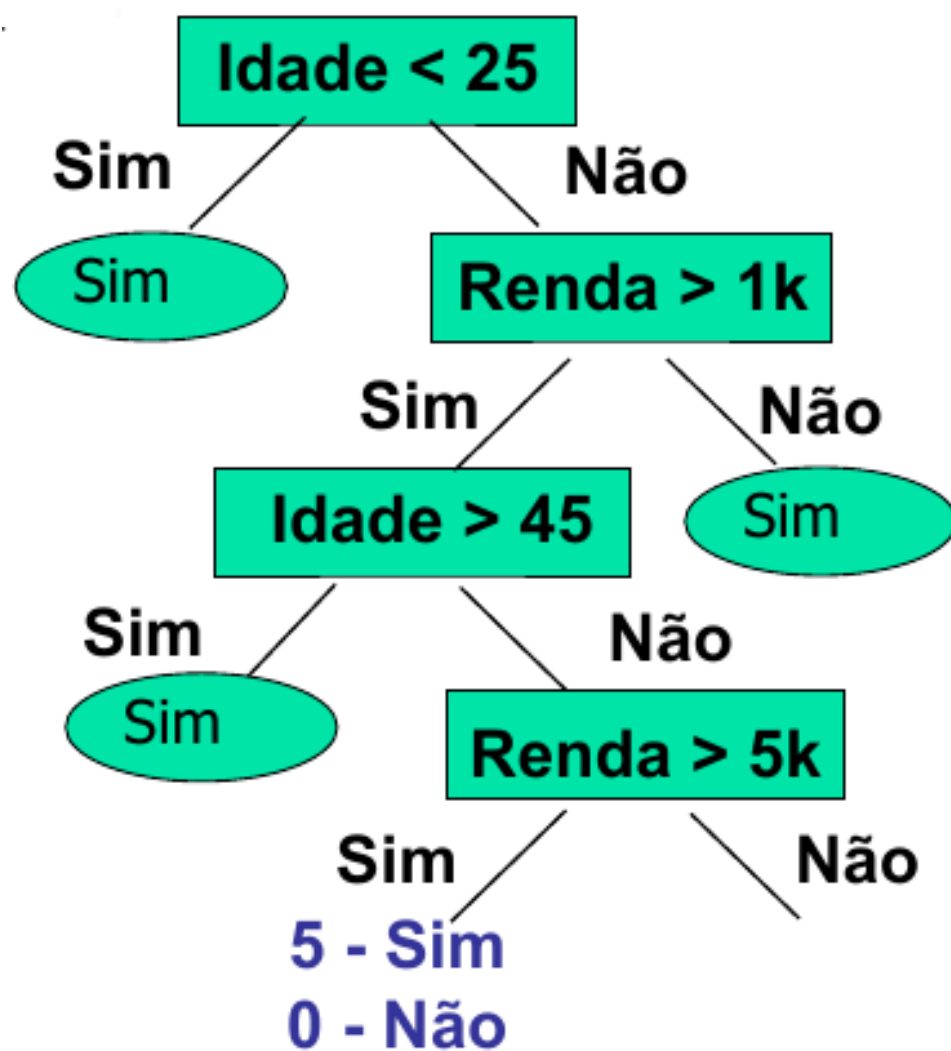
Indução ₍₉₎



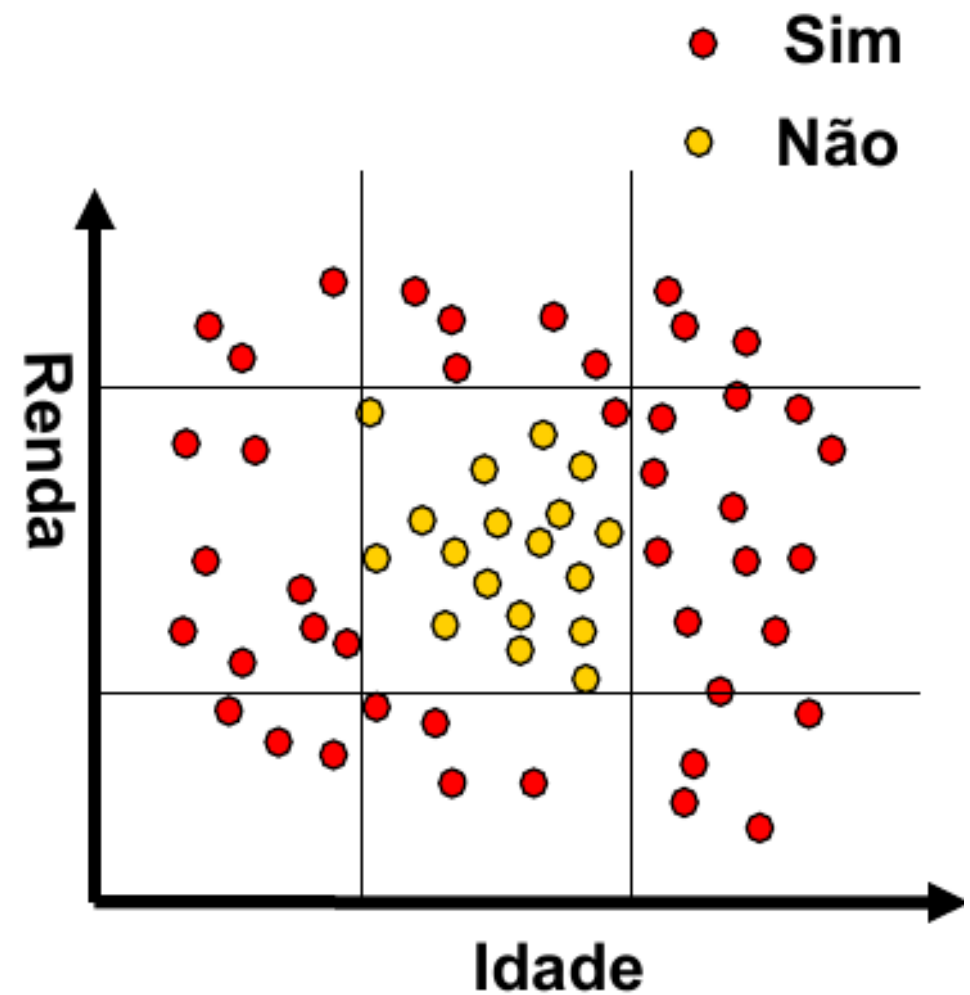
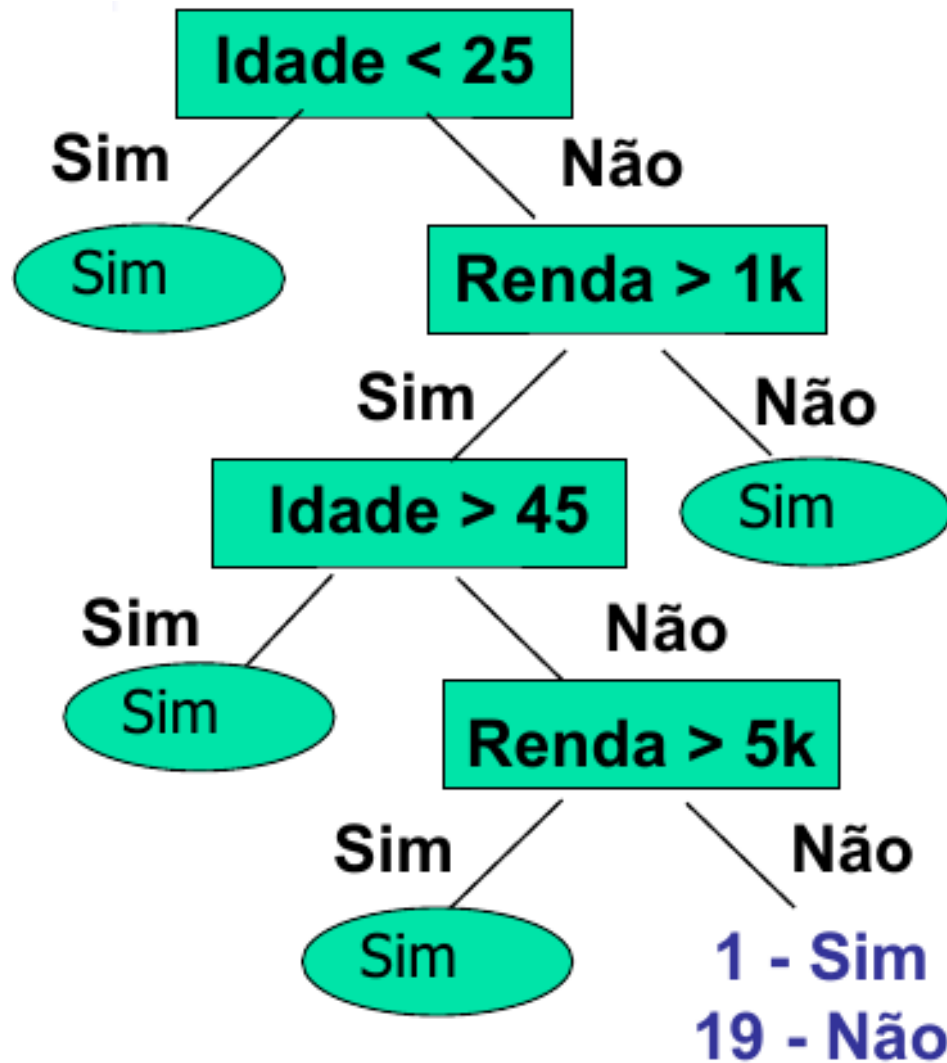
Indução₍₁₀₎



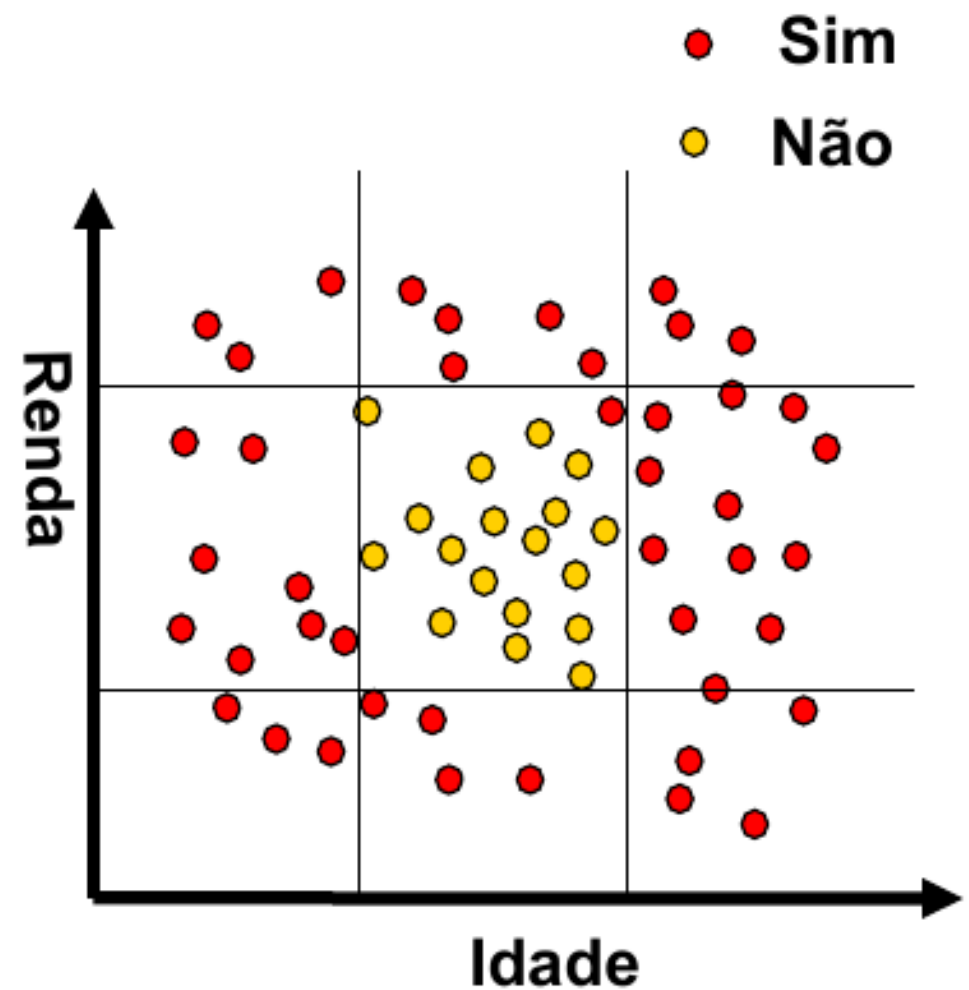
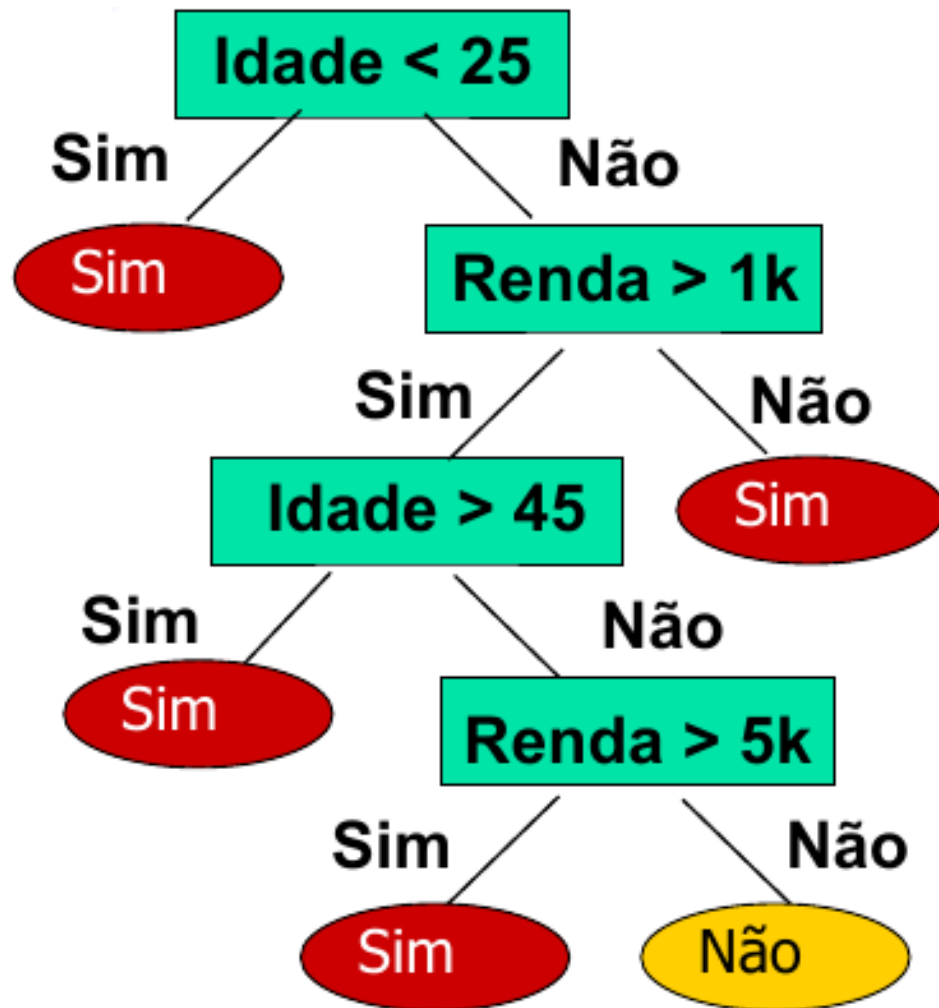
Indução ₍₁₁₎



Indução ₍₁₂₎



Indução ₍₁₃₎



Vantagens e Desvantagens

- Vantagens
 - Indica quais atributos são mais importantes para a classificação
 - Atributos de entrada podem ser numéricos ou categóricos
- Desvantagem: hipercubos são paralelos aos eixos
 - Tem mais dificuldade de analisar classes cujos atributos são proporcionais ou inversamente proporcionais

Algoritmo C4.5 (J48)

- Indução simples por busca em profundidade
- Usa ganho de informação
- Ordena atributos contínuos em cada nó
- Todos os dados precisam caber na memória principal
- Inadequado para grandes conjuntos de dados
 - Ordenação será feita fora da memória principal

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 a decisão para esses nós
 - Reduz capacidade de generalização

Overfitting (2)

- Navalha de Ockham (Ockham's razor)
 - Quanto mais simples a solução, melhor
 - Preferir as hipóteses mais simples
 - Quando hipótese mais simples explica os dados, é pouco provável que seja coincidência
 - Explicação dos dados por uma hipótese mais complexa pode ser apenas uma coincidência
 - Árvore de decisão pode ser simplificada por poda

Poda de Árvore

- 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

Algoritmo Simples de Poda

Percorrer a árvore em profundidade

Para cada nó i de decisão

E_i = erro no nó

ES = soma dos erros nos nós descendentes

Se $E_i \leq ES$

Então nó E_i é transformado em nó folha

Usar conjunto de validação para a poda

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

Exercício ₍₂₎

- Usando medidas Gini e GiniP, induzir uma árvore de decisão capaz de distinguir:
 - Pacientes potencialmente saudáveis (slide anterior)
 - Pacientes potencialmente doentes
- Testar a árvore para novos casos
 - (Luis, não, não, pequenas, sim)
 - (Laura, sim, sim, grandes, sim)

Conclusão

- Introdução
- Algoritmo de Hunt
- Medidas para selecionar divisão de atributos
- Ponto de referência
- Critério de parada
- Espaço de hipóteses
- Poda

Agradecimentos/referências

- Notas de aula do Prof. André de Carvalho (USP)