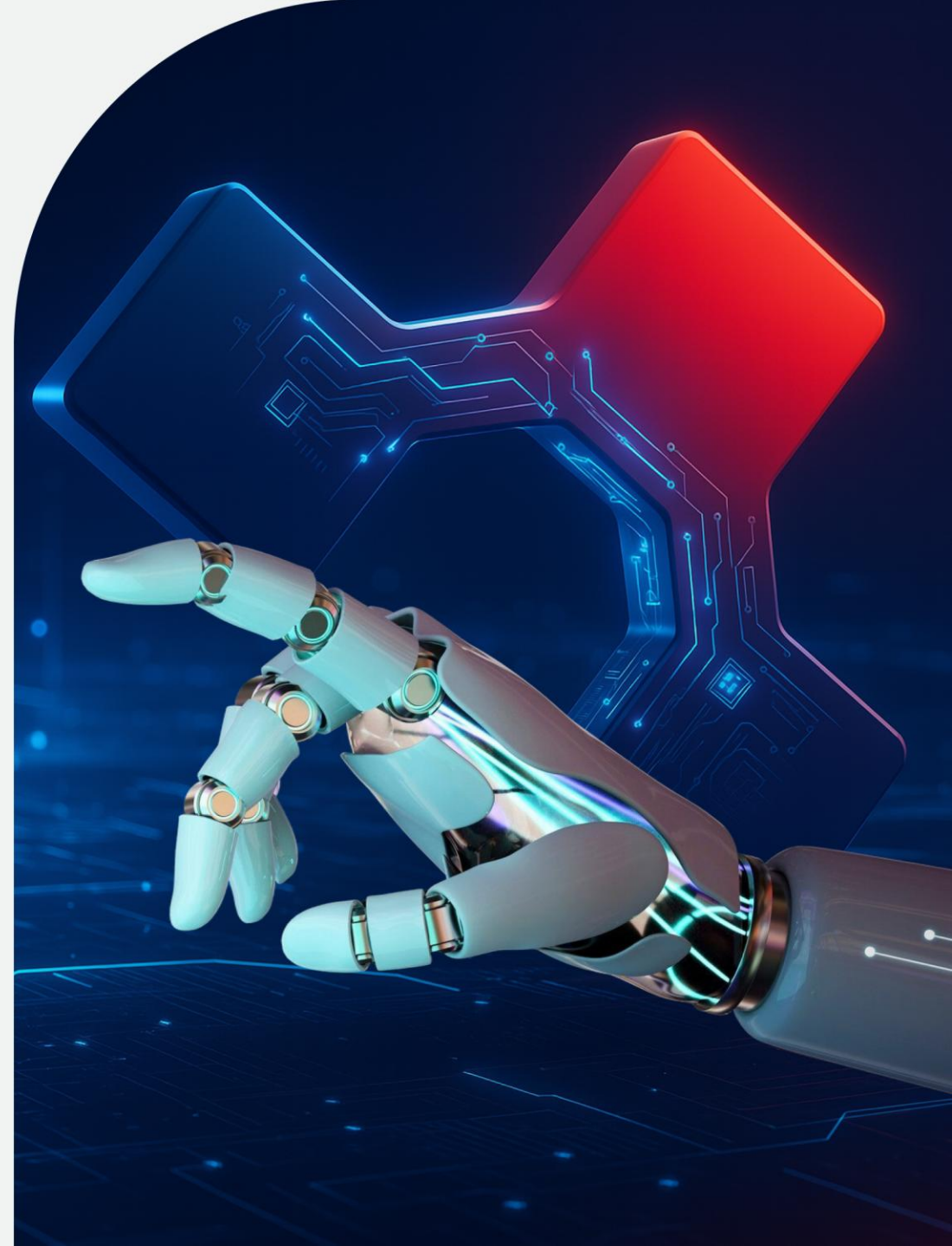




Núcleo de Capacitação em Inteligência Artificial





# Main Challenges of Machine Learning

Insufficient Quantity of Training Data,  
Nonrepresentative Training Data, Poor-Quality Data,  
Irrelevant Features, Overfitting the Training Data,  
Underfitting the Training Data

# Algoritmos Ruins



```
divs[1];  
var atpos=inputs[i].indexOf("@");  
var dotpos=inputs[i].lastIndexOf(".");  
if (atpos<1 || dotpos<atpos+2 || dotpos==inputs[i].length-1 ||  
document.getElementById("errEmail").innerHTML !=  
else  
document.getElementById(div).innerHTML +=  
inputs[i].lastIndexOf(".")
```



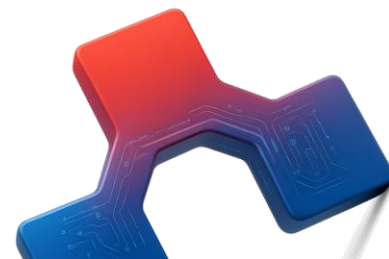
NCIA



FOXCONN



FPFtech



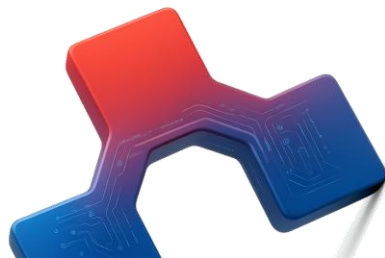
# Algoritmos Ruins



```
var atpos=inputs[i].indexOf("e");  
var dotpos=inputs[i].lastIndexOf(".");  
if (atpos<1 || dotpos<atpos+2 || dotpos>inputs[i].length-1)  
document.getElementById("errmsg").innerHTML+=  
else  
document.getElementById(div).innerHTML+=
```

```
var atpos=inputs[i].indexOf("@");
var dotpos=inputs[i].lastIndexOf(".");
if (atpos<1 || dotpos<atpos+2 || dotpos>inputs[i].length-1)
document.getElementById("errEmail").innerHTML += "Invalid email address  
";
else
document.getElementById(div).innerHTML += "Valid email address  
";
```

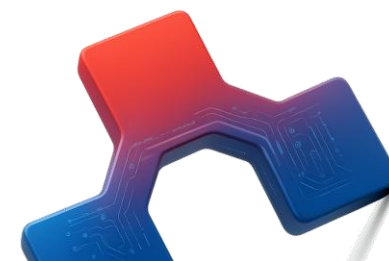
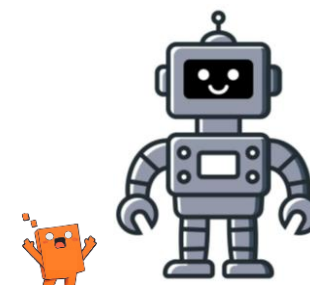
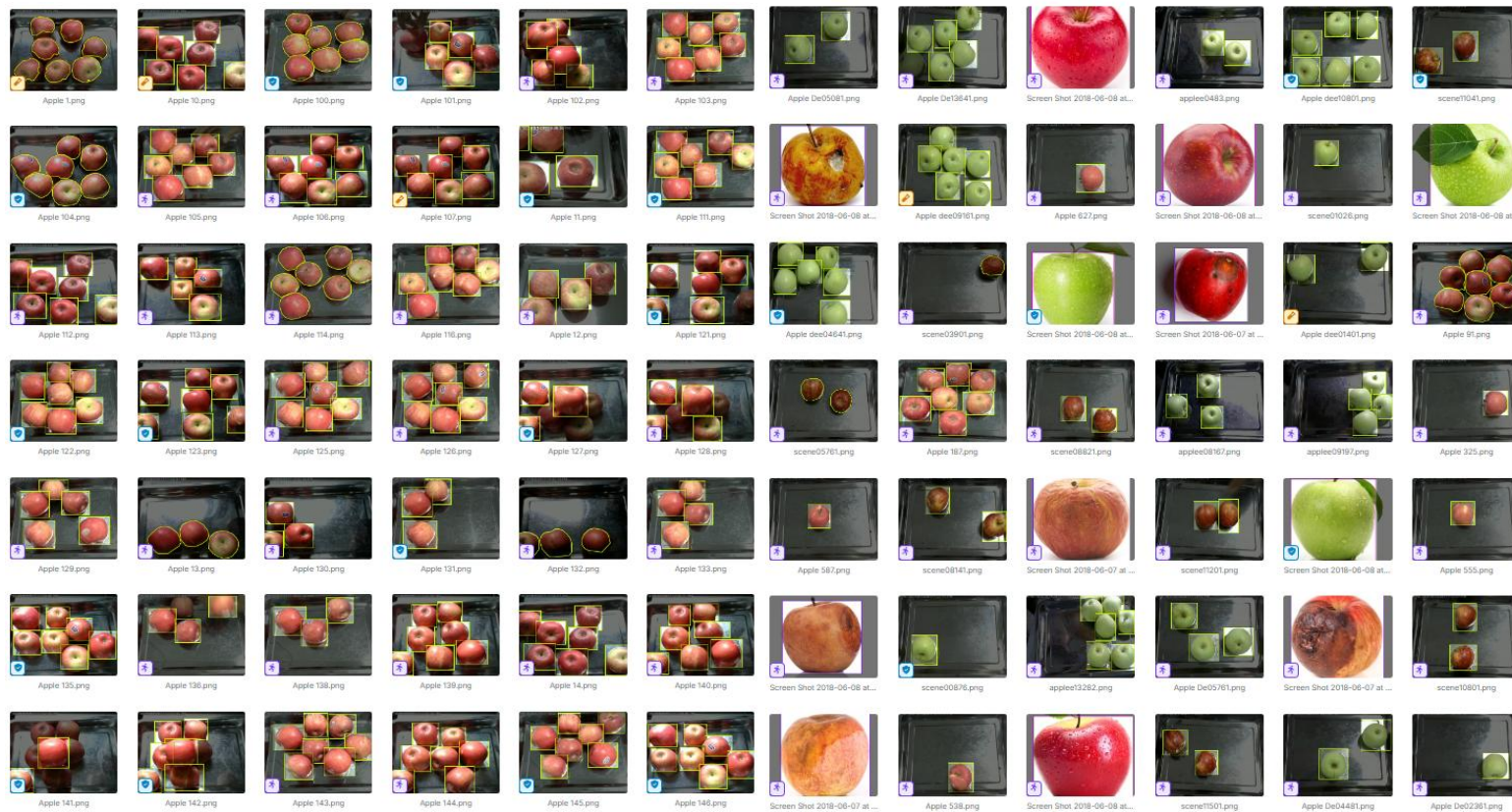
# Dados Ruins

Three blue dice are arranged in a row. Each die has an orange face with a white 'X' shape, which is composed of four small orange pips (ones) and one larger orange pip (four), totaling a value of five. The dice are set against a white background with a thick orange border.

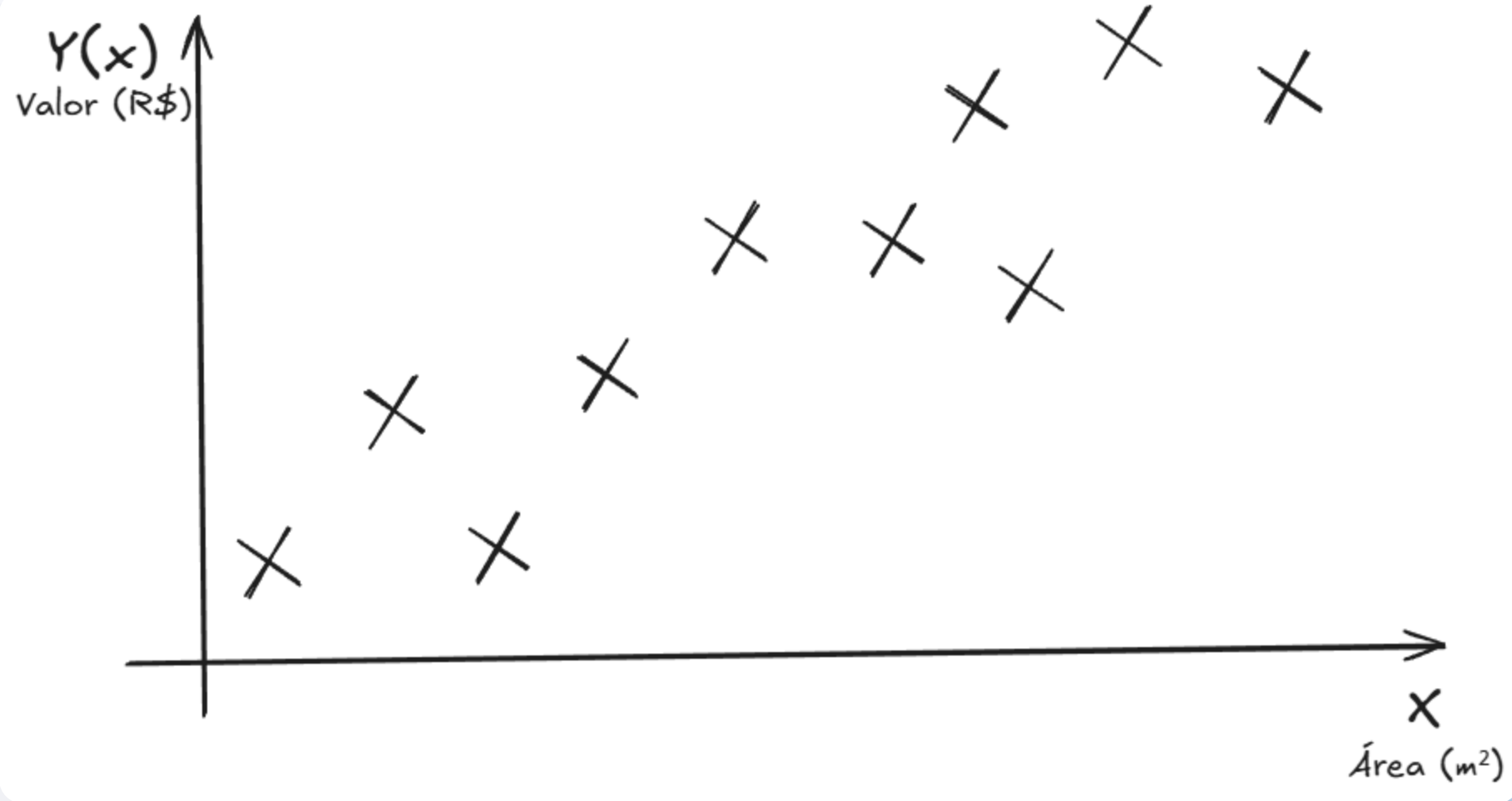


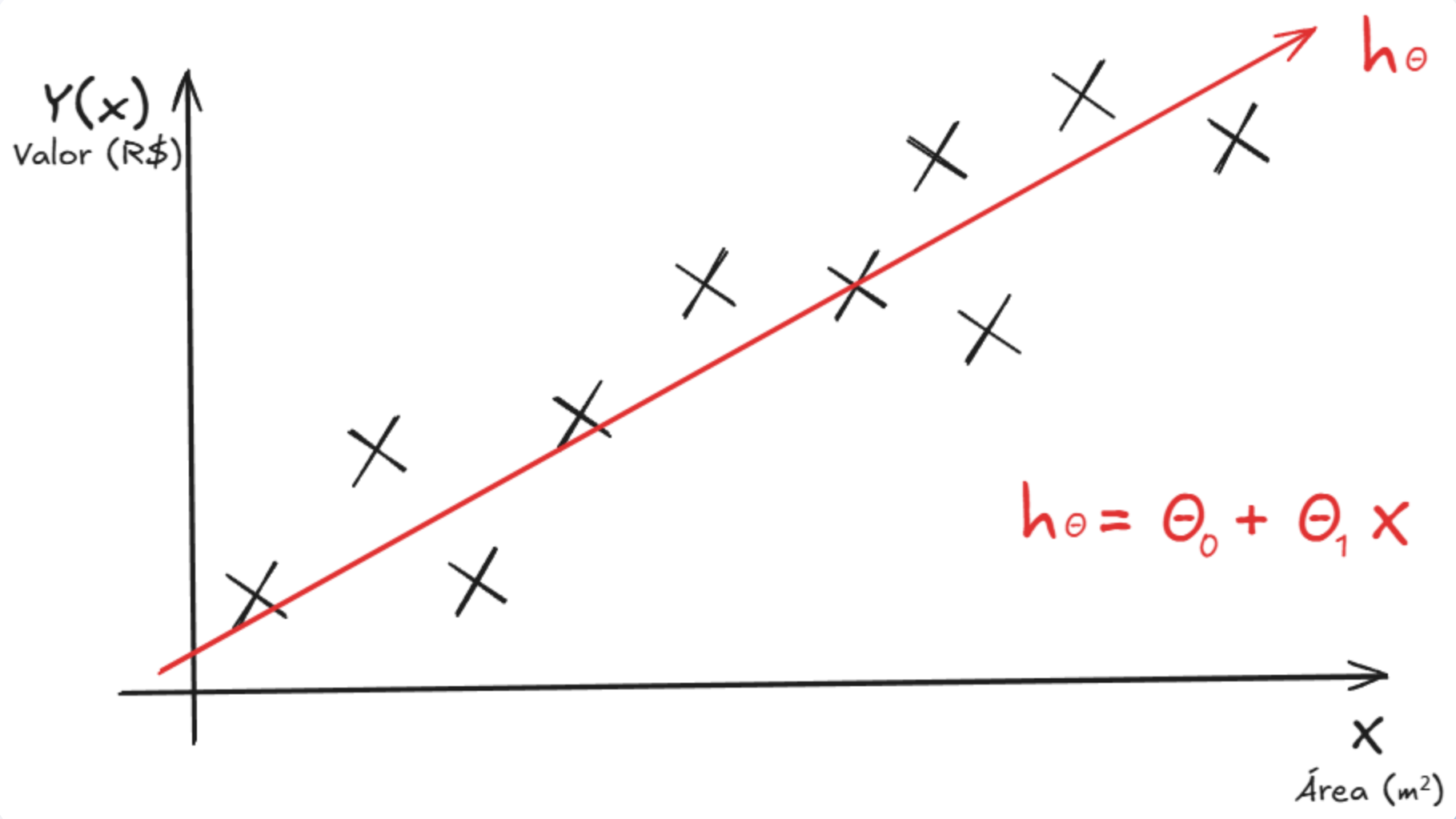


# Insufficient Quantity of Training Data

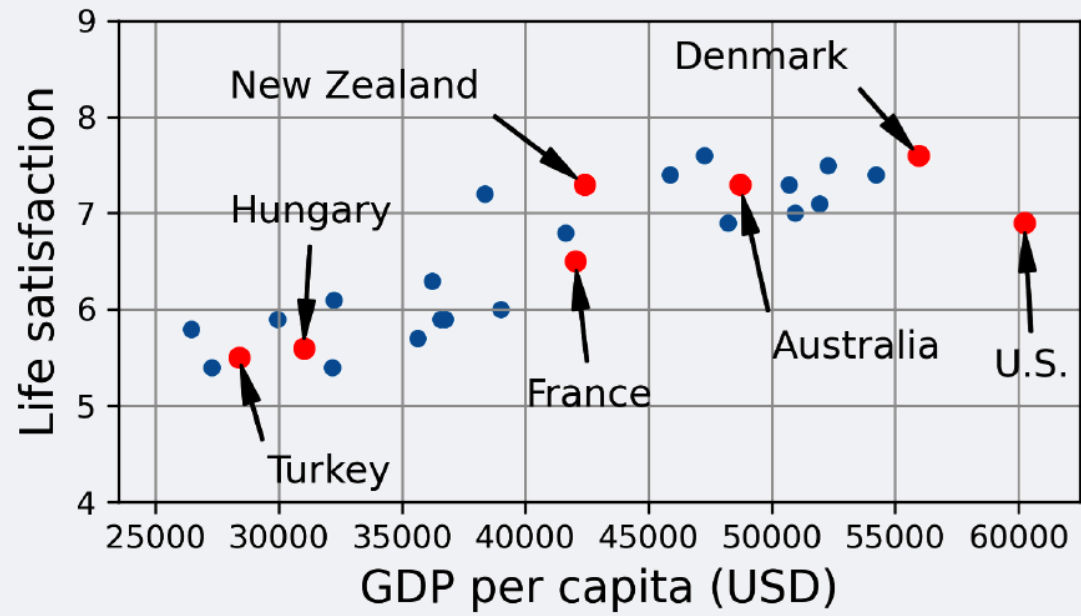


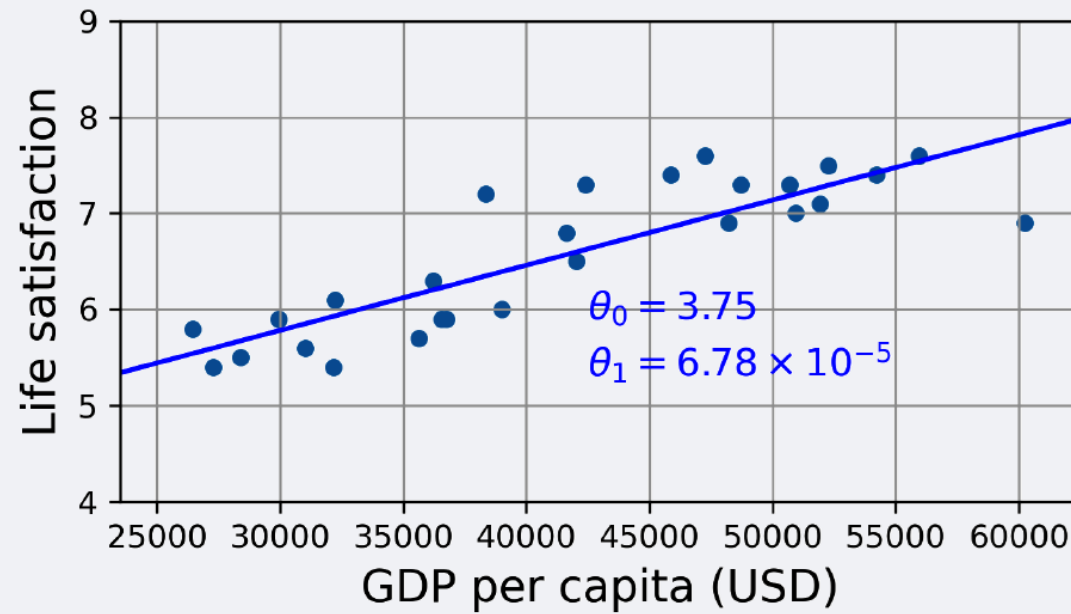
# Nonrepresentative Training Data

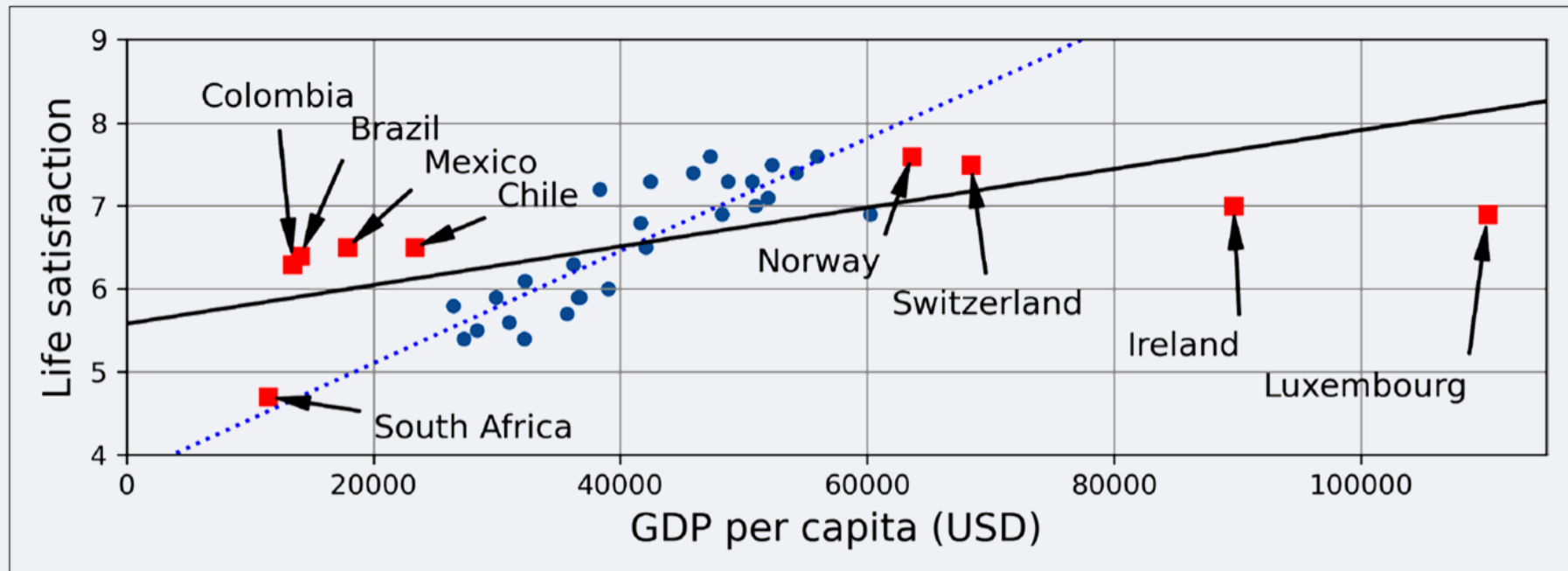








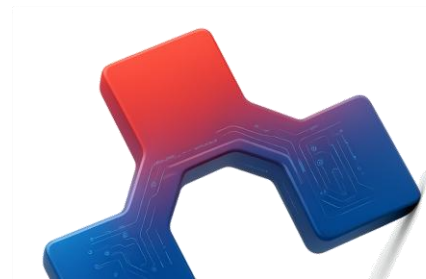






**Generalização** em aprendizado de máquina é a capacidade essencial de um modelo de **performar bem** em **dados novos** e nunca vistos, após ter sido treinado em um conjunto de dados limitado.

Um modelo que generaliza bem não apenas "decorou" os exemplos de treinamento, mas **aprendeu os padrões** verdadeiros e subjacentes, permitindo-lhe fazer **previsões** precisas e úteis em situações do **mundo real**.



Para que um modelo de Machine Learning generalize bem para novos dados, é crucial que seu conjunto de treinamento seja **representativo** para aquele **universo de dados**.

Um modelo que previa a **satisfação** com a vida a partir do **PIB**: ao omitir países **muito ricos** e **muito pobres**, apresentou uma forma linear aparentemente satisfatória mas quando submetido às novas amostras mostrou-se inadequado.

A falta de representatividade pode ocorrer por **ruído de amostragem** (amostras pequenas) ou por **viés de amostragem** (método de coleta falho), que é um problema sério mesmo em amostras grandes.





# Poor-Quality Data

Problema: Dados com erros, outliers, ruído, ausentes...

Solução: Pré-processamento de dados



# Irrelevant Features

Um modelo só consegue aprender se os dados de treino contiverem **atributos relevantes**.

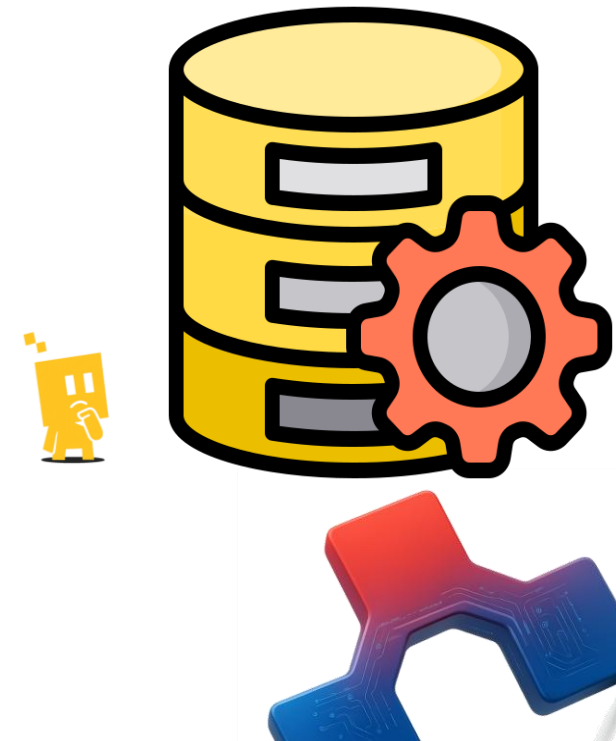
**Engenharia de Atributos:** É o processo de selecionar e criar o melhor conjunto de atributos (features) para o modelo, sendo uma parte crítica para o sucesso de um projeto de Machine Learning.

Etapas Principais:

**Seleção de Atributos:** Escolher os atributos mais úteis entre os que já existem.

**Extração de Atributos:** Combinar atributos existentes para criar um novo mais informativo.

**Criação de Atributos:** Coletar novos dados para gerar atributos que ainda não existem.



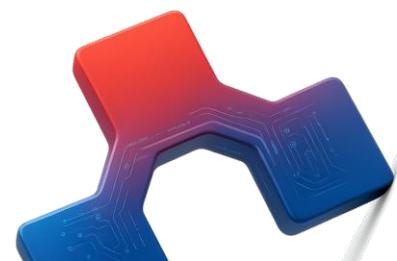


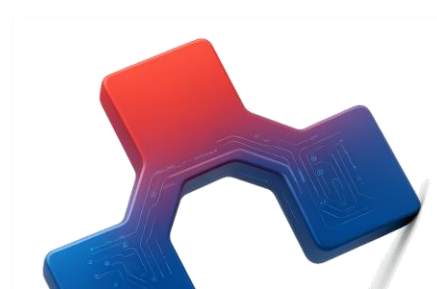
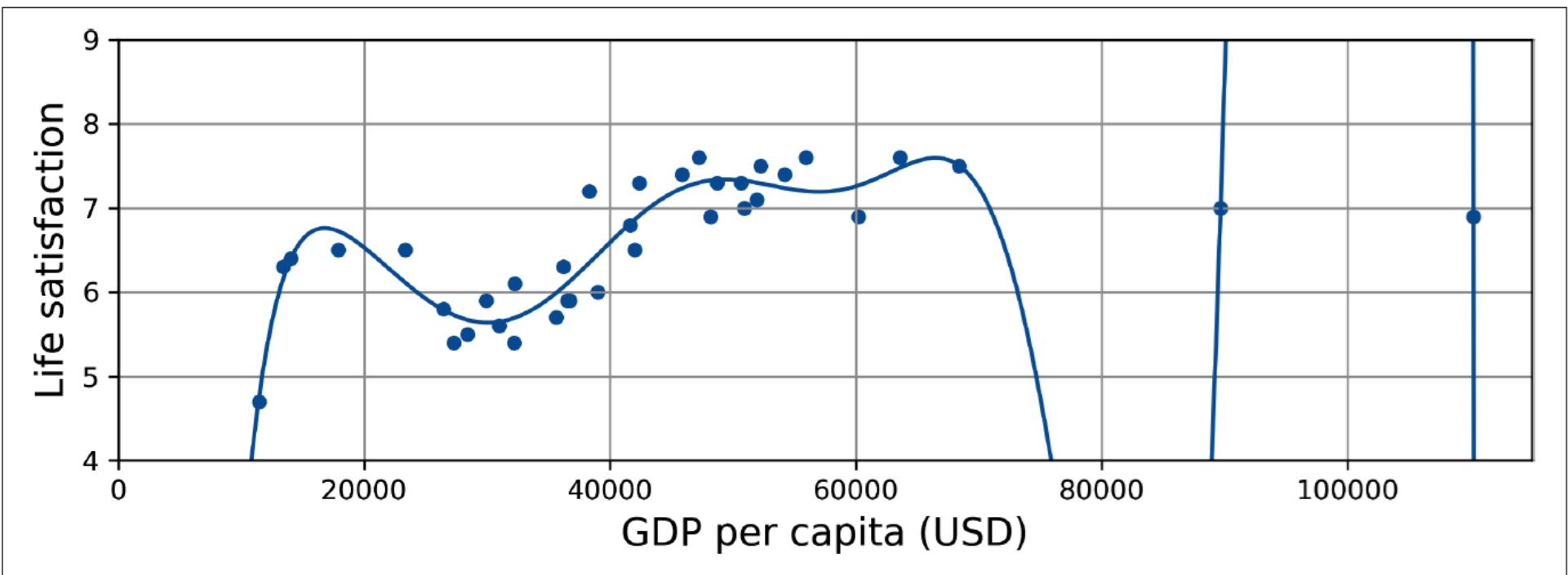
# Overfitting the Training Data

**Overfitting** (Sobreaajuste): O Risco de Aprender Demais

Ocorre quando o modelo se ajusta perfeitamente aos dados de treino, mas falha ao fazer previsões com novos dados. Ele **memoriza o ruído** em vez de aprender o padrão real.

Exemplo Prático: O modelo "aprende" que todos os países com a letra "W" no nome têm alta satisfação com a vida (ex: New Zealand, Sweden). Obviamente, essa é uma regra falsa que não se aplicará a outros casos.





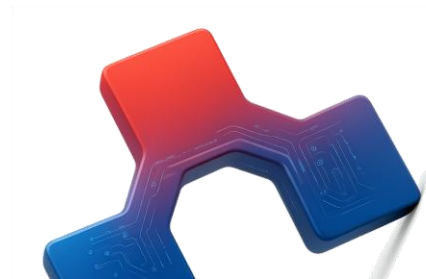


# Underfitting - Subajuste dos dados de treinamento

**Underfitting** (Subajuste): Quando o Modelo é Simples Demais

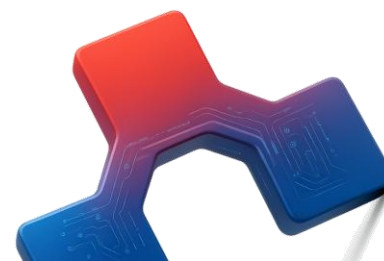
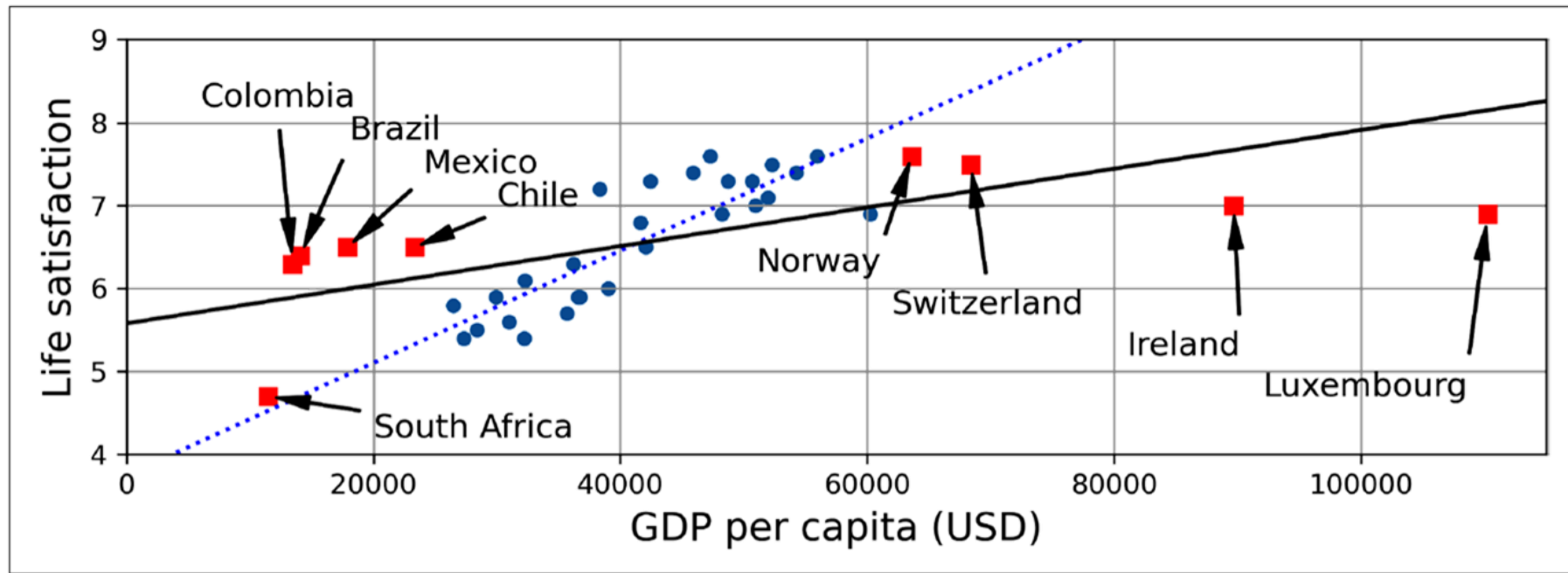
É o oposto de overfitting. Ocorre quando o seu modelo é **muito simples** para aprender a estrutura e os padrões reais presentes nos dados.

Exemplo Prático: Usar um modelo de linha reta (linear) para descrever algo complexo como a satisfação com a vida. A realidade não é uma linha reta.

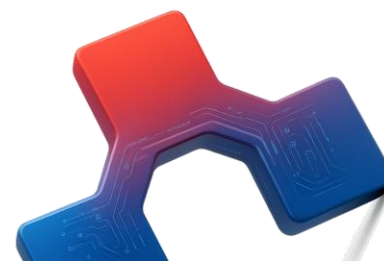
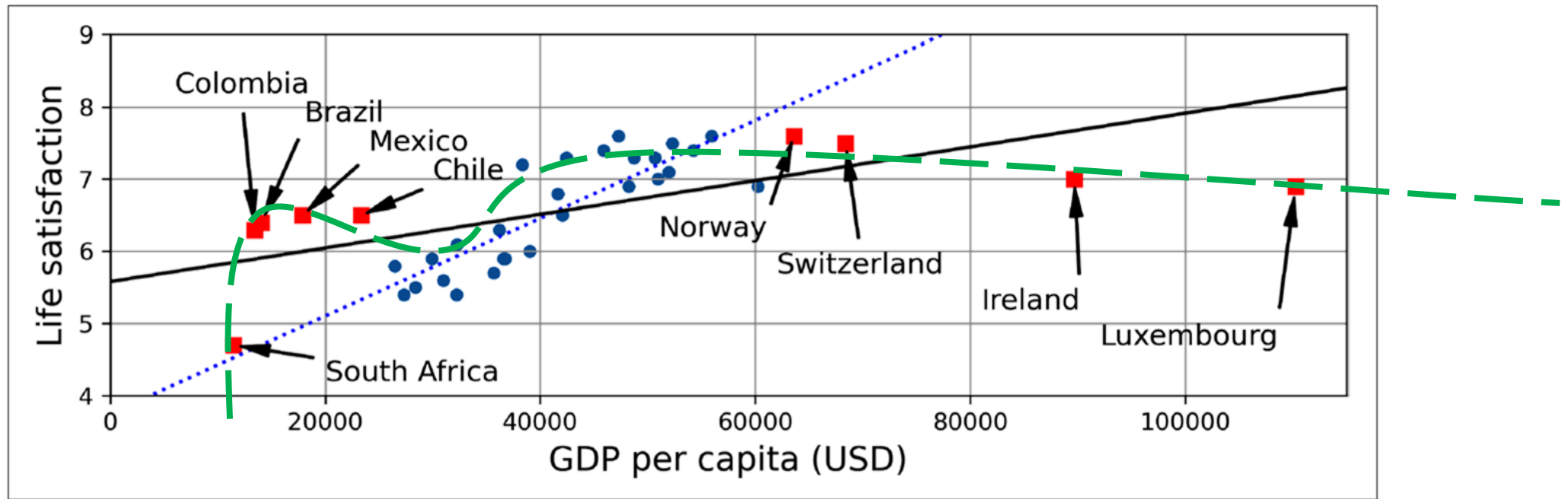


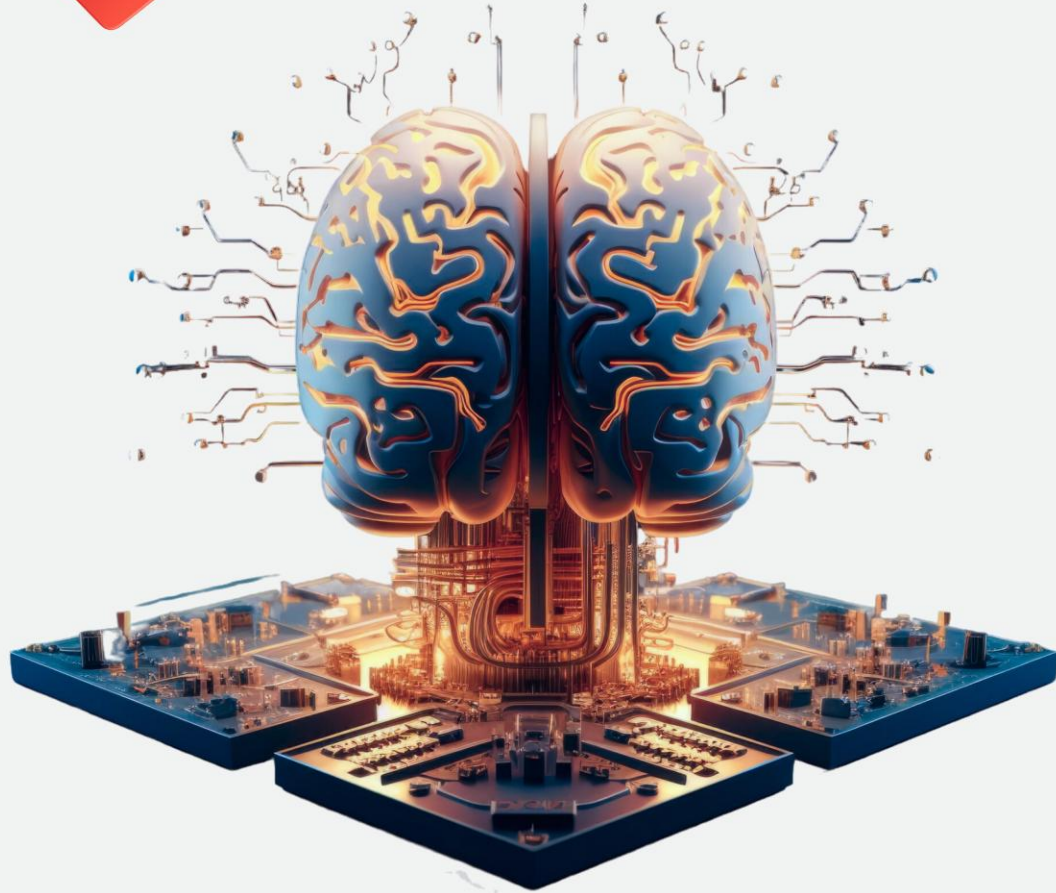


Regressão Linear **não** seria uma abordagem suficiente para **esse** universo de dados



Regressão Linear **não** seria uma abordagem suficiente para **esse** universo de dados





# Testing and Validating

Hyperparameter Tuning and Model Selection, Data Mismatch



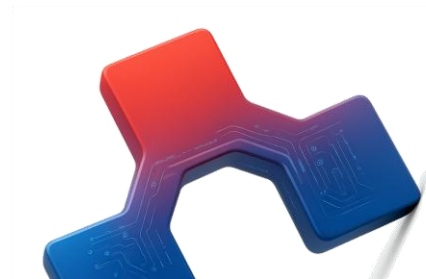
# Testando e validando

Para saber se um modelo é bom sem arriscar em produção, divida seus dados em **conjunto de treino** e **conjunto de teste**.

**Treine** o modelo com o conjunto de treino; **avale** o desempenho com o conjunto de teste.

O erro no teste (erro de generalização) mostra como o modelo se sairá com dados que nunca viu.

Se o modelo vai bem no treino, mas mal no teste, ele está com **overfitting**.





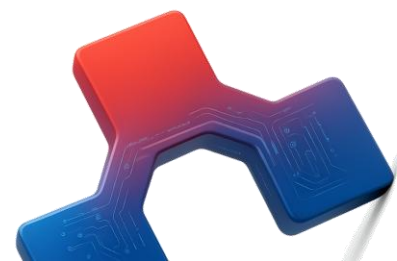
Data Base



NCIA



FOXCONN®

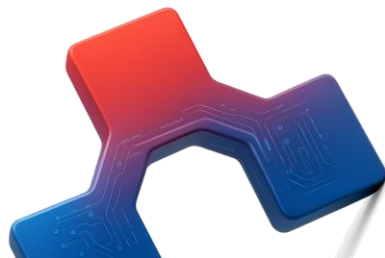






Data Base

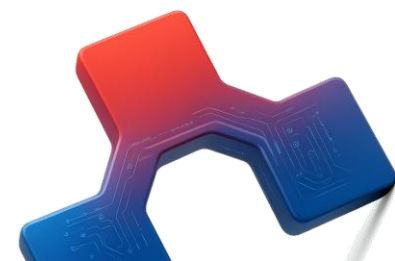
class\_index x\_center  
y\_center width height





Data Base

class\_index x\_center  
y\_center width height





Data Base



NCIA



FOXCONN®

