



**Universidade Federal da Paraíba**

Coordenação do Curso de Ciência de Dados e  
Inteligência Artificial



# Viabilidade da aprendizagem

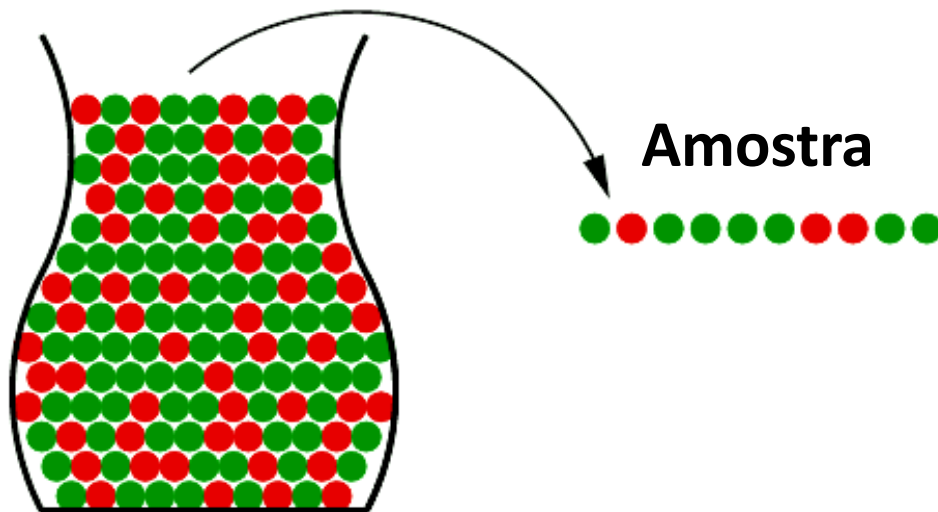
Prof. Dr. Bruno Pessoa

# Roteiro

- Relação entre amostra e população
- Desigualdade de Hoeffding
- Conexão com Aprendizagem de Máquina
- Hoeffding para múltiplas hipóteses
- Modelos baseados em distâncias
- Algoritmo k-NN

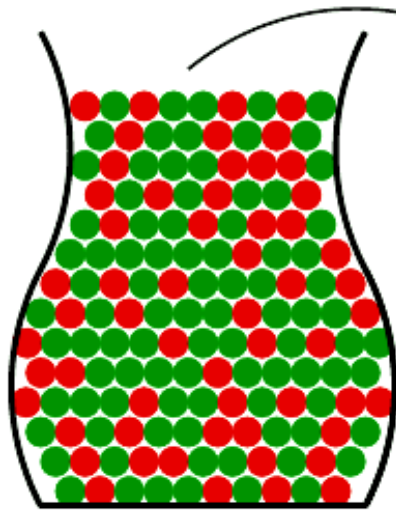
# Relação entre amostra e população

- Considere um pote com bolas verdes e vermelhas, onde:
  - $P(\text{pegar uma vermelha}) = L$
  - $P(\text{pegar uma verde}) = 1 - L$



- O valor de  $L$  é desconhecido.
- A fração de vermelhas na amostra é  $l$ .
- O tamanho da amostra é  $N$ .

# Relação entre amostra e população



$L$  = probabilidade de pegar uma bola vermelha

Amostra



$l$  = fração de bolas vermelhas

A amostra determina o que ocorre na população ou ocorre o inverso?

A amostra dá alguma pista sobre o que ocorre na população?

# Teoria da Probabilidade

- Em uma amostra grande o suficiente,  $l$  provavelmente se aproxima de  $L$ .

- Matematicamente,

$$P(|l - L| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

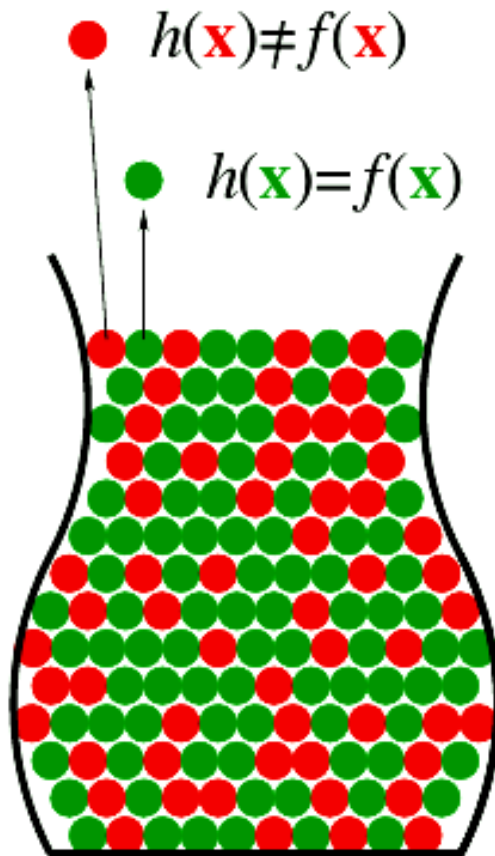
- Conhecida como **desigualdade de Hoeffding**.

# Desigualdade de Hoeffding

$$P(|\textcolor{red}{l} - \textcolor{red}{L}| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

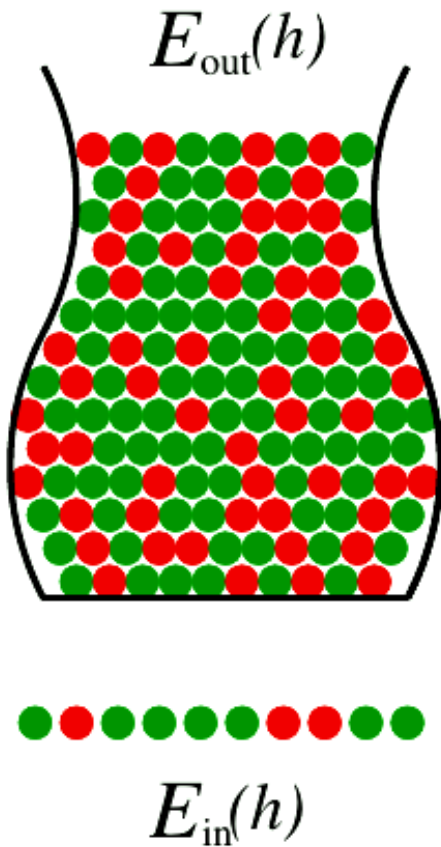
- **Considerações:**
  - O lado direito da inequação, que atua como um **limitante**, não depende de  $\textcolor{red}{L}$ .
  - $\textcolor{red}{l}$  é um componente aleatório da desigualdade.
  - Há claramente um tradeoff entre  $\epsilon$  e  $N$ .
  - Necessidade de aumentar o tamanho da amostra.

# Conexão com Aprendizagem de Máquina



- O ato de pegar uma amostra no pote equivale a escolher um hipótese aleatória para classificar os dados.
- $L$  é desconhecida assim como a função alvo  $f: X \rightarrow Y$ .
- Cada bola diz respeito a um ponto  $x \in X$ .
- Uma **bola verde** significa que o ponto foi classificado de forma correta por uma hipótese  $h(x)$ , e uma bola vermelha o contrário.
- $L$  e  $l$  correspondem à fração de bolas classificadas incorretamente na população e na amostra, respectivamente.
- A amostra consiste nos dados usados no treinamento.

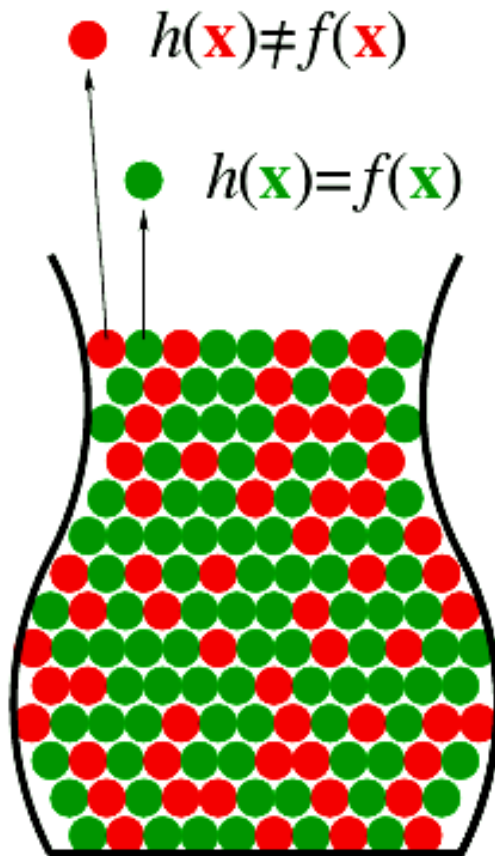
# Conexão com Aprendizagem de Máquina



- $l$  é será substituído por  $E_{in}(h)$ , denotando o erro dentro da amostra para uma hipótese  $h$ .
- $L$  é será substituído por  $E_{out}(h)$ , denotando o erro fora da amostra para uma hipótese  $h$ .



# Conexão com aprendizagem de máquina

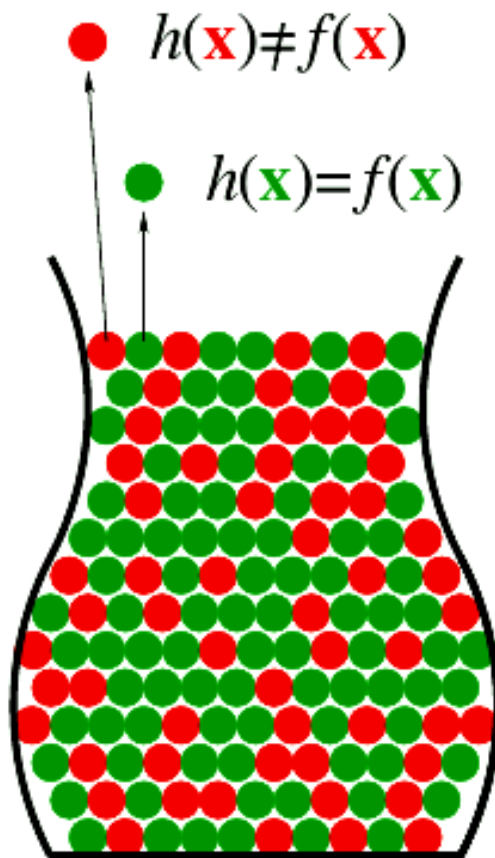


Como resultado:

$$P(|E_{in}(h) - E_{out}(h)| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

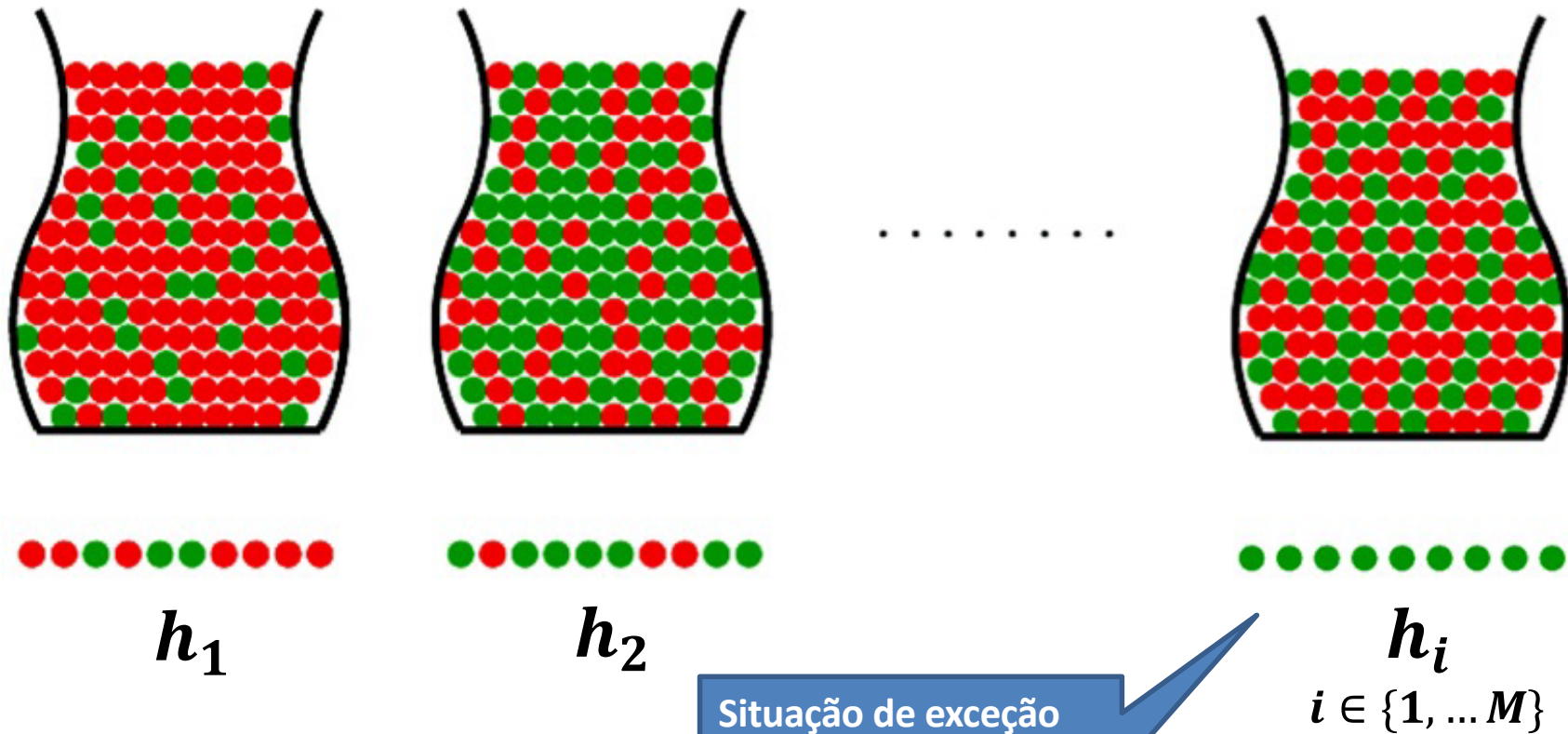
Encontramos um alicerce matemático para a viabilidade da aprendizagem?

# Conexão com aprendizagem de máquina



- A desigualdade de Hoeffding dá suporte matemático para a retirada aleatória de apenas uma amostra.
- Portanto, a hipótese  $h(x)$  deve ser **fixada** a priori, e não aprendida.
- Para essa hipótese especificamente,  $l$  generaliza  $L$ . Contudo, pode-se generalizar um **aprendizado equivocado**.
- A desigualdade de Hoeffding **não funciona** para um algoritmo que trabalha com **várias hipóteses**.

# Hoeffding para múltiplas hipóteses



# Um pouco de probabilidade

**Problema 1:** Se você lançar duas moedas duas vezes, qual a probabilidade de se obter duas caras?

Eventos:

$M_1$ : Obter duas caras com a moeda 1.

$M_2$ : Obter duas caras com a moeda 2.

Espaço amostral:

$A = \{(\textcolor{red}{K}, \textcolor{red}{K}), (K, C), (C, K), (C, C)\}$

$$P(M_1) = P(M_2) = \textcolor{red}{0,25}$$

# Um pouco de probabilidade

Continuação:

Probabilidade da união

$$\begin{aligned}P(M_1 \cup M_2) &= P(M_1) + P(M_2) - P(M_1 \cap M_2) \\&= P(M_1) + P(M_2) - P(M_1) \cdot P(M_2)\end{aligned}$$

$$= \frac{1}{4} + \frac{1}{4} - \frac{1}{16} = \frac{7}{16} \cong \mathbf{0,44}$$

# Um pouco de probabilidade

**Problema 2:** Sabendo-se que a probabilidade de obter-se 10 caras, ao lançar **uma** moeda 10 vezes, é  $\frac{1}{2^{10}} \cong \mathbf{0,1\%}$ , qual a probabilidade do mesmo evento no lançamento de 1000 moedas?

$$P(M_1 \cup M_2 \dots \cup M_{1000}) = \left( \sum_{i=1}^{1000} P(M_i) \right) - \textit{expr}$$

$$P(M_1 \cup M_2 \dots \cup M_{1000}) \leq \sum_{i=1}^{1000} P(M_i)$$

$$\mathbf{R \cong 63\%}$$

# Hoeffding para múltiplas hipóteses

Seja  $g$  a **hipótese final**, selecionada por um algoritmo de aprendizagem qualquer, em um **conjunto de hipóteses**  $H$ , onde  $|H| = M$ . Ao aplicar Hoeffding ao algoritmo de aprendizagem, obtemos:

$$\begin{aligned} P(|E_{in}(g) - E_{out}(g)| > \epsilon) &\leq P( \begin{aligned} &|E_{in}(h_1) - E_{out}(h_1)| > \epsilon \\ &\cup |E_{in}(h_2) - E_{out}(h_2)| > \epsilon \\ &\dots \\ &\cup |E_{in}(g) - E_{out}(g)| > \epsilon \\ &\dots \\ &\cup |E_{in}(h_M) - E_{out}(h_M)| > \epsilon \end{aligned} ) \end{aligned}$$

$$= \sum_{i=1}^M P( |E_{in}(h_i) - E_{out}(h_i)| > \epsilon ) - \text{expr}$$

# Hoeffding para múltiplas hipóteses

$$P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq \sum_{i=1}^M P(|E_{in}(h_i) - E_{out}(h_i)| > \epsilon)$$

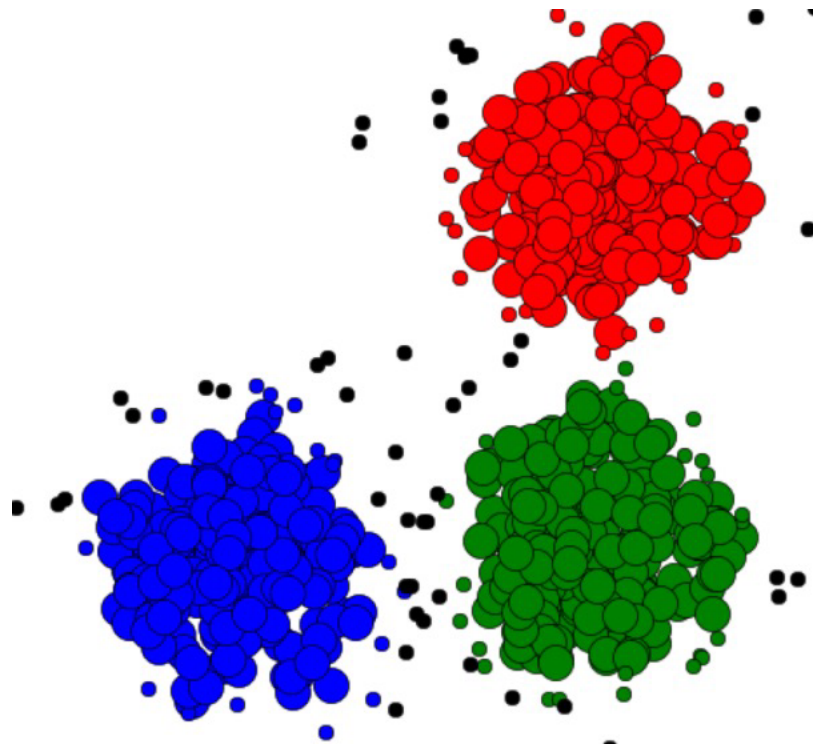
$$\leq \sum_{i=1}^M 2e^{-2\epsilon^2 N}$$

$$P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq 2Me^{-2\epsilon^2 N}$$

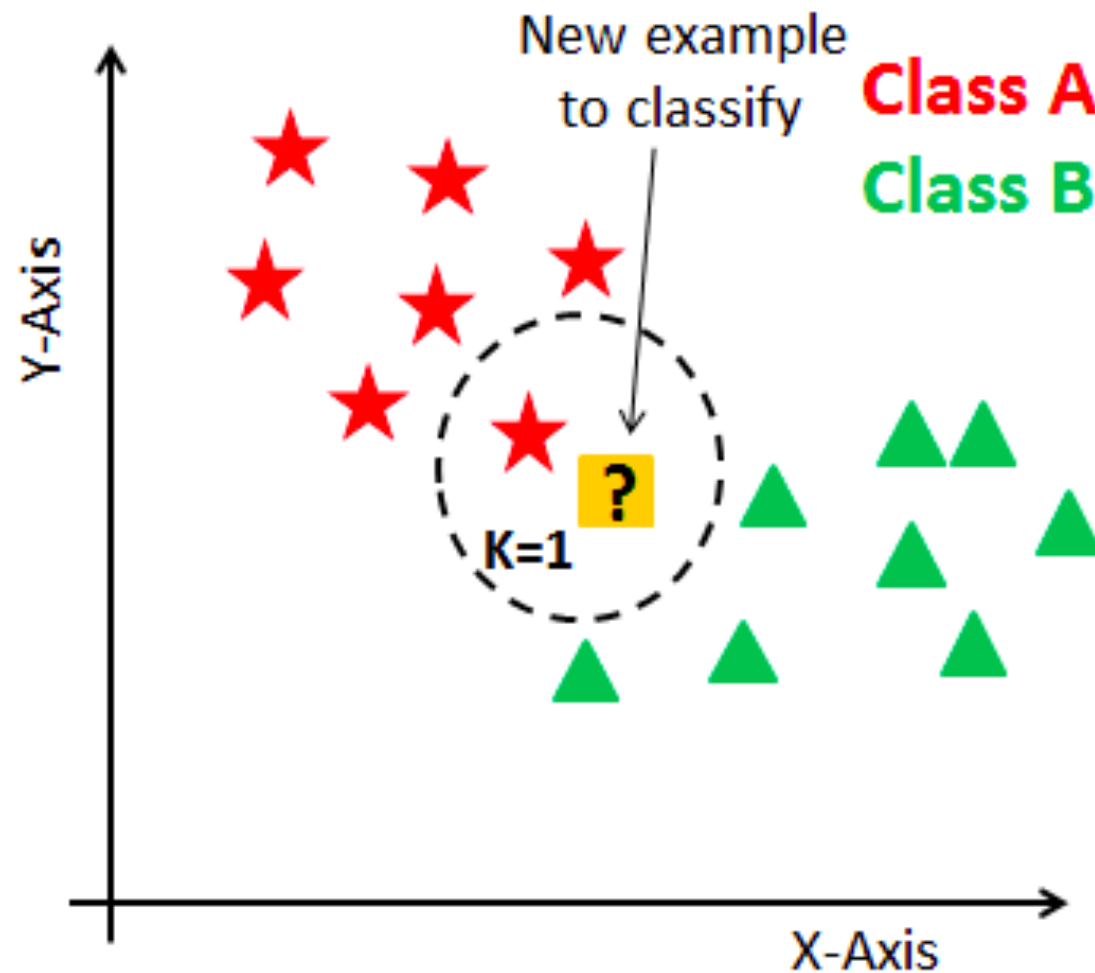


# Modelos baseados em distâncias

## Algoritmo do vizinho mais próximo (*k-Nearest Neighbors*)



# Modelos baseados em distâncias



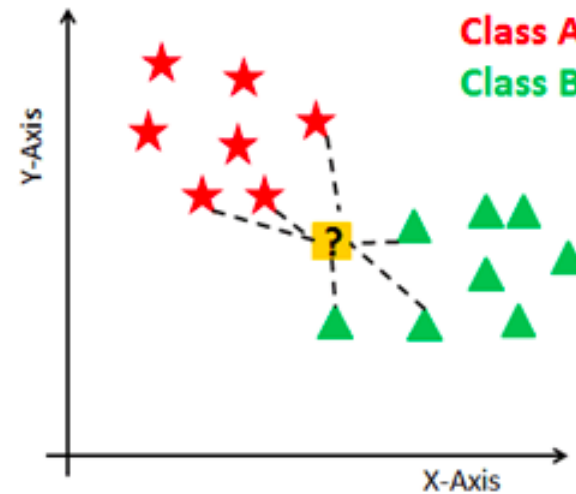
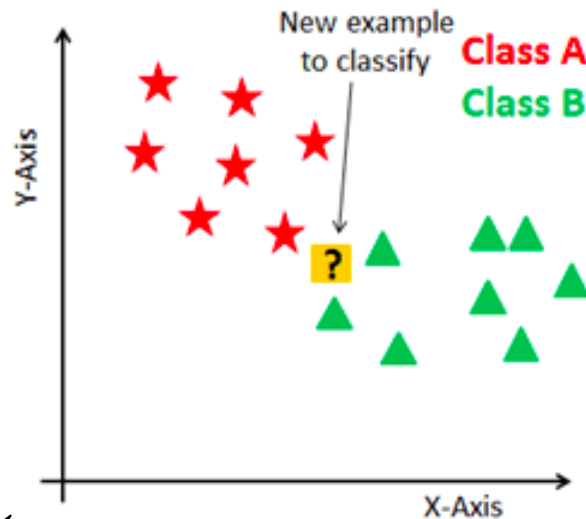
# Modelos baseados em distâncias

- Modelos que consideram a **proximidade de objetos** para a realização de predições.
- Supõe-se que objetos similares tendem a se concentrar em uma mesma região do espaço de entrada.
- **Não há aprendizado** no sentido de indução de uma função hipótese.
  - Há uma memorização do conjunto de treinamento.

# Algoritmo k-NN

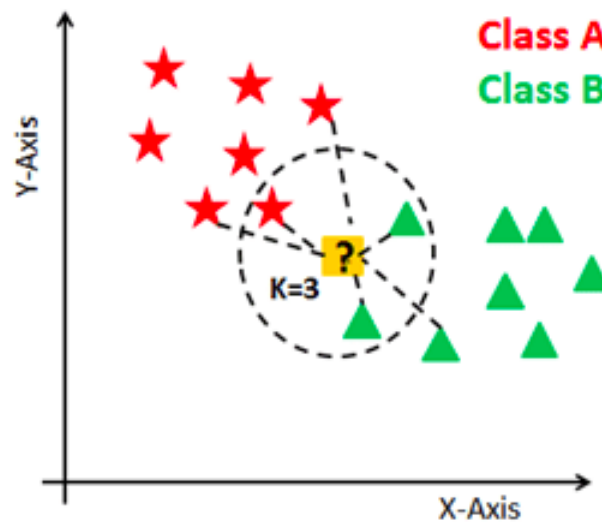
- **Pode ser dividido em três etapas:**
  - Calcular as distâncias;
  - Encontrar os  $K$  vizinhos mais próximos;
  - Eleger o rótulo mais adequado.

# Algoritmo k-NN



Cálculo das distâncias

Dados iniciais



Seleção de vizinhos mais próximos e eleição do rótulo.

# Distâncias

Sejam  $\mathbf{x}_i$  e  $\mathbf{x}_j$  dois vetores em  $\mathbb{R}^d$ , as distâncias entre eles são dadas por:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{l=1}^d (x_i^l - x_j^l)^2}$$

**Euclidiana**

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{l=1}^d |x_i^l - x_j^l|$$

**Manhattan**

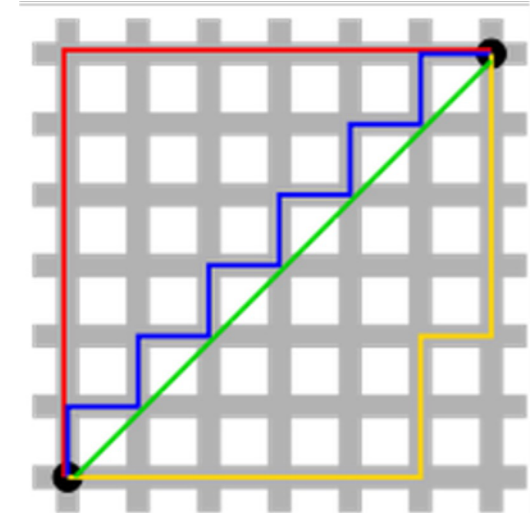
$$d(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_{l=1}^d |x_i^l - x_j^l|^p \right)^{1/p}$$

**Minkowski**

# Distâncias

- Generalização das distâncias:

$$d(x_i, x_j) = \left( \sum_{l=1}^d |x_i^l - x_j^l|^p \right)^{1/p}$$



- Para  $p = 1$ , temos a distância de Manhattan.
- Para  $p = 2$ , a distância euclidiana.
- Para  $p = \infty$ , a dimensão dominante é destacada.

# Distâncias

- **Considerações:**

- A distância euclidiana é mais **sensível** a pequenas modificações do que a distância de Manhattan.
  - Não é adequada para dimensões com muito ruído.
- Um valor de  **$p$  alto** pode dar muita ênfase a dimensões com outliers.



# Algoritmo k-NN

- **Considerações**

- A escolha do **k** não é trivial.
  - Normalmente um valor pequeno e ímpar.
  - A integração com algoritmos evolutivos é uma alternativa.
- Considerado um algoritmo preguiçoso.
  - Maior parte da computação ocorre no momento da classificação.

# Algoritmo k-NN

- **Vantagens**

- Algoritmo intuitivo e simples de implementar.
- Consiste em um algoritmo incremental.
  - Novos exemplos são adicionados sem gerar esforço computacional adicional.
- Poucos parâmetros a serem ajustados.

# Algoritmo k-NN

- **Desvantagens**

- Necessário recalcular as distâncias para cada novo ponto a ser rotulado.
- Susceptível a atributos redundantes e irrelevantes.
- Com o aumento no número de dimensões, há um salto na magnitude das distâncias.
  - A distância do vizinho mais próximo aproxima-se da do mais afastado.
- Necessidade de normalização dos valores das dimensões.

# Aplicações

- Reconhecimento facial.
- Identificação de padrões de fraude na utilização de cartões de crédito.
- Identificação de padrões de compra em lojas varejistas.
- Sistemas de recomendação.
- Benchmark para modelos mais sofisticados.

# Referências bibliográficas

- Abu-Moustafa, Y.S.; Magdon-Ismael, M.; Lin, H-S.  
*“Learning from data”*. AMLBook, 2012.
- Faceli, K.; Lorena, A.C.; Gama, J.; Carvalho, A.C.P.L.F.  
*“Inteligência Artificial Uma Abordagem de Aprendizado de Máquina”*. LTC, 2011.
- Notas de aula do prof. Abu-Moustafa.