



BY

“Prevejo um tempo em que os humanos serão para as máquinas, o que os cães são para nós. E eu estou torcendo pelas máquinas” (Claude Shannon).

Árvores de Decisão

Paulo Ricardo Lisboa de Almeida



Classificação Supervisionada

Em aprendizado de máquina, um modelo de classificação prediz a classe de determinada instância.

A instância pode ser uma foto, um sinal de áudio, um vetor com dados, ...

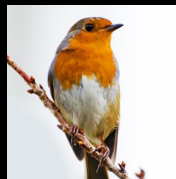
A classe pode ser um item que desejamos reconhecer (ex.: essa é uma foto de um cachorro, pessoa, pássaro, ...).

Classificação Supervisionada

Em aprendizado de máquina, um modelo de classificação prediz a classe de determinada instância.

A instância pode ser uma foto, um sinal de áudio, um vetor com dados, ...

A classe pode ser um item que desejamos reconhecer (ex.: essa é uma foto de um cachorro, pessoa, pássaro, ...).

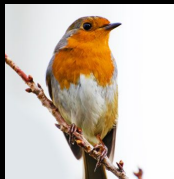


Classificação Supervisionada

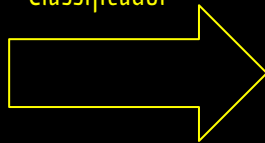
Em aprendizado de máquina, um modelo de classificação prediz a classe de determinada instância.

A instância pode ser uma foto, um sinal de áudio, um vetor com dados, ...

A classe pode ser um item que desejamos reconhecer (ex.: essa é uma foto de um cachorro, pessoa, pássaro, ...).



Classificador



Classificação Supervisionada

Modelo Supervisionado.

Um modelo supervisionado aprende automaticamente com **dados rotulados de treino**.

Exemplo: são dadas imagens juntamente com seus rótulos criados por humanos para que o modelo possa gerar suas representações internas.

Problema

Considere um problema onde os dados de pressão arterial e nível de glicose são medidos para uma pessoa.

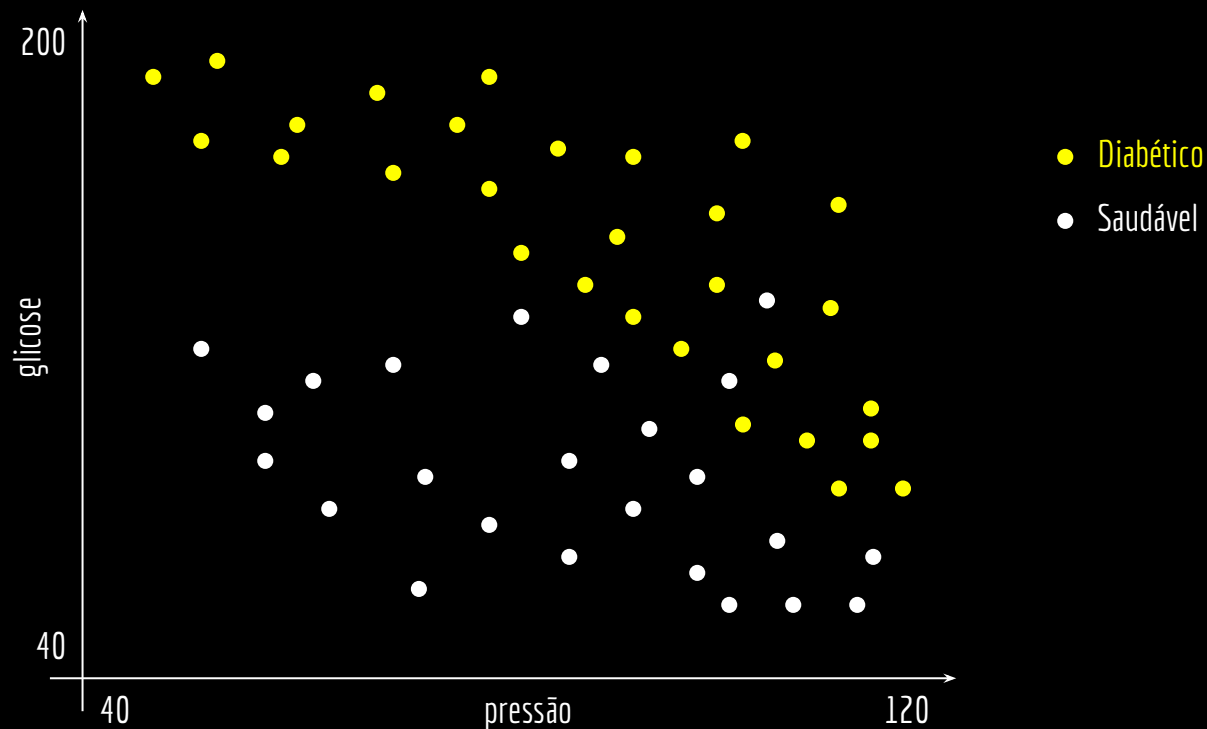
Esses dados podem ser lidos, por exemplo, via um *smart-watch*.

Desejamos classificar as pessoas entre “com diabetes” (classe 1) e “saudáveis” (classe 0).

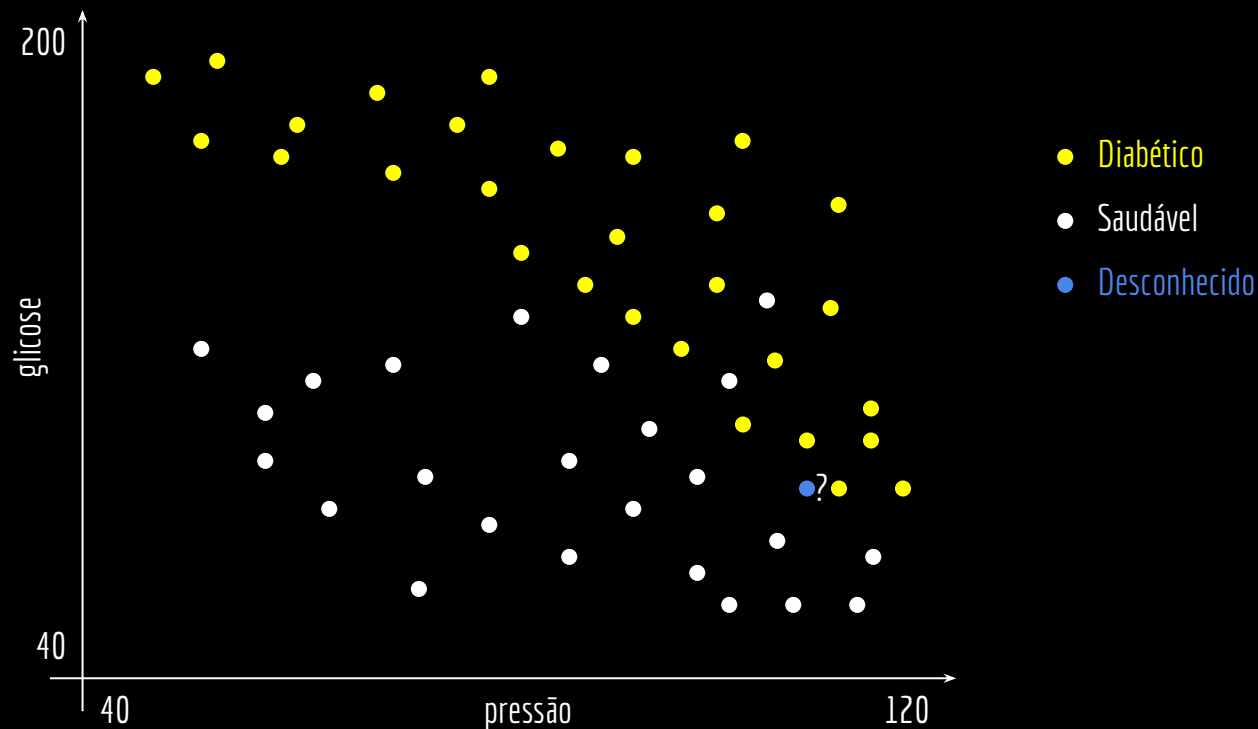
Temos então um problema binário.

Atenção: essa é uma simplificação da realidade. Classificar pessoas entre diabéticas e saudáveis não é nem de longe tão simples.

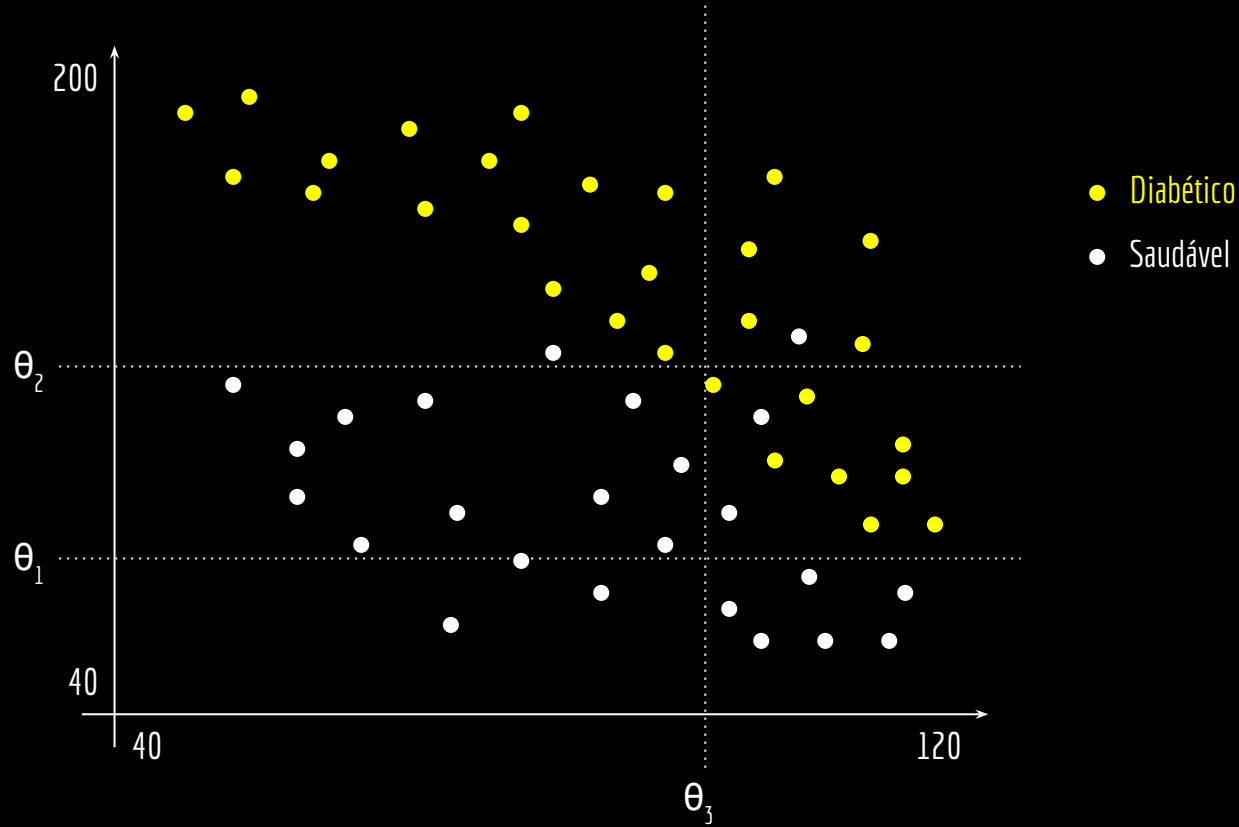
Dataset de Treinamento



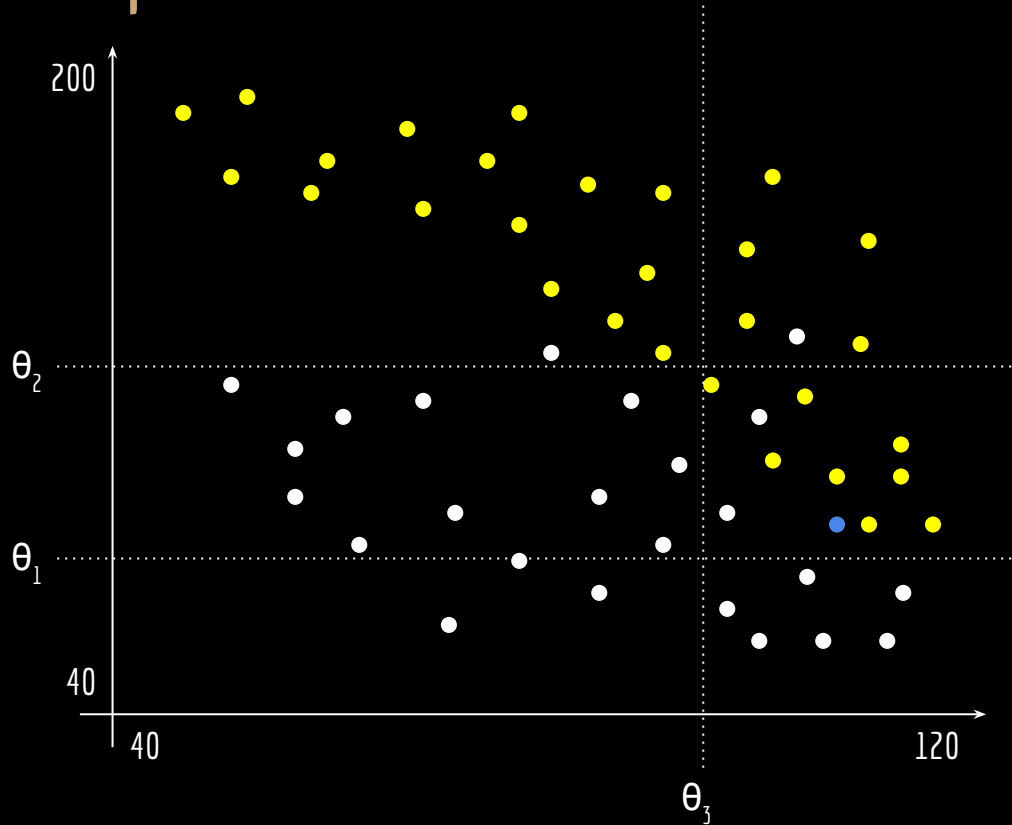
Dataset de Treinamento



Dividindo os dados



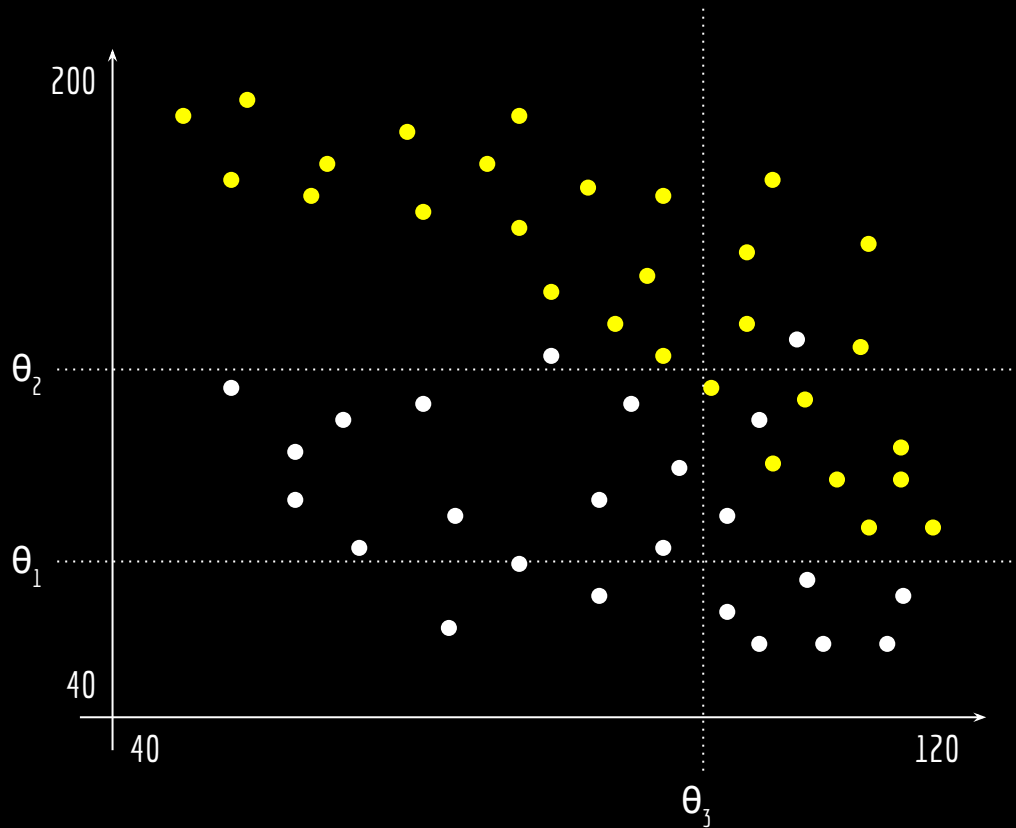
Classificando



- Diabético
- Saudável
- Classificado como Diabético

É maior que θ_3 no eixo x, e está entre θ_1 e θ_2 no eixo y, então está na área considerada como diabético.

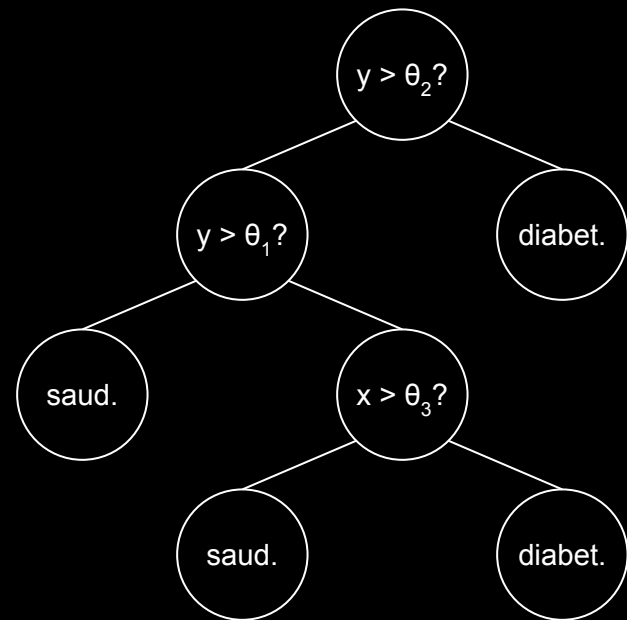
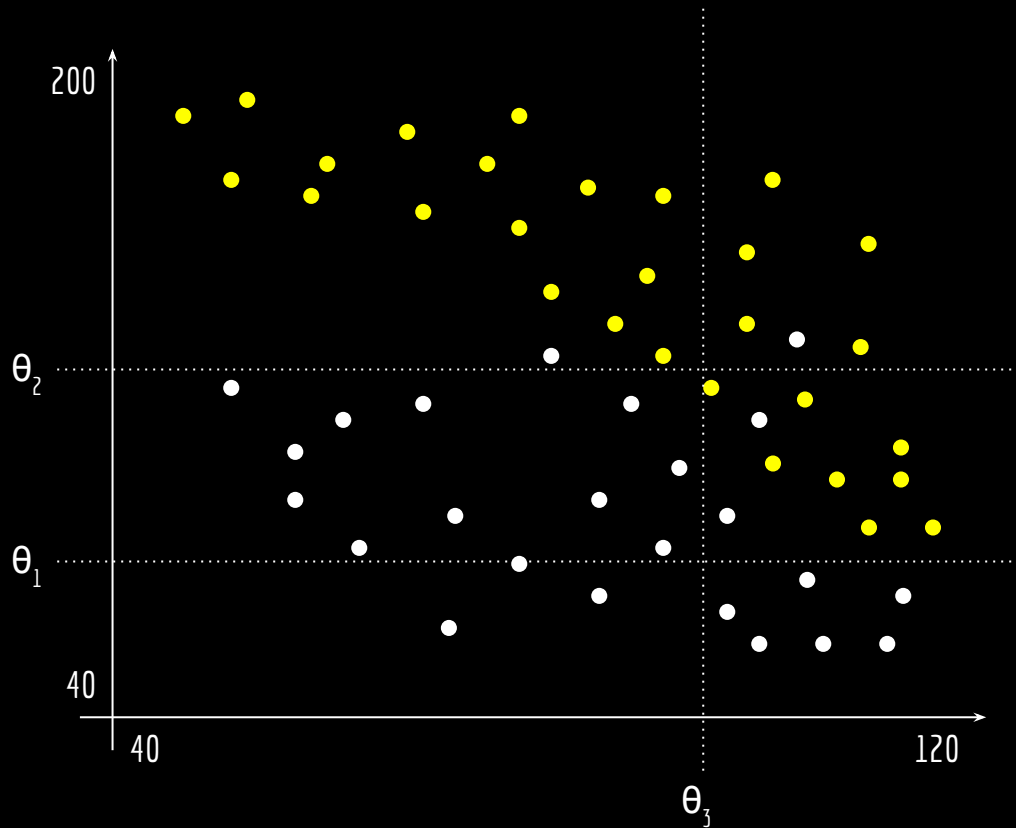
Árvore de Decisão



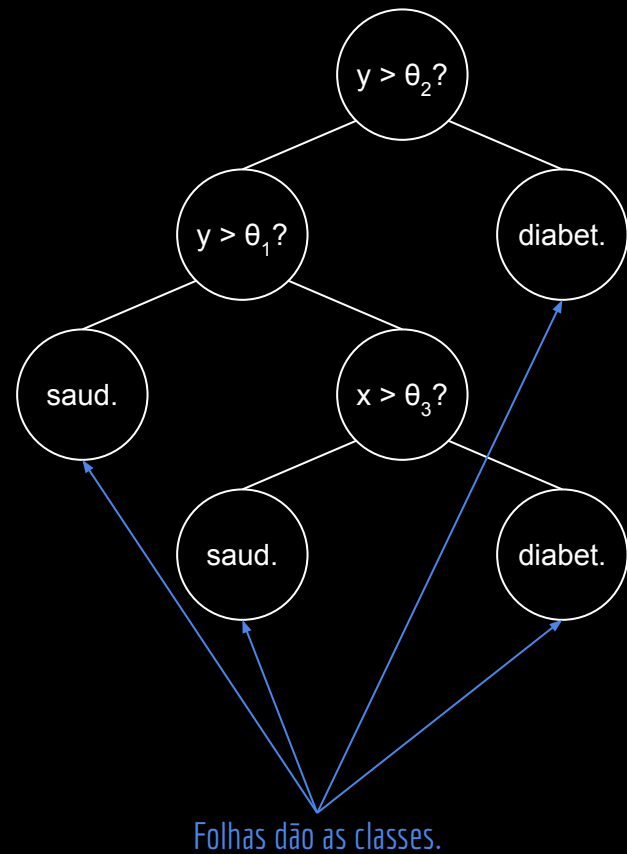
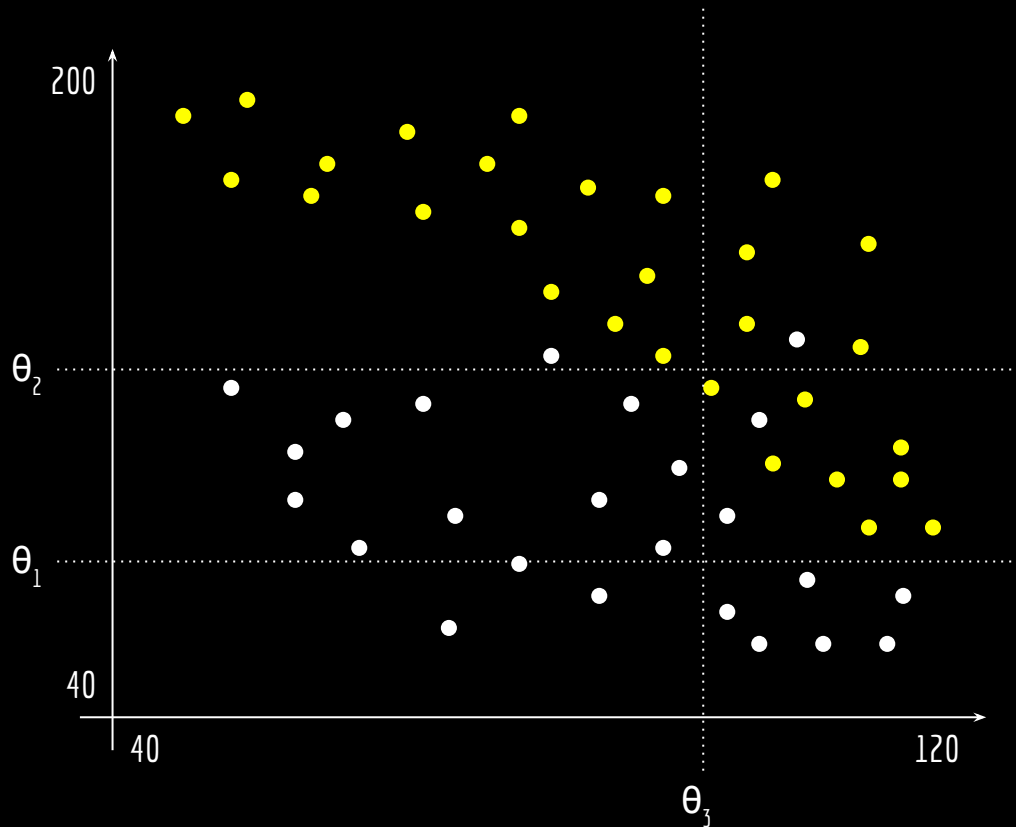
O problema pode ser modelado como uma árvore binária.

Algoritmo de **Árvore de Decisão**.

Árvore de Decisão



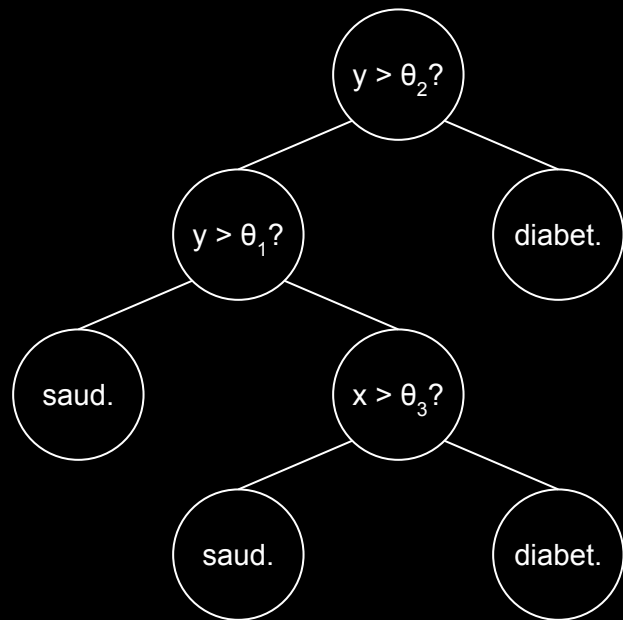
Árvore de Decisão



Criando a árvore

Precisamos criar a árvore.

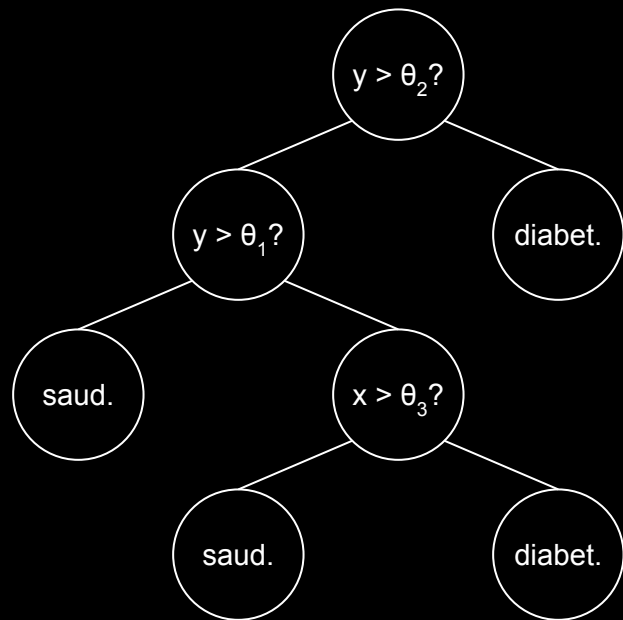
O que precisamos definir?



Criando a árvore

Precisamos criar a árvore.

- Definir a estrutura da árvore.
 - Ex.: por qual θ precisamos começar a busca?
- Definir os limiares (valores para cada θ).



Algoritmo

Definir a estrutura ótima da árvore é computacionalmente inviável mesmo para problemas pequenos.

Geralmente usamos algum algoritmo guloso.

Exemplos de algoritmos clássicos.

- CART
- ID.3
- C4.5 (Implementação Open Source no Weka se chama J48)

Vetores de característica

O exemplo anterior possuía apenas duas características.

Glicose e pressão.

Problema bidimensional.

No mundo real, geralmente encontramos problemas com mais dimensões.

Os dados que descrevem as instâncias de um problema são chamados de características.

Uma instância tem um vetor de características $x = [x_1, x_2, \dots, x_n]$, e uma classe $y \in [y_1, y_2, \dots, y_k]$.

Onde n é o número de características que descrevem o problema, e k é o número de classes possíveis.

Vetores de característica

Suponha que as características do problema são: *numGravidez*, *Glicose*, *Pressão*, *espessuraPele*, *Insulina*, *BMI*, *DPF* e *Idade*.

A classe y pode ser $0 = \text{saudável}$, e $1 = \text{diabetes}$.

Um vetor de características para treinamento pode ser:

$x = [6, 148, 72, 35, 0, 33.6, 0.627, 50, 1]$ e $y = 1$.

Temos agora um problema com 9 dimensões.

Dados de treinamento

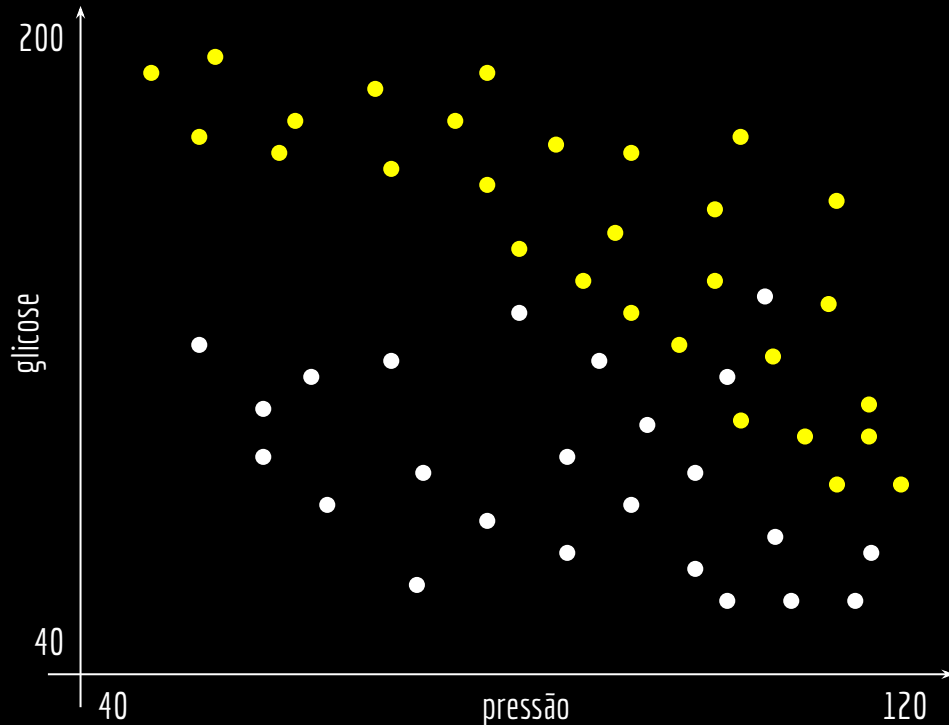
Supondo que temos k dados para treinamento, temos uma matriz de $k \times n$, onde cada linha da matriz representa um vetor de treinamento. E um vetor y^T de tamanho k .

Matriz de $k \times n$ onde cada linha é um vetor de características.

	n										
k	6, 148, 72, 35, 90, 33.6, 0.627, 50	1									
	1, 85, 66, 29, 100, 26.6, 0.351, 31	0									
	8, 183, 64, 0, 110, 23.3, 0.672, 32	1									
	1, 89, 66, 23, 94, 28.1, 0.167, 21	0									
	0, 137, 40, 35, 168, 43.1, 2.288, 33	1									
	5, 116, 74, 0, 101, 25.6, 0.201, 30	0									
...											

Vetor y^T com os rótulos das instâncias.

Impureza

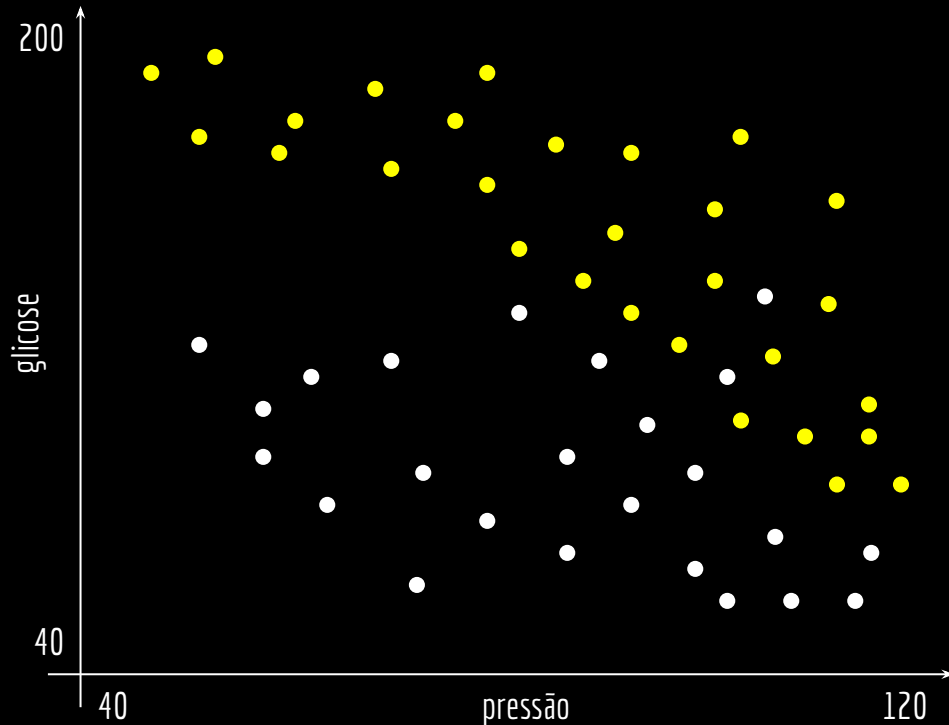


Em um problema **binário**, podemos calcular a impureza de determinada região como:

$$i(N) = P(y_1)P(y_2)$$

Onde $P(y_x)$ é a proporção de itens da classe 0 ou 1 na região.

Impureza



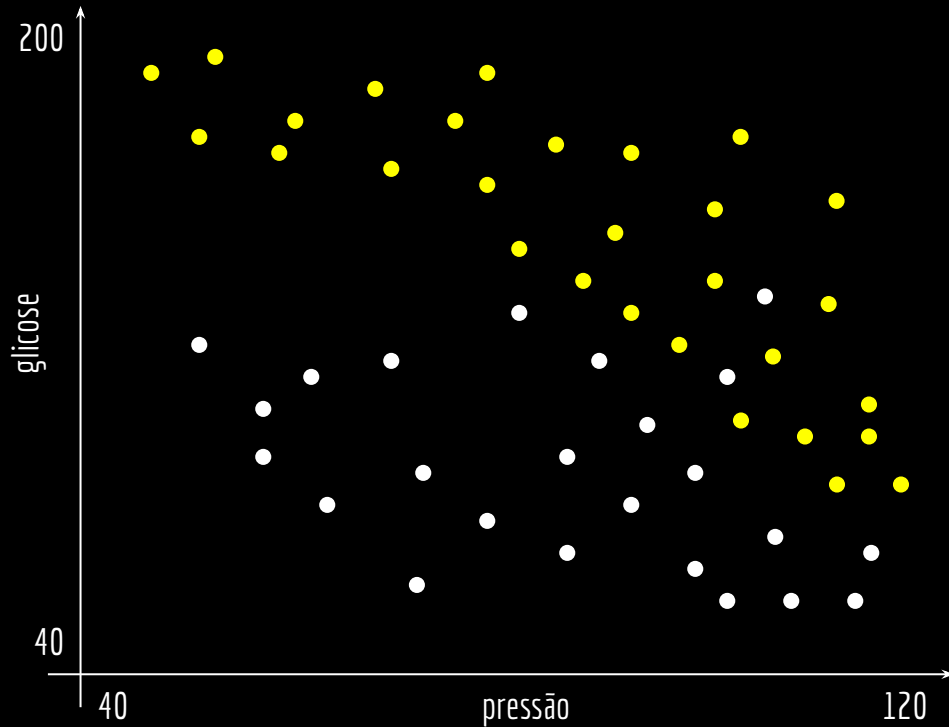
Em um problema **binário**, podemos calcular a impureza de determinada região como:

$$i(N) = P(y_1)P(y_2)$$

Onde $P(y_x)$ é a proporção de itens da classe 0 ou 1 na região.

Uma região pura terá $i(N) = 0$.

Impureza



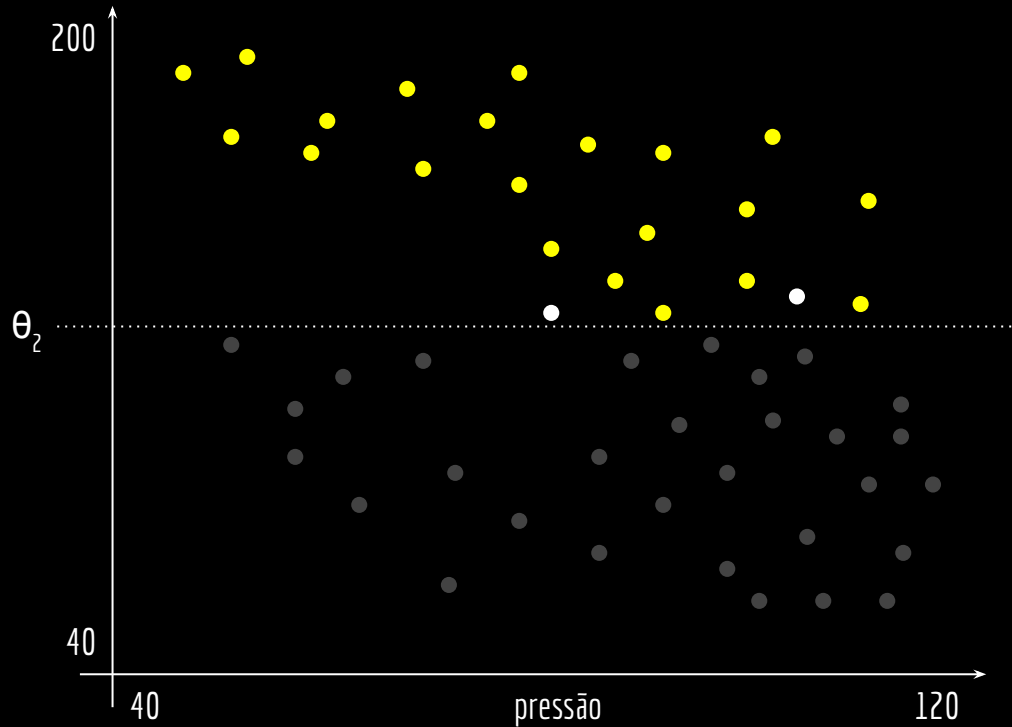
Contagens

● 28

● 24

$$i(N) = (28:52) * (24:52) = 0,25$$

Impureza



Contagens

● 21

● 2

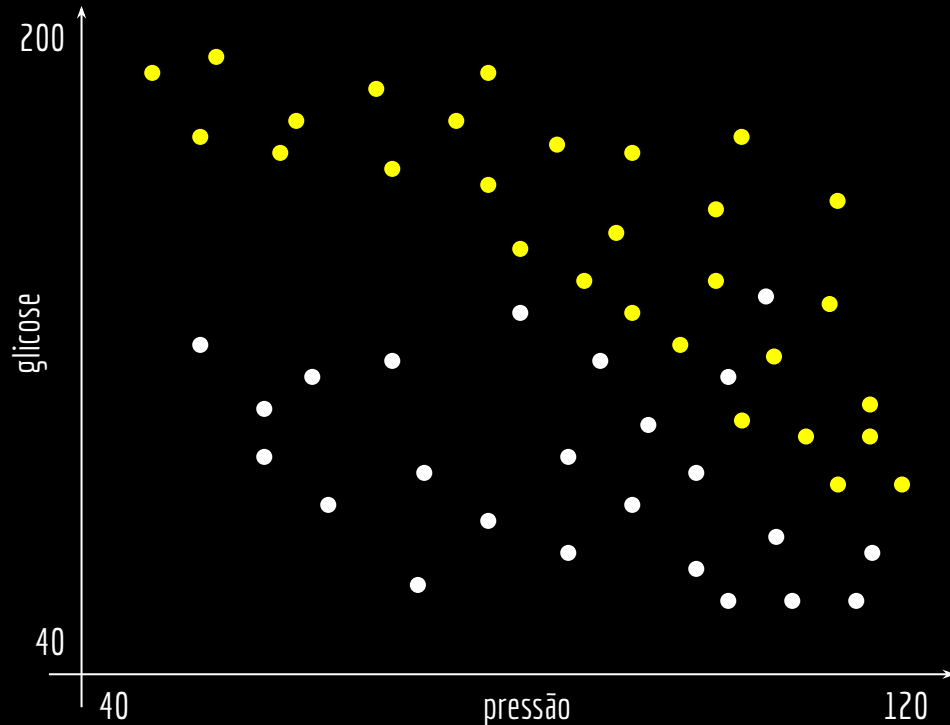
$$i(N) = (21 \div 23) * (2 \div 23) = 0,08$$

Índice Gini

A impureza pode ser generalizada para casos onde temos duas ou mais classes, criando o **Índice Gini**:

$$i(N) = \sum_{i \neq j} P(y_i)P(y_j) = 1 - \sum_j P^2(y_j)$$

Algoritmo - Ideia



Adicionar um novo nodo na árvore.

Para cada característica c

Para cada limiar de decisão possível

Calcular o Índice Gini

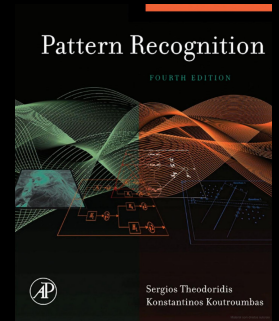
Escolher como a característica e limiar que levam ao menor Índice Gini.

Continuar criando nodos enquanto reduções no Índice Gini sejam significativas.

Atenção

Atenção: essa é uma simplificação de um dos muitos algoritmos possíveis que podem ser usados para se criar a árvore de decisão.

Veja uma versão completa de um algoritmo, e algumas provas em Theodoridis e Koutroumbas (2008).
Quarta Edição, Seção 4.20.



Árvores de Decisão

Árvores de decisão são propensas a overfitting (sobreajuste).

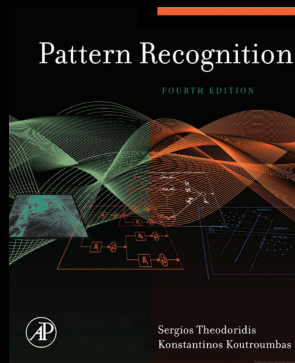
- O modelo pode se tornar muito complexo, com muitas fronteiras.
- No pior caso, cada região contém uma única instância.
- Problemas de generalização.

Exercícios

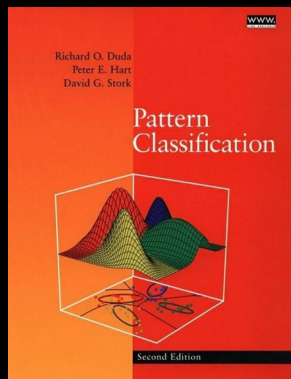
1. Faça testes com o programa de Árvore de Decisão disponibilizado no Moodle.
 - a. Use o dataset para previsão de Diabetes disponibilizado.
 - i. O dataset é do mundo real, e está disponível em www.kaggle.com/datasets/uciml/pima-indians-diabetes-database
 - ii. Cuidado ao analisar os resultados! O dataset é desbalanceado (possui mais itens da classe saudável do que da classe diabetes). Isso pode enviesar a acurácia.
 - b. Baixe outros datasets, e faça testes.

Referências

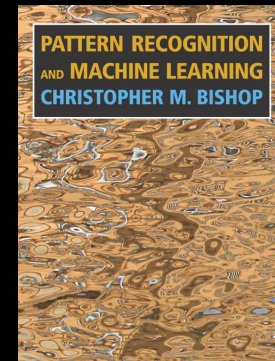
Theodoridis, S., Koutroumbas, K. Pattern Recognition. 4a ed. Elsevier Science. 2008.



Stork, D. G., Hart, P. E., Duda, R. O. Pattern Classification. Wiley. 2012.



Bishop, C. M. Pattern Recognition and Machine Learning. Springer. 2007.



scikit-learn.org/stable/modules/tree.html

www.kaggle.com/datasets/uciml/pima-indians-diabetes-database

www.datacamp.com/tutorial/decision-tree-classification-python

Licença

Esta obra está licenciada com uma Licença [Creative Commons Atribuição 4.0 Internacional](https://creativecommons.org/licenses/by/4.0/).