



<https://developers.google.com/machine-learning/crash-course/training-and-test-sets/splitting-data>

# Test Set

**O volume de dados é suficientemente grande para trazer resultados estatisticamente significativos**

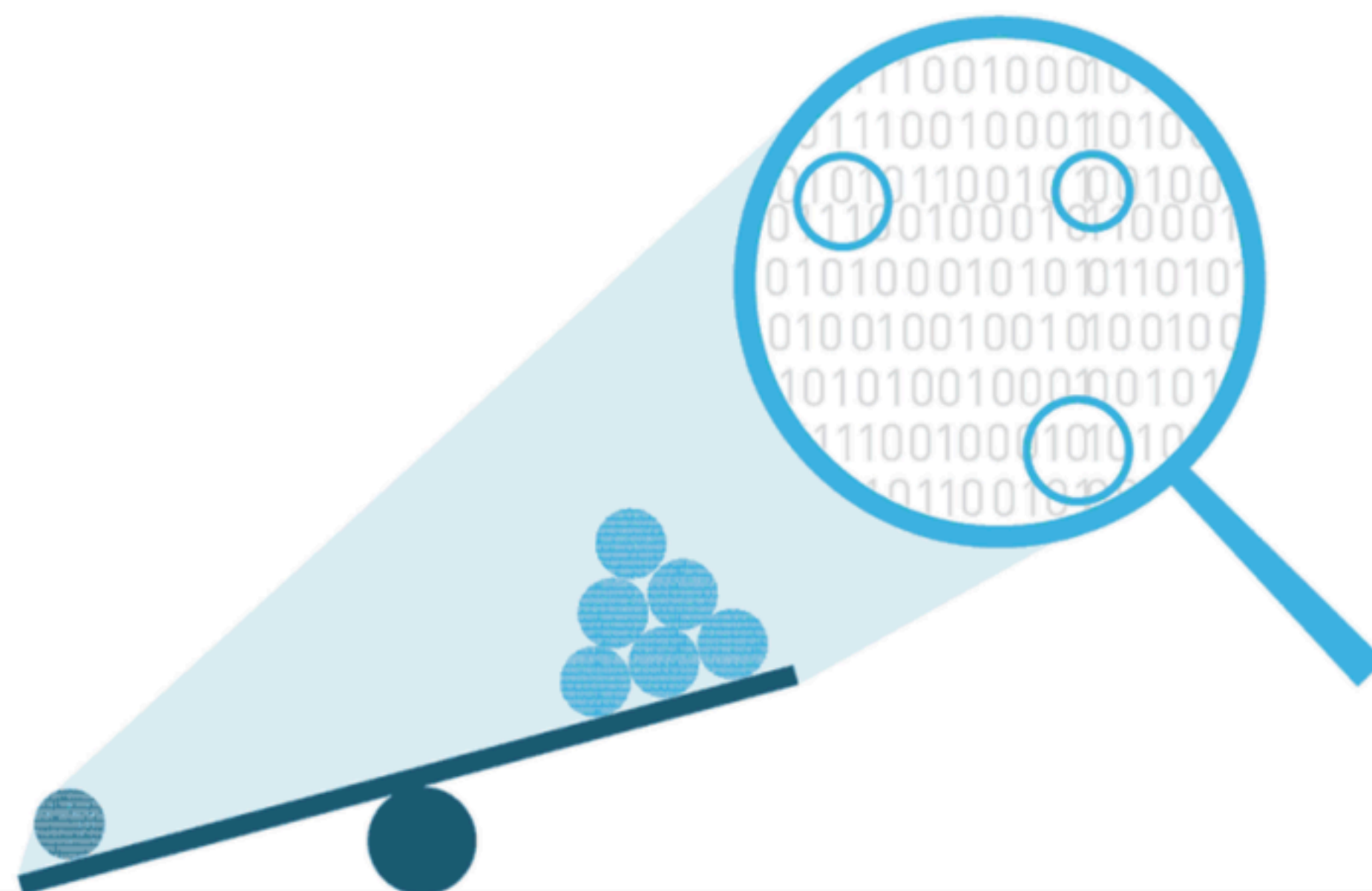
**Representa o dataset como um todo (desbalanceamento)**

BLOG

# Como lidar com dados desbalanceados?

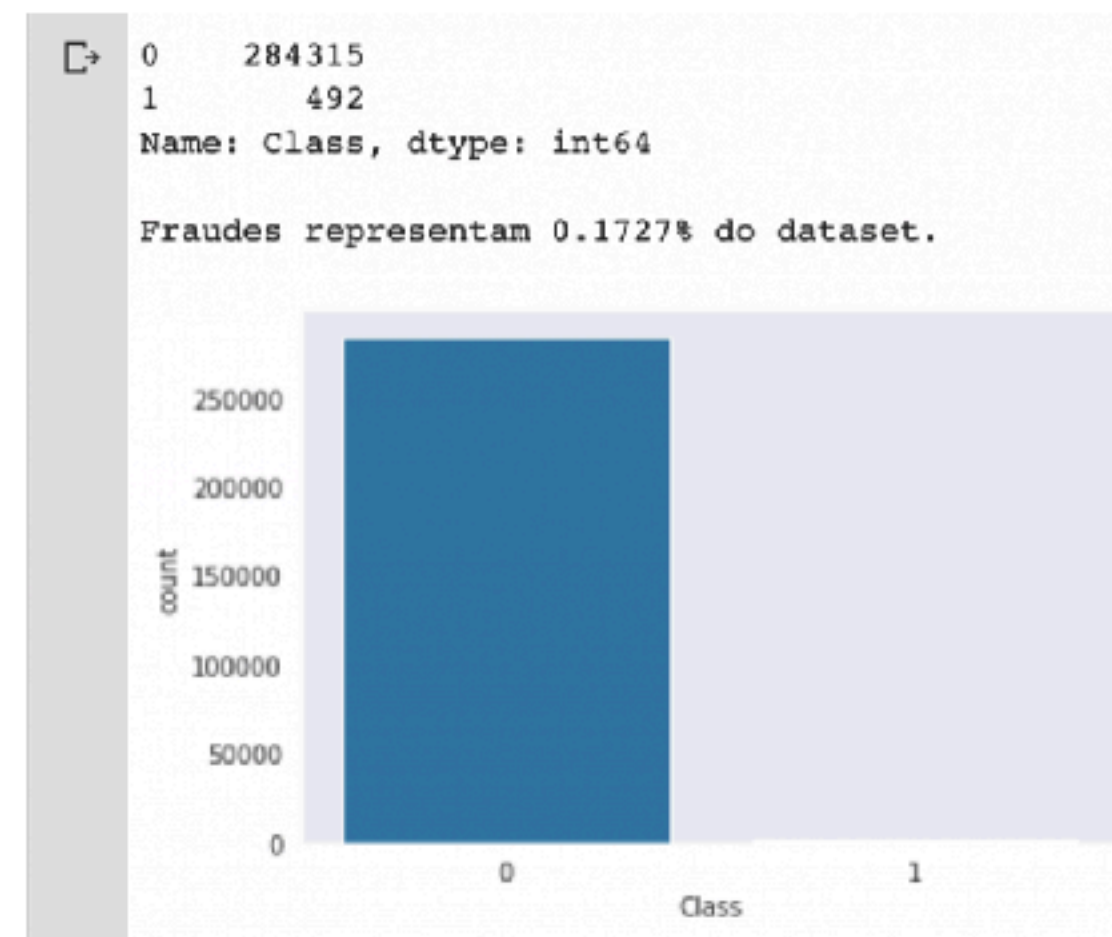


Escrito por Carlos Melo  
em 24/12/2019



dados. Mas conseguimos ver a barra da nossa variável `df['Class']` para instâncias de fraude (`df['Class'] == 1`).

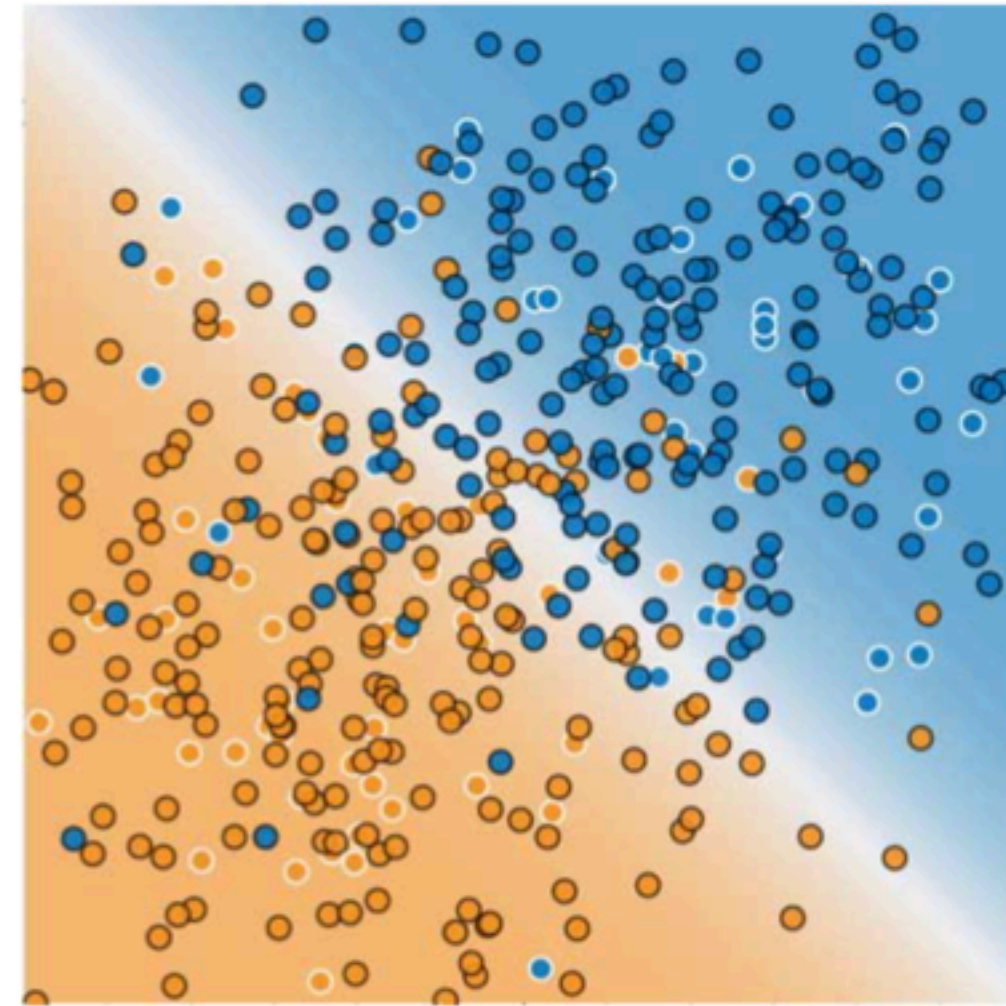
```
1 file_path = "https://www.dropbox.com/s/b44o3t3ehmnx2b7/creditcard.csv?dl=1"
2
3 # importar os dados para um dataframe
4 df = pd.read_csv(file_path)
5
6 # ver o balanceamento das classes
7 print(df.Class.value_counts())
8 print("\nFraudes representam {:.4f}% do dataset.\n".format((df[df.Class == 1].shape[0] / df.shape[0]) *
9
10 # plotar gráfico de barras para as Classes
11 sns.countplot('Class', data=df);
```



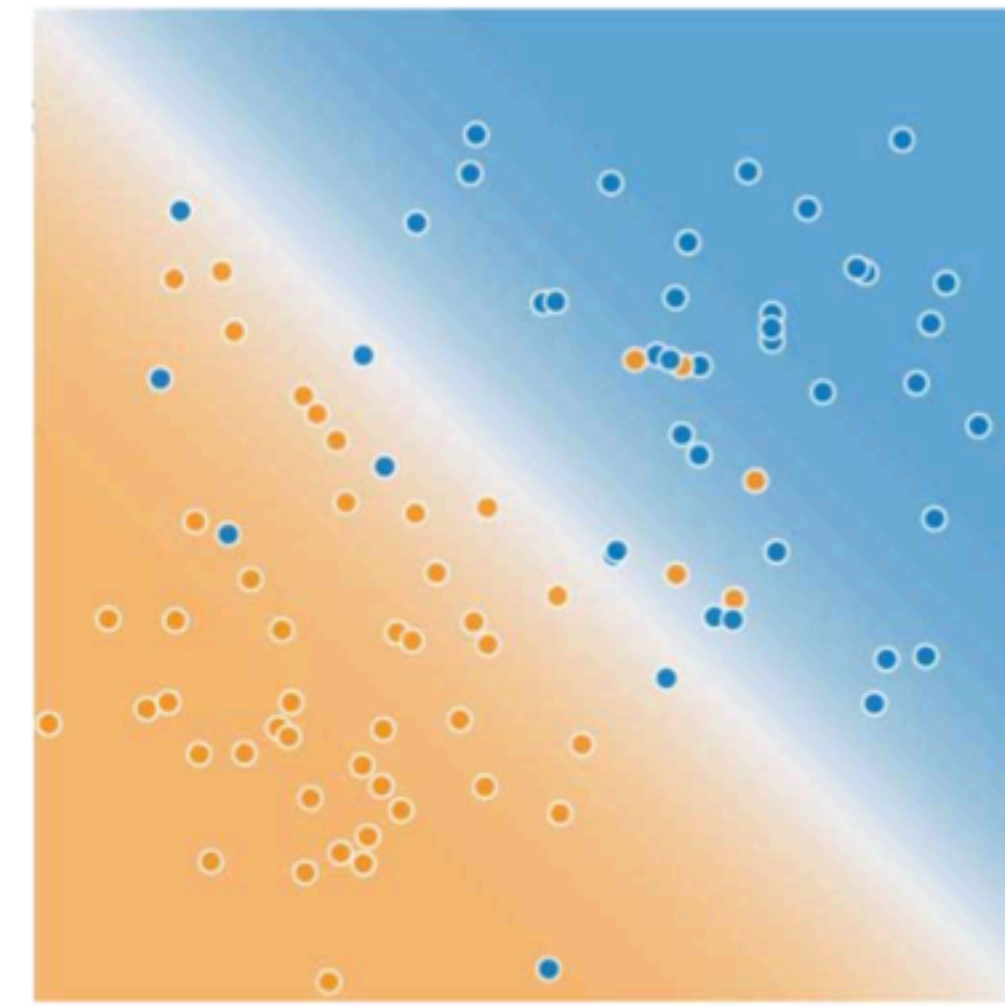
Para você ter um ideia real das consequências dessa situação, vou construir dois modelos de Regressão Logística. No nosso primeiro modelo, separei as variáveis `X` e `y` normalmente e dividi entre conjuntos de treino e teste, como é praxe em *machine learning*.

Sem maiores ajustes, treinei o modelo usando o método `fit(X_train, y_train)` e fiz a previsão de valores em cima do conjunto de teste (`X_test`). Na sequência, plotei a matriz de confusão e o relatório de classificação.

```
1 # separar variáveis entre X e y
2 X = df.drop('Class', axis=1)
3 y = df['Class']
```



Training Data

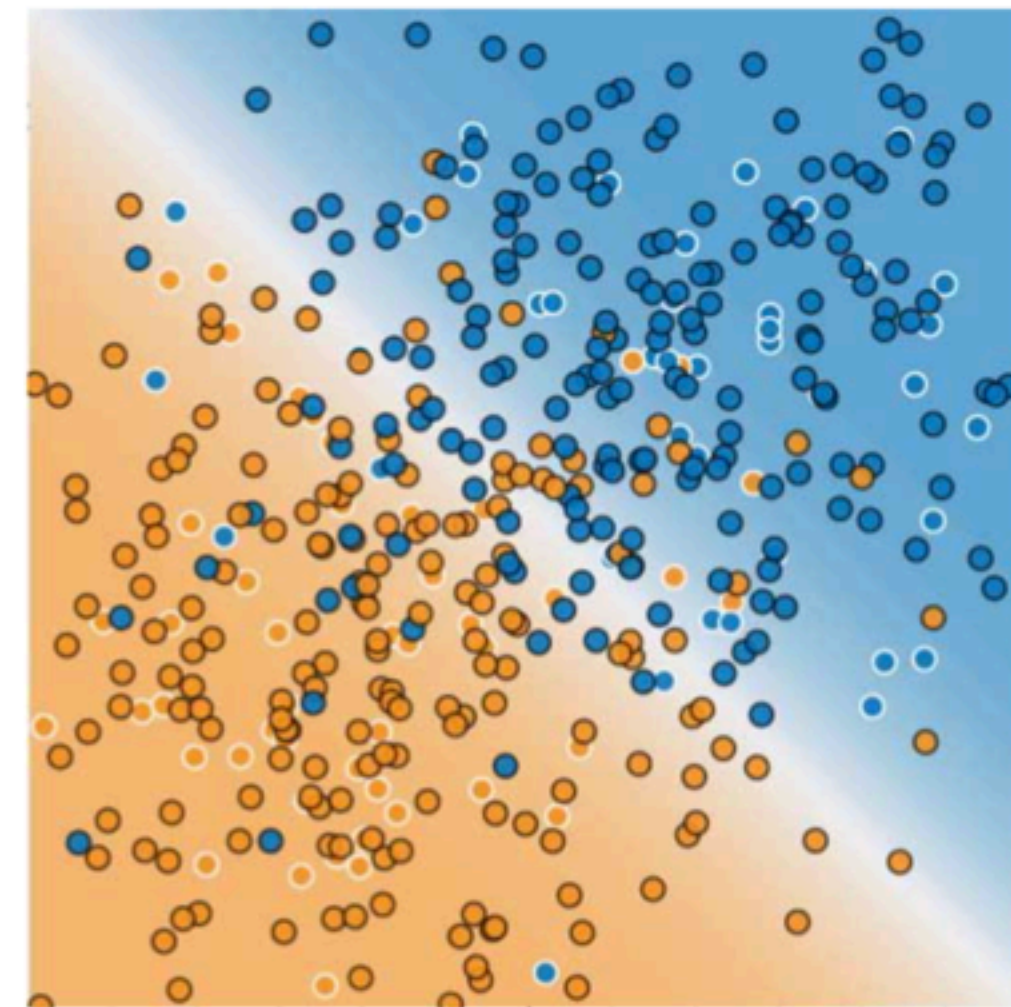


Test Data

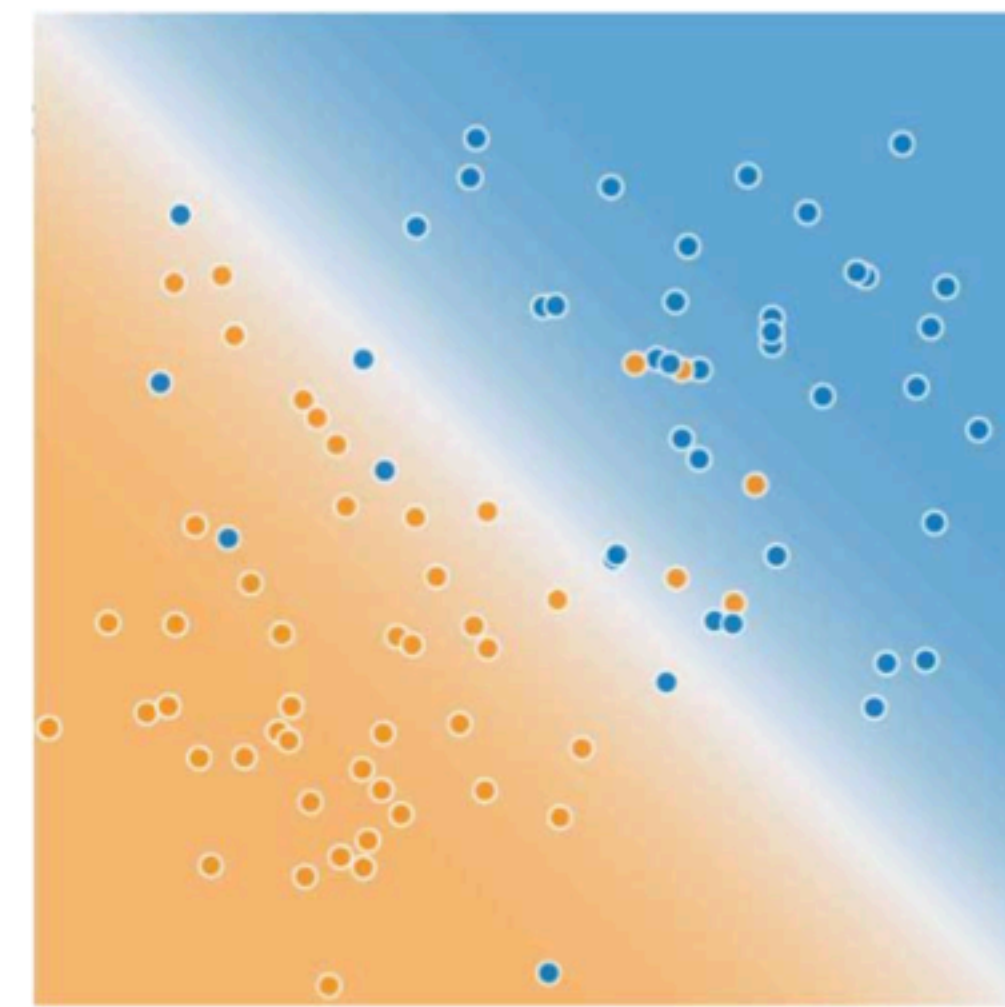
<https://developers.google.com/machine-learning/crash-course/training-and-test-sets/splitting-data>



Um bom modelo tem  
desempenho bom  
nos datasets de  
**treino e teste**



Training Data



Test Data

<https://developers.google.com/machine-learning/crash-course/training-and-test-sets/splitting-data>

**NUNCA** treine  
seu modelo nos  
dados de  
treino.

**Acurácia alta é  
um indicativo  
de **overfitting**.  
Verifique seus  
dados de teste.**