Mineração de Dados Aula 3

Especialização em Ciência de Dados e suas Aplicações





Classificação

Classificação: Definição



- Dado uma coleção de registros (treino)
 - Cada registro é caracterizado pela tupla (x,y), onde
 x é um conjunto de atributos e y é o rótulo da classe

Tarefa:

 Aprender um modelo que mapeia cada conjunto de atributo x em uma das classes predefinidas y

Exemplos



Tarefa	Conjunto de atributos (x)	Rótulos das classes (y)
Categorização de emails	Características (features) extraídas das mensagens e cabeçalho	spam ou regular
Identificação de tumor	Features extraídas resonâncias magnéticas	Malígno ou benígno
Categorizar textos	Features extraídas do texto	Cultura, Lazer, Esportes e Educação

Abordagem geral

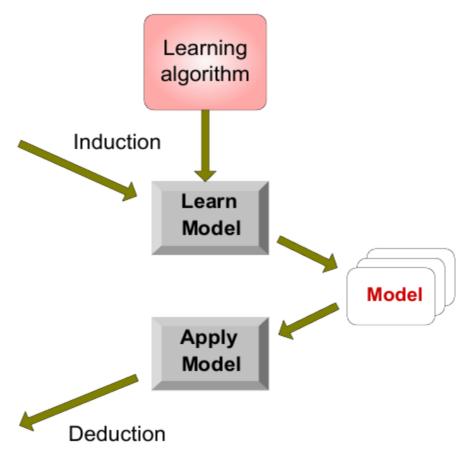


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Técnicas de classificação



- Métodos baseados em árvore de decisão
- Naïve Bayes e Redes Bayesianas
- Redes Neurais Artificiais

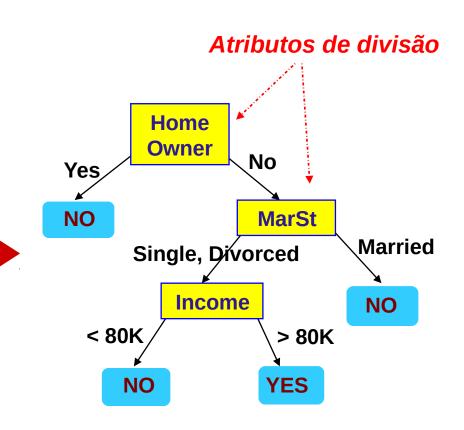
— ...

Exemplo de árvore de decisão



Categórico Contínuo

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Dados de treinamento

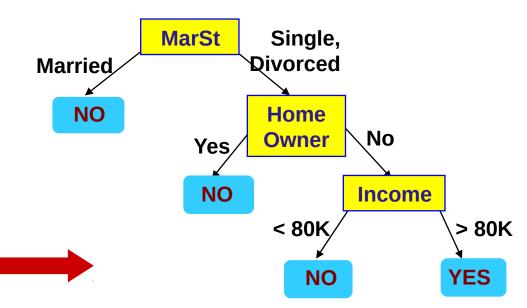
Modelo: Árvore de decisão

Outra árvore de decisão





ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Pode haver mais de uma árvore para o mesmo conjunto de dados

Dados de treinamento

Abordagem geral

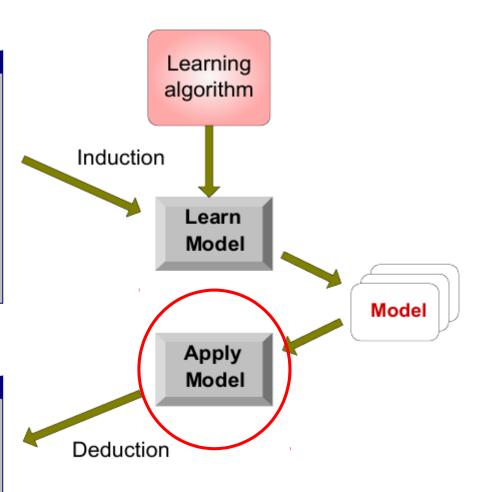


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

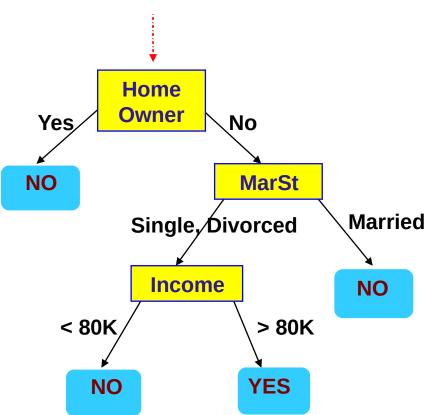
Test Set



Aplicar o modelo aos dados de teste Urp



Iniciar da raiz da árvore



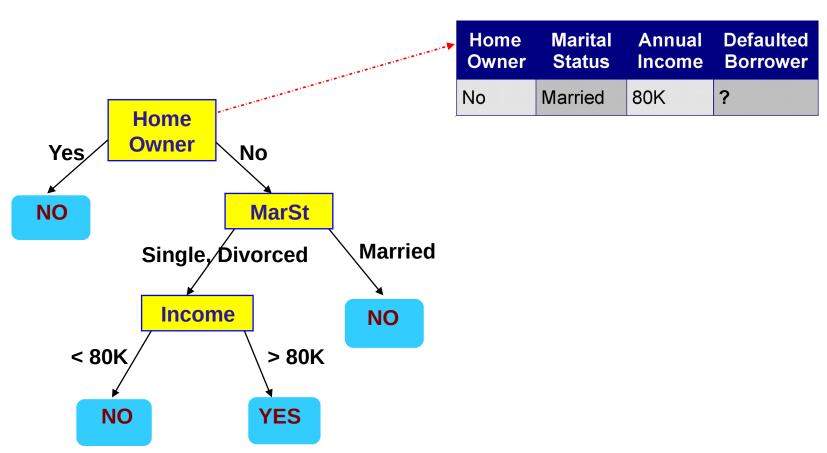
Dados de teste

			Defaulted Borrower
No	Married	80K	?

Aplicar o modelo aos dados de teste Urp



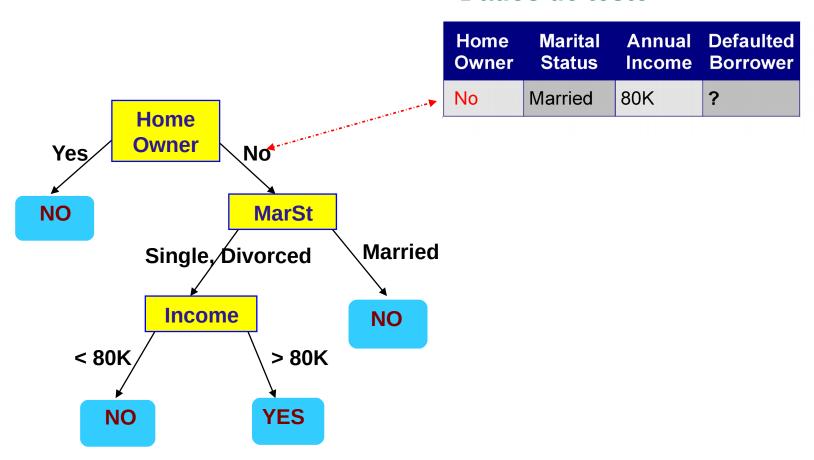




Aplicar o modelo aos dados de teste Urpe



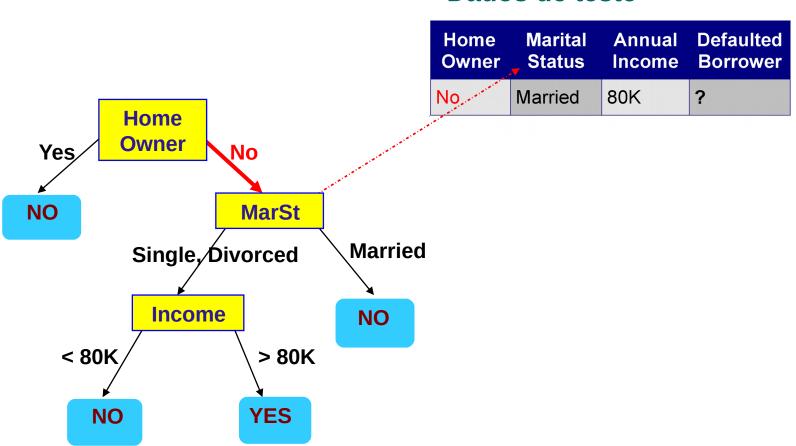
Dados de teste



Aplicar o modelo aos dados de teste Urper

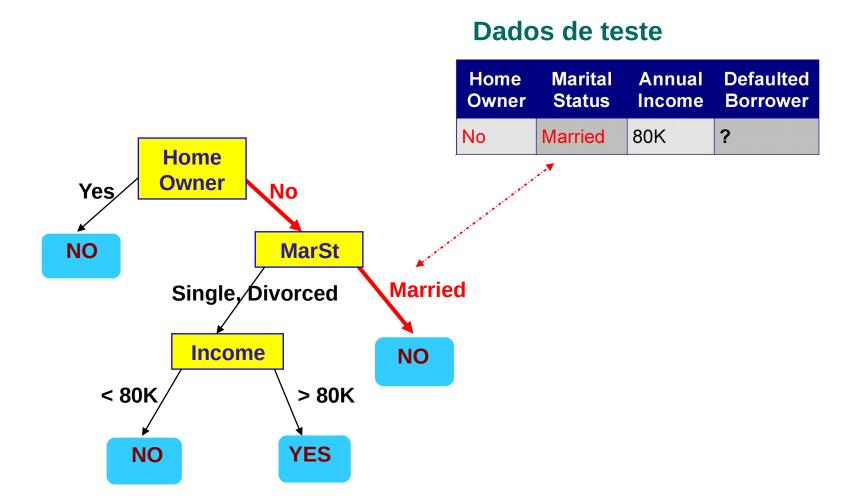






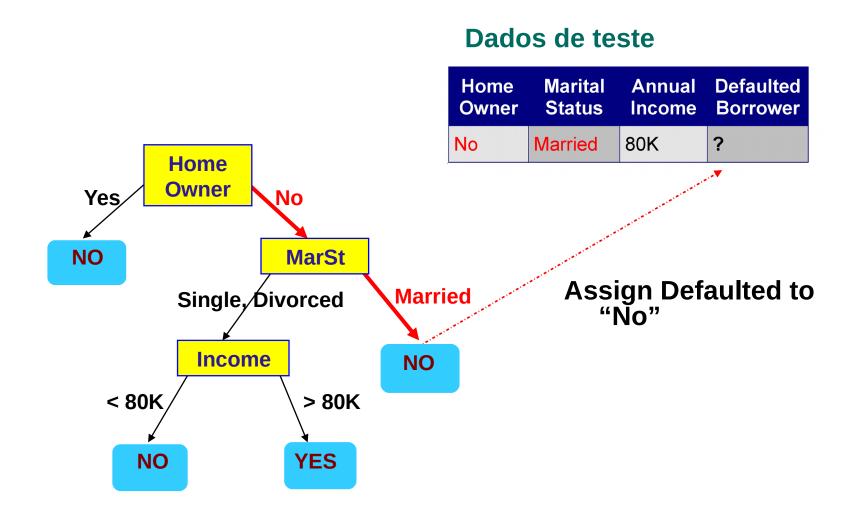
Aplicar o modelo aos dados de teste Urpe





Aplicar o modelo aos dados de teste Urpr





Árvore de decisão

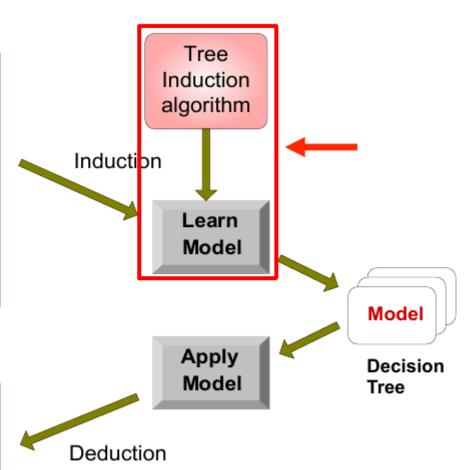


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Indução da árvore de decisão

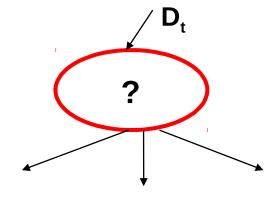


- Vários algoritmos :
 - Algoritmo de Hunt (um dos primeiros)
 - CART
 - ID3, C4.5, C.5
 - SLIQ, SPRINT



- Seja D_t o conjunto de treino que alcança o nó t
- Procedimento geral:
 - Se D_t contém registos que pertencem à mesma classe y_t, então t é um nó folha rotulado como y_t
 - Se D_t contém registros que pertencem a mais de uma classe, usar um atributo para dividir os dados em subconjuntos menores. Recursivamente aplicar o procedimento para cada subconjunto.

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes





Defaulted = No

(7,3)

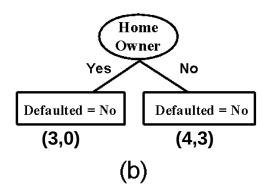
(a)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



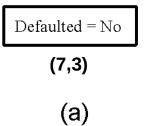
Defaulted = No **(7,3)**

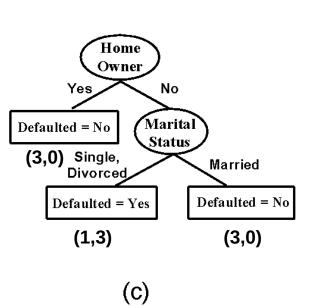
(a)

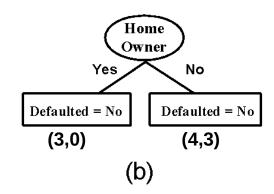


ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



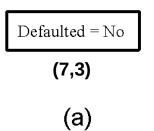


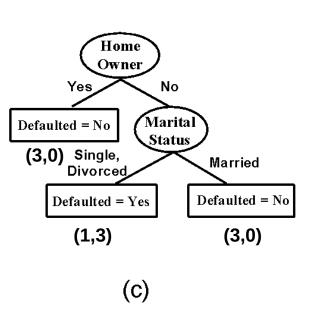


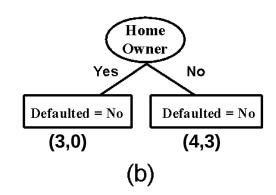


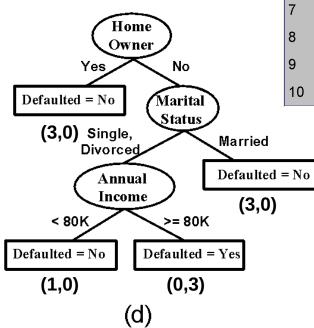
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes











ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Como determinar a melhor divisão UTr



- Estratégia gulosa:
 - Nós com distribuição de classes mais **puras** são preferidos
- Precisa de uma medida de impureza:

C0: 5

C1: 5

Alto grau de impureza

C0: 9

C1: 1

Baixo grau de impureza

Medidas de impureza



- Índice de Gini
- Entropia
- Erro de classificação

Encontrando a melhor divisão

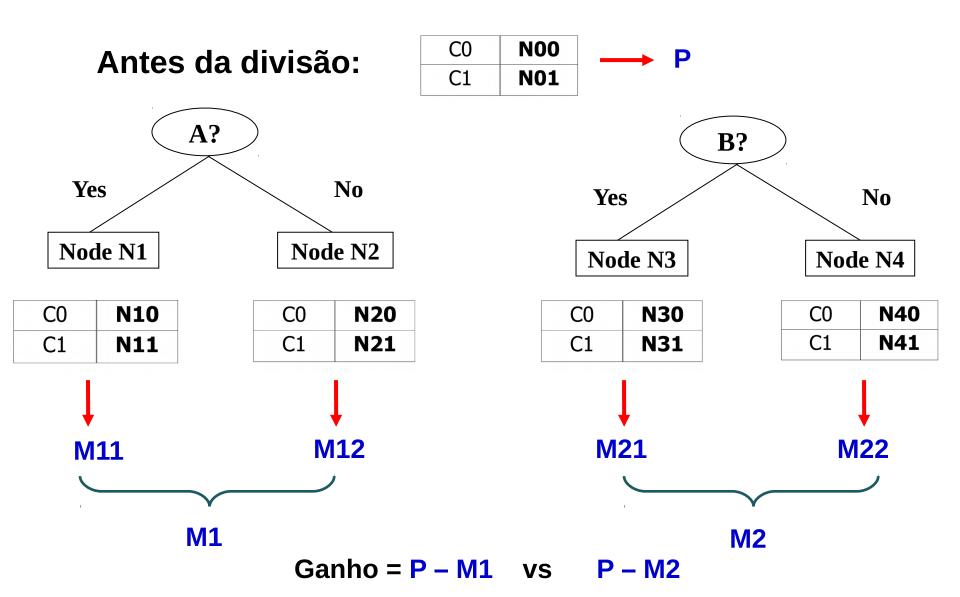


- 1. Computar a impureza (P) antes de dividir
- 2. Computar a impureza (M) depois de dividir
 - Computar a impureza para cada nó filho
 - M é a impureza ponderada para os filhos
- 3. Escolha o atributo que produz o maior ganho

Ganho = P - M

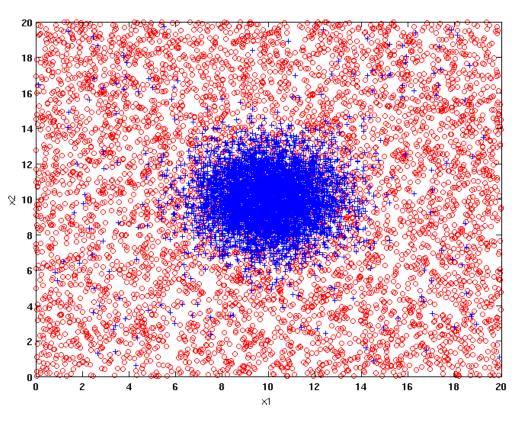
Encontrando a melhor divisão





Exemplo





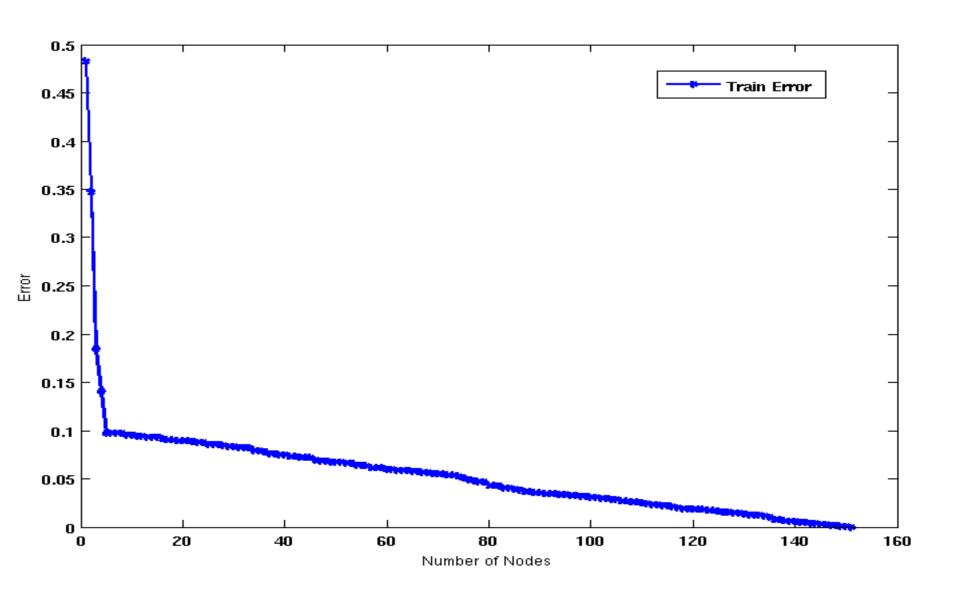
Problema com duas classes:

- +: **5200** instances
- ullet 5000 instances generated from a Gaussian centered at (10,10)
- *200 noisy instances added
- o: 5200 instances
- · Generated from a uniform distribution

10 % para treino e 90% para teste

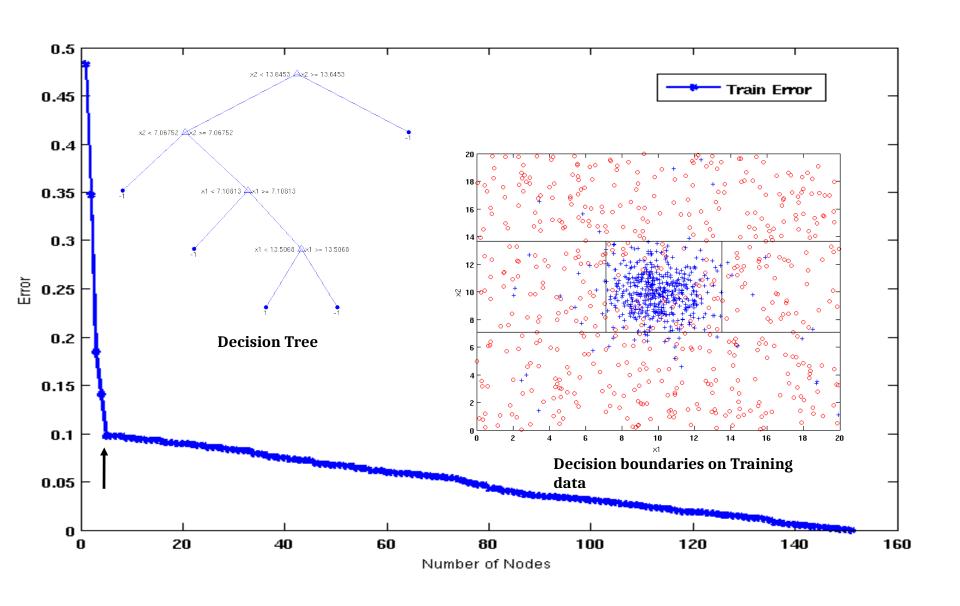
Aumentando o número de nós





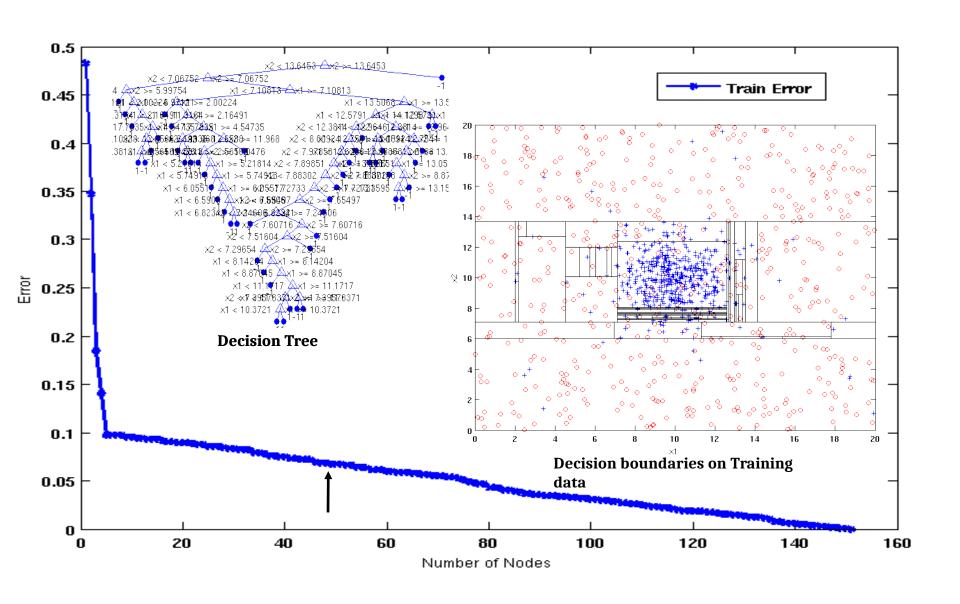
Árvore com 4 nós





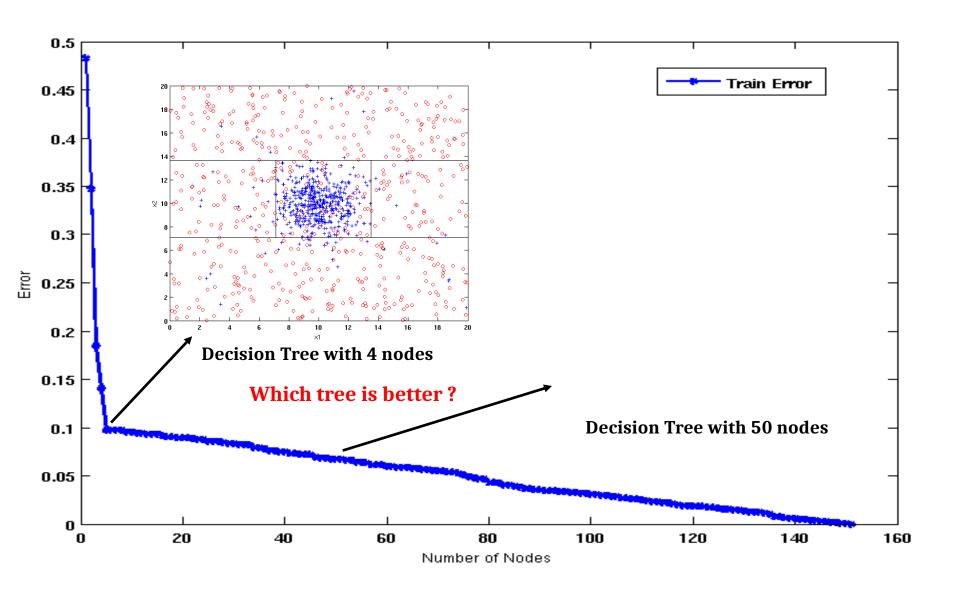
Árvore com 50 nós





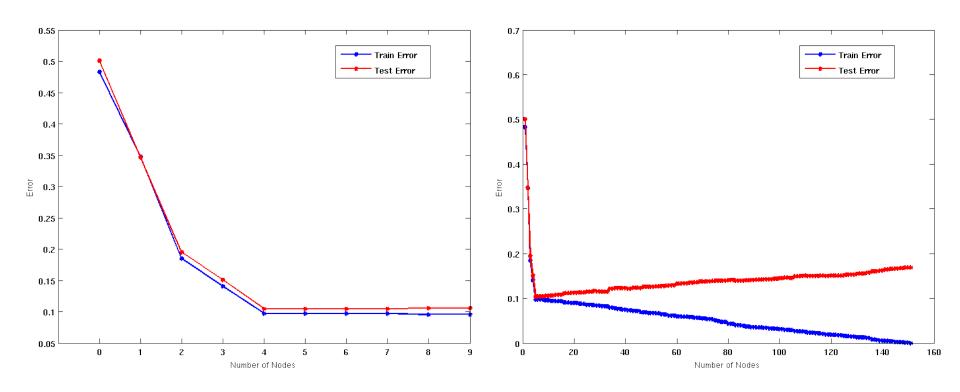
Qual a melhor?





Model Overfitting

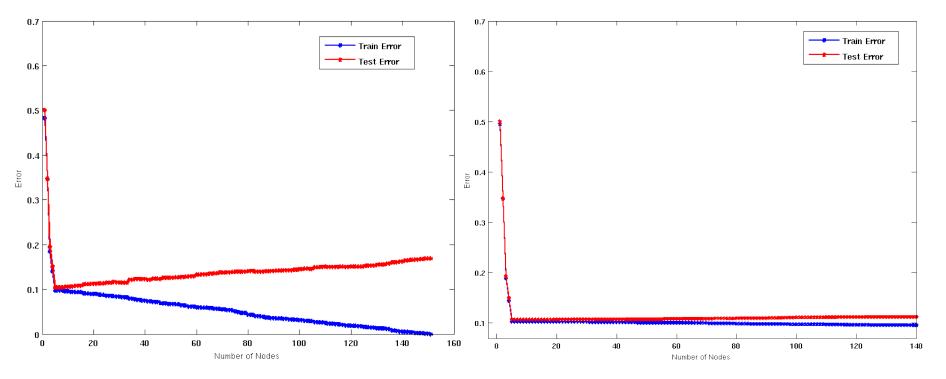




Underfitting: o modelo é muito simples, os erros de treino e teste são altos Overfitting: o modelo é muito complexo, o erro de treino é pequeno, mas o erro de teste é alto

Model Overfitting





Usando o dobro de instâncias de treino

- Se o conjunto de treino é subrepresentado, erros de teste aumentam e erros de treino diminuem ao aumentar os nós
- Aumentando o tamanho do conjunto de treino reduz a diferença entre os erros de teste e treino em um determinado número de nós

Métricas de avaliação



Exemplo: acc(M) = 90%

YES = tem-câncer (4 pacientes)
NO = não-tem-câncer (500 pacientes)

Classes "nãobalanceadas"

- Classificou corretamente **454** pacientes que não tem câncer
- Não acertou nenhum dos que tem câncer

Bom classificador?

acc = # classificação correta / total de classificações

Métricas de avaliação

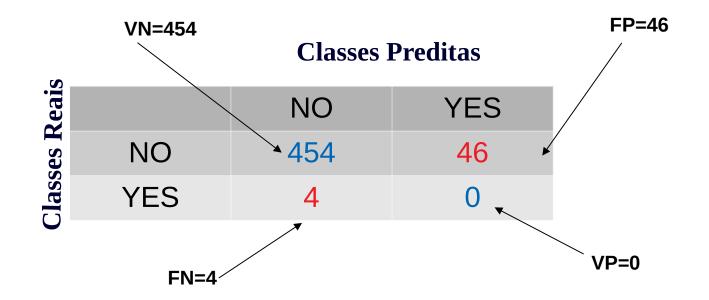


Exemplo: acc(M) = 90%

YES = tem-câncer (4 pacientes)
NO = não-tem-câncer (500 pacientes)

Classes "nãobalanceadas"

- Classificou corretamente **454** pacientes que não tem câncer
- Não acertou nenhum dos que tem câncer



Métricas de avaliação



% pacientes classificados <u>corretamente</u> como positivos dentre todos os que <u>realmente</u> são positivos

% pacientes classificados <u>corretamente</u> com câncer dentre todos os que foram classificados <u>com câncer</u>

Precisão e Recall: medidas originadas em *Recuperação de Informação* utilizadas em classificação, quando se lida com "classes não-balanceadas"

Exemplo prático 1



- Dataset Iris
 - Três tipos de plantas (classes):
 - Setosa
 - Versicolour
 - Virginica
 - Quatro atributos
 - Sépala: largura e comprimento
 - Pétalas: largura e comprimento



Courtesy of USDA NRCS Wetland Science Institute.

Exemplo prático 1

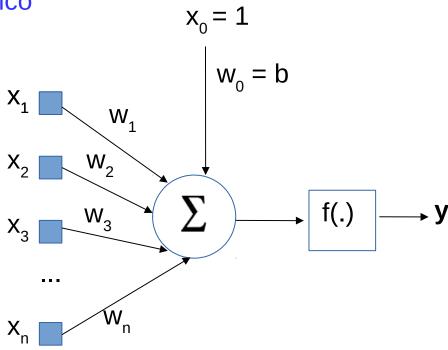


arvoreSimples.r

Redes Neurais Artificiais





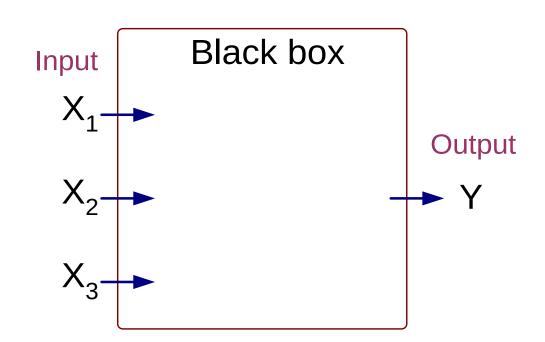


Neurônio artificial do tipo perceptron x = entrada (dendrito) e w = peso sináptico função somatória = corpo celular f() = função de ativação que gera a

f(.) = função de ativação que gera a saída no axônio **y**

Artificial Neural Networks (ANN)

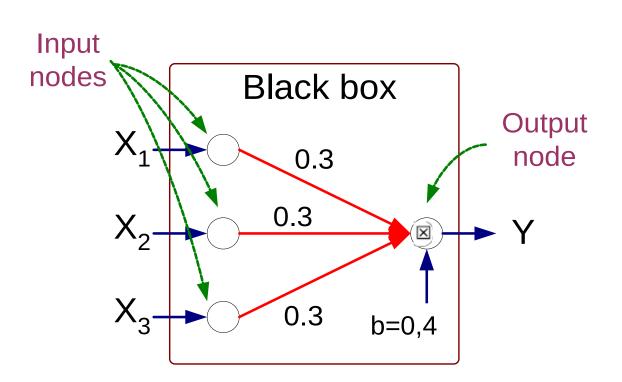
X ₁	X ₂	X_3	Υ
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1



Saída Y é 1 se pelo menos duas das três entradas são iguais a 1

Artificial Neural Networks (ANN)

X_1	X_2	X_3	Υ
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1



$$Y = degrau(0.3X_1 + 0.3X_2 + 0.3X_3 - 0.4)$$

onde degrau(v) =
$$\begin{cases} 1 \text{ se } v >= 0 \\ -1 \text{ se } v < 0 \end{cases}$$

Degrau = função degrau

Função lógica AND



Tabela da verdade

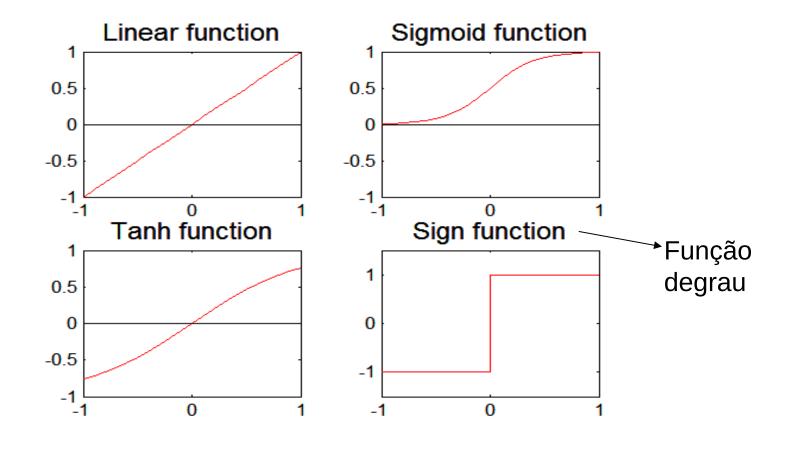
X ₁	X ₂	У
0	0	0
0	1	0
1	0	0
1	1	1

O neurônio projetado para resolver o problema deve:

- Utilizar como entrada x₁ e x₂
- Ponderar as entradas com os pesos sinápticos
- Realizar um somatório
- Aplicar uma função de ativação para produzir uma saída ŷ que deve ser igual à y

Vários tipos de funções de ativação **UTr**PR





Função lógica AND



Para o exemplo definimos a função sign (degrau):

- Ela produzirá saída **1** quando o campo induzido for > 0 e **-1** caso contrário

A regra de processamento do neurônio pode ser definida como:

- cálculo do sinal que entra no neurônio o é

$$v_o = \sum_{i=1}^{n} (x_i * w_{oi}) + b_o$$

- x_i é o valor de entrada
- a saída do neurônio é $y_0 = f(v_0)$

Função lógica AND (aprendizado) Tr

Inicializar os pesos $(w_0, w_1, ..., w_n)$

Repetir

Para cada instância de treino (x_i, y_i)

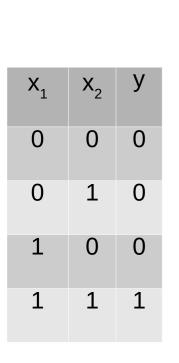
Computar $f(w, x_i)$

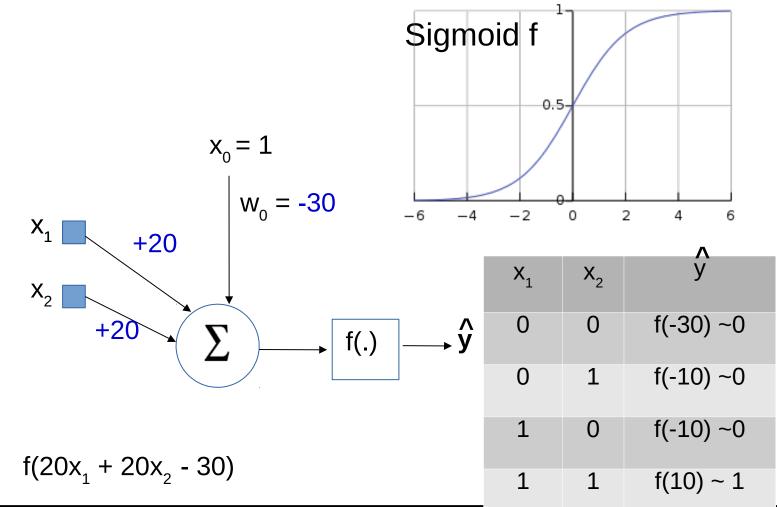
Atualizar os pesos: Taxa de aprendizado
$$w^{(k+1)} = w^{(k)} + \lambda \left[y_i - f(w^{(k)}, x_i) \middle| x_i \right]$$
 Erro

Até que a condição de parada seja atingida

Função lógica AND (aprendizada) Tr

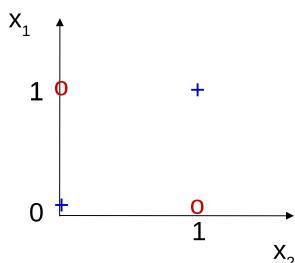
Exemplo com outra função de ativação:







X ₁	X ₂	У
0	0	1
0	1	0
1	0	0
1	1	1

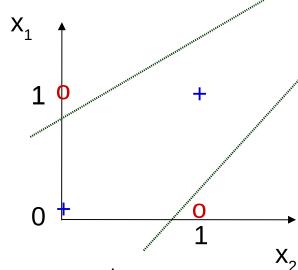


Não é linearmente separável

Ou seja, classes que podem ser separadas por uma reta (ou um hiperplano)



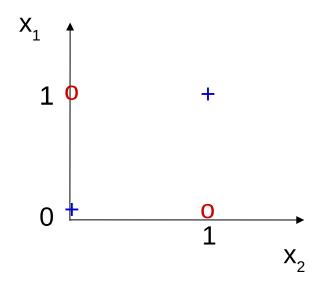
X ₁	X ₂	У
0	0	1
0	1	0
1	0	0
1	1	1

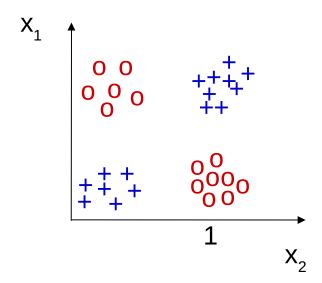


Não é linearmente separável Ou seja, classes que podem ser separadas por uma reta (ou um hiperplano)

Precisamos combinar mais de um neurônio (possibilita combinar retas)







Problema à esquerda é uma versão simplificada (mais fácil de analisar) do problema a direita

Rede neural de múltiplas camadas



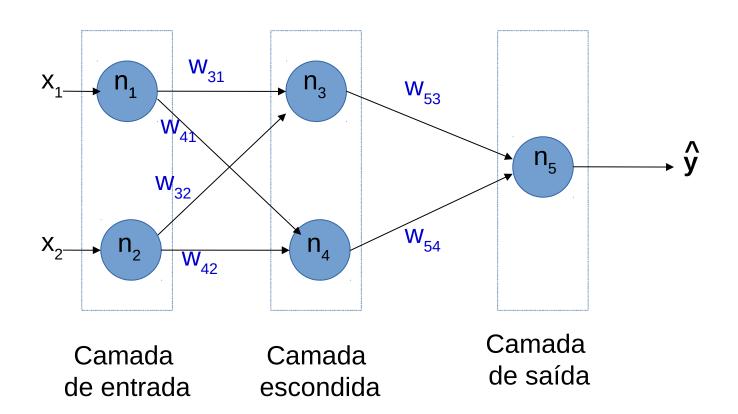
Possibilidade de combinação de neurônios em uma rede de múltiplas camadas

Possibilita a obtenção de estruturas mais complexas

Pode ser útil na resolução de tarefas que envolvem superfícies de decisão não-lineares

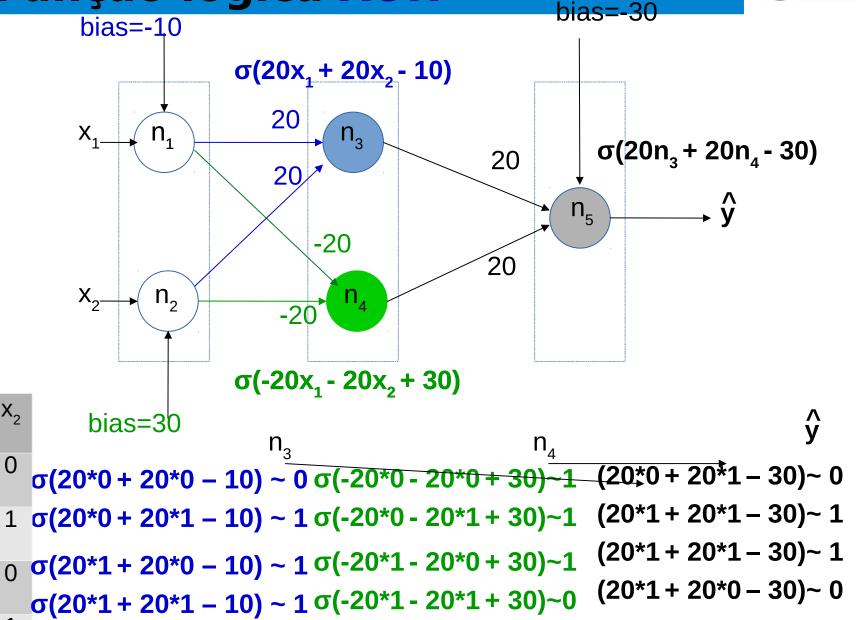
Rede neural de múltiplas camadas





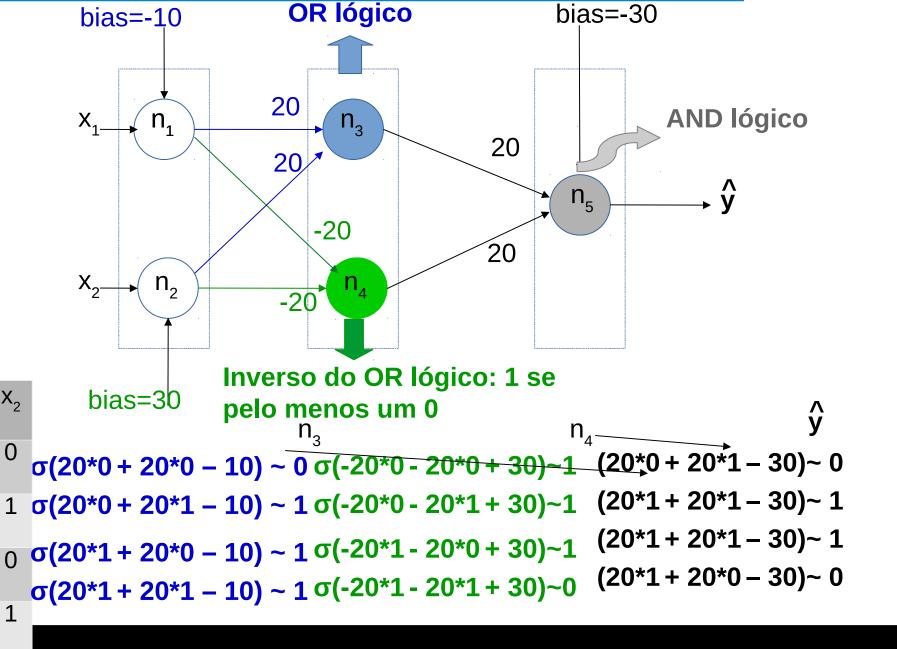
0





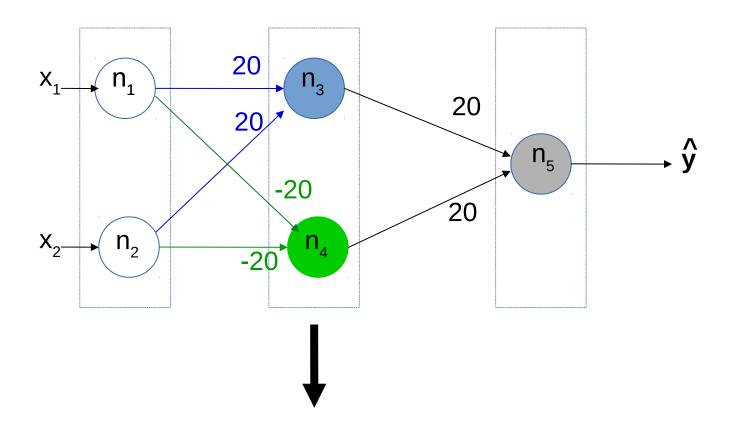
0





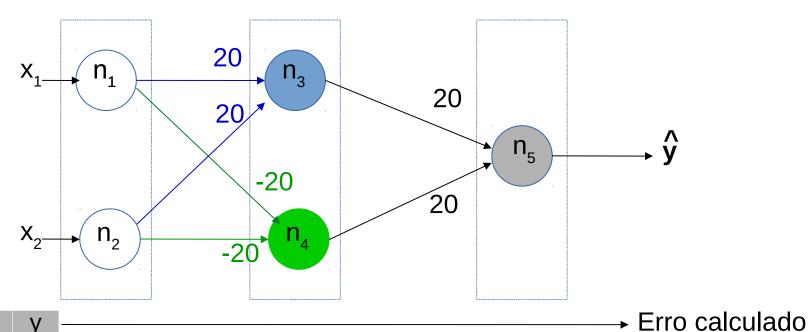
bias=-30





Não há como calcular o erro nos neurônios da camada escondida de forma direta, pois não existe a resposta deseja para tais neurônios.





Erro calculad	у	X ₂	X ₁
Erro propagado	1	0	0
Pesos reajustados levemente	0	1	0
O processo é repetido para todos as entradas e saídas até que o erro seja pequeno ou outra condição imposta.	0	0	1
Após esse processo a rede é considerada treinada	1	1	1

Vantagens e Desvantagens



- Fase de treinamento demorada. As iterações no *backpropagation* podem ocorrer centenas de milhares de vezes.
- Fraca interpretabilidade
- Muitos parâmetros determinados empiricamente
- Elevada tolerância a ruídos
- Resultados tendem a ser bons

Exemplo prático



redeNeuralSimples.r creditset.csv

Metodos de avaliação



Alguns dos métodos:

Holdout

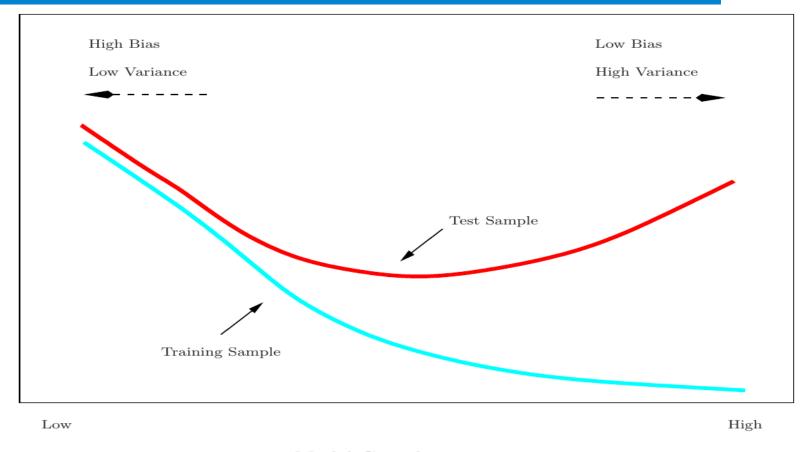
Reserva k% para treino (100-k)% para teste

Cross validation

Particionar os dados em k subconjuntos disjuntos k-fold: treinar em k-1 partições, testar na remanescente

Erro de treino vs erro de teste UTPR



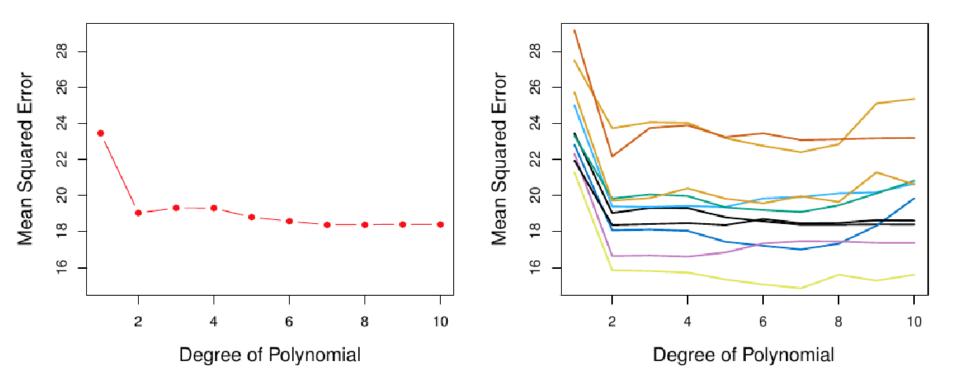


Model Complexity

Não podemos controlar o erro de teste com o de treino

Divisão em duas partes





- Esquerda: uma única divisão. Direita: várias.
- Variabilidade muito grande.

k-folds



Divide o dataset em partes iguais

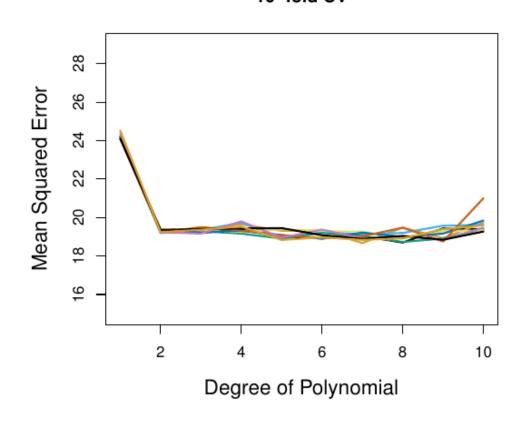
1	2	3	4	5
Valida ção	Treino	Treino	Treino	Treino

k-folds



· Variação do resultado fica menor.





K= 5 ou 10 são escolhas típicas

Exercício prático



Fazer o exercício 3 disponível no moodle.

Agradecimentos



Alguns slides foram derivados/inspirados em:

- Livro Introduction to Data Mining - Tan, Steinbach, Kumar. 2018.