

Métricas de Classificação e Validação Cruzada

Accuracy, Matriz de Confusão, F1-Score e Cross-Validation (com exemplos no Colab)

Objetivos da aula

O que você deve conseguir fazer ao final 🤖

- Entender o que é Accuracy e quando ela pode enganar
- Ler uma Matriz de Confusão e identificar TP, TN, FP e FN
- Calcular e interpretar Precision, Recall e F1-Score
- Aplicar Validação Cruzada (k-fold) para estimar generalização
- Gerar gráficos no Colab para comunicar resultados



Regra de ouro: uma única métrica raramente conta a história inteira.

Exercício Prático

Classifique as pessoas exibidas a seguir em uma das classes: “feliz” ou “triste”

A



B



C



D



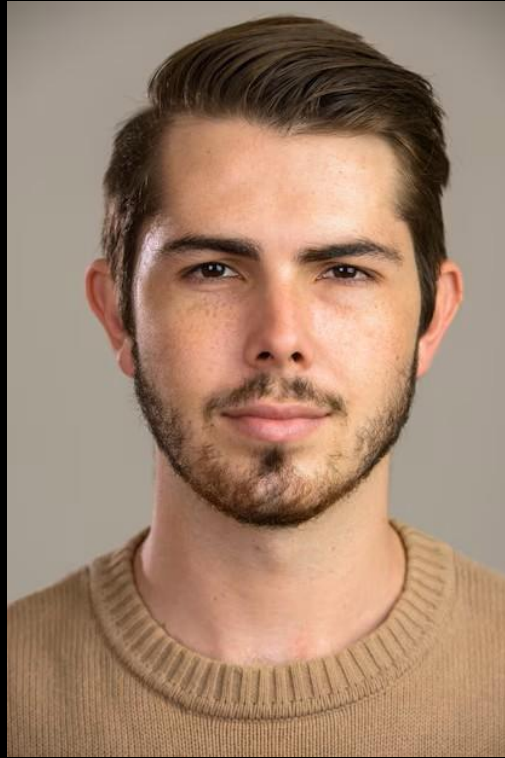
E



F



G

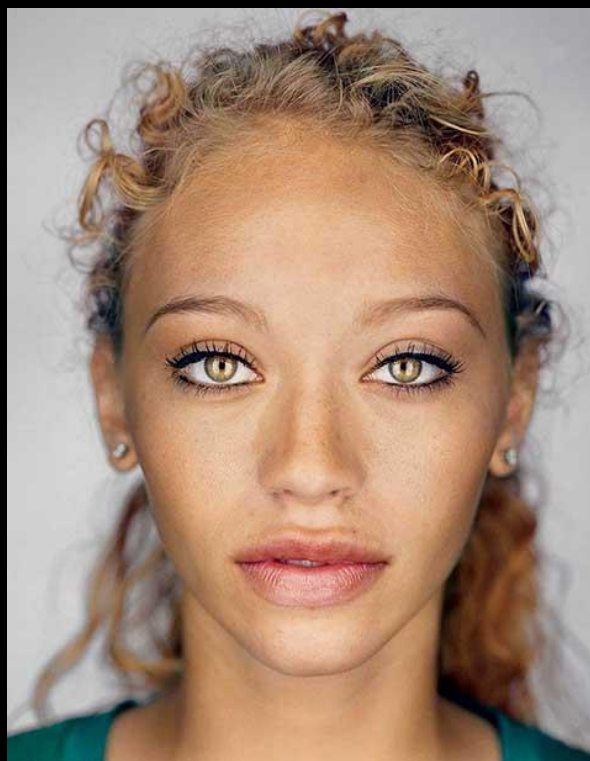


H





J



K



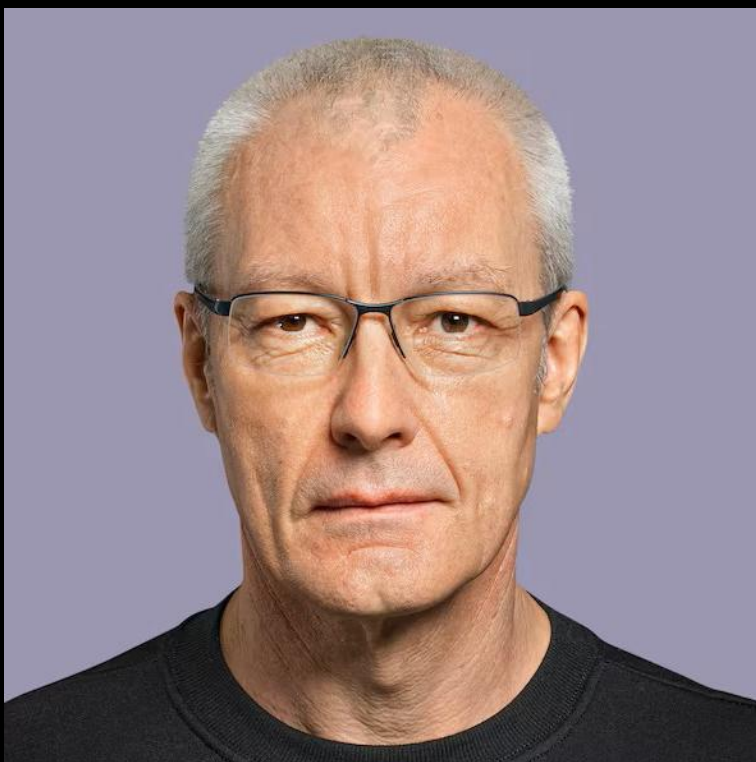
L



M



N



O



P



Q



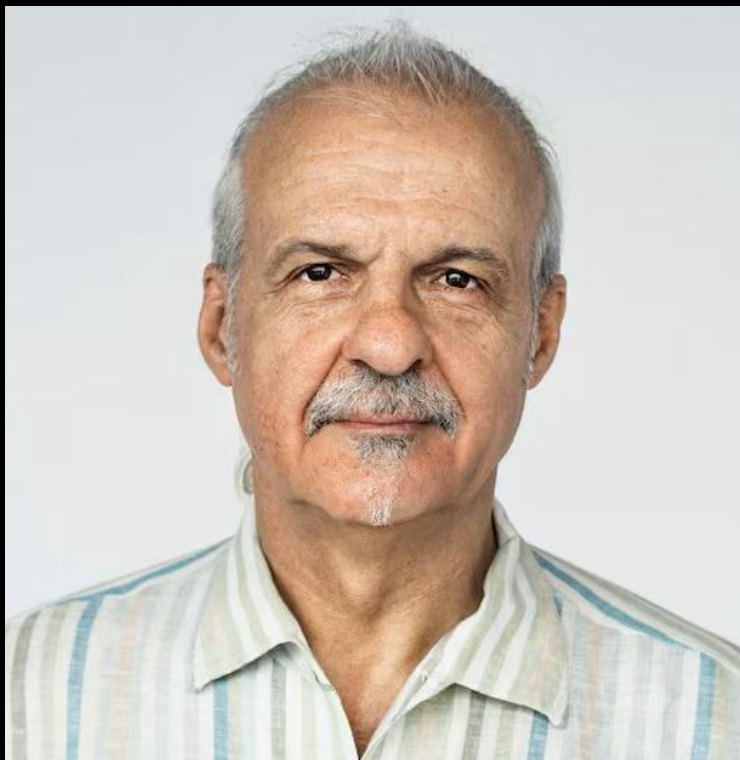
R



S



T



Exercício Prático

Anote a classe real de cada exemplo como: “feliz” ou “triste”

GABARITO

A



FELIZ

B



TRISTE

C



FELIZ

D



TRISTE

E



TRISTE

F



FELIZ

G



FELIZ

H



FELIZ

I



TRISTE

J



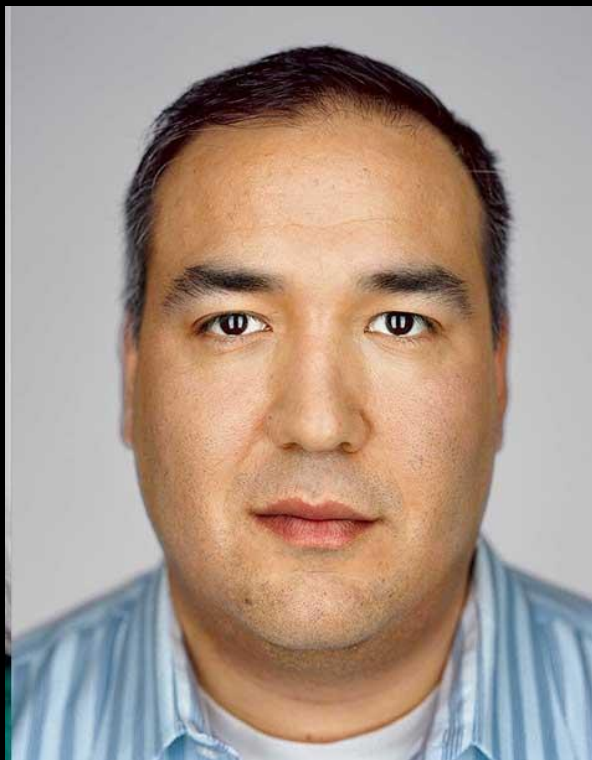
TRISTE

K



FELIZ

L



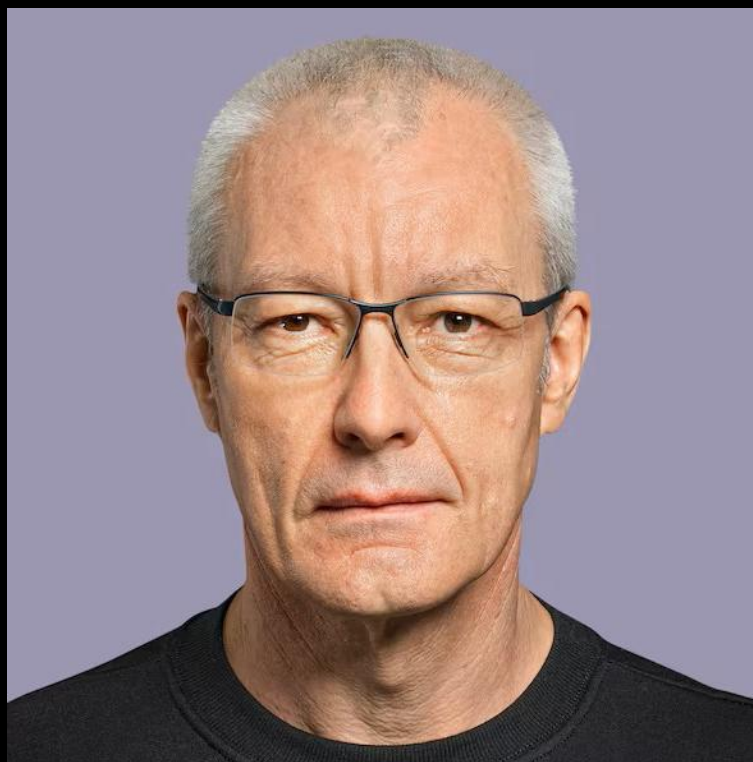
FELIZ

M



TRISTE

N



TRISTE

O



FELIZ

P



FELIZ

Q



TRISTE

R



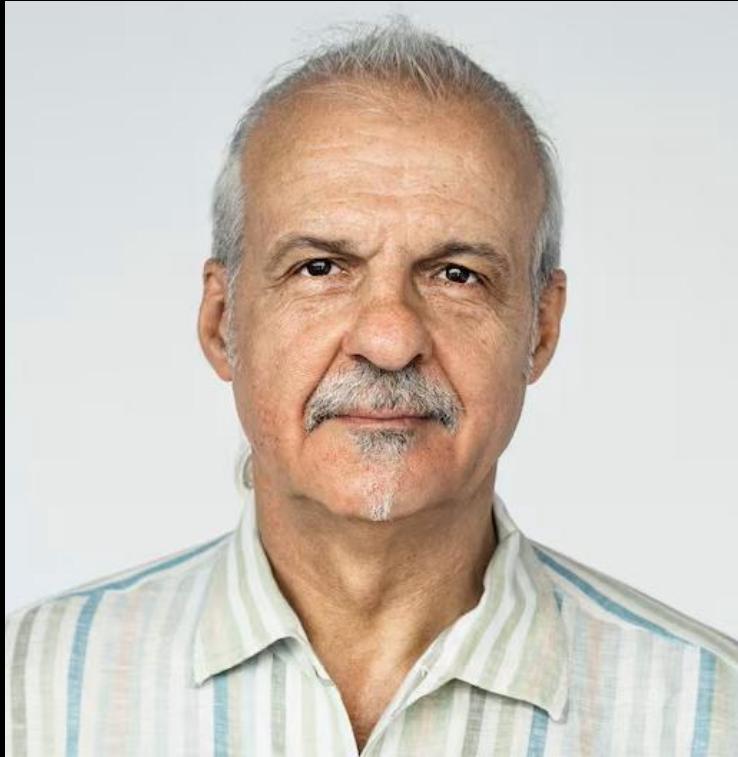
TRISTE

S



TRISTE

T



FELIZ

Matriz de Confusão

Matriz de Confusão

Entendendo as nomenclaturas

Tenho um sistema treinado em reconhecer Maças.
Eu apresento para meu sistema uma Maça.



Sistema retorna: Maça !



VERDADEIRO
POSITIVO
(VP)

ACERTOU !

Sistema retorna: Não é Maça !



FALSO
NEGATIVO
(FN)

ERROU !

Eu apresento pro sistema um abacaxi.



Sistema retorna: Maça !



FALSO
POSITIVO
(FP)

ERROU !

Sistema retorna: Não é Maça !



VERDADEIRO
NEGATIVO
(VN)

ACERTOU !




Exercício Prático

Monte a sua matriz de confusão



Accuracy (Acurácia)

A porcentagem de acertos no conjunto avaliado

	DOENTE 	SAUDÁVEL 
TESTE POSITIVO 	(VP)  Verdadeiro positivo	(FP)  Falso positivo
TESTE NEGATIVO 	(FN)  Falso negativo	(VN)  Verdadeiro negativo

1. Está doente e o teste é positivo: **verdadeiro positivo** → (VP)
2. Está saudável e o teste é negativo: **verdadeiro negativo** → (VN)
3. Está saudável e o teste é positivo: **falso positivo** → (FP)
4. Está doente e o teste é negativo: **falso negativo** → (FN)

$$\text{ACURÁCIA: } \frac{\text{VP} + \text{VN}}{\text{TOTAL (VP+VN+FN+FP)}}$$

Accuracy (Acurácia)

A porcentagem de acertos no conjunto avaliado

Definição

$$\text{Accuracy} = (\text{VP} + \text{VN}) / (\text{VP} + \text{VN} + \text{FP} + \text{FN})$$

Quando funciona bem

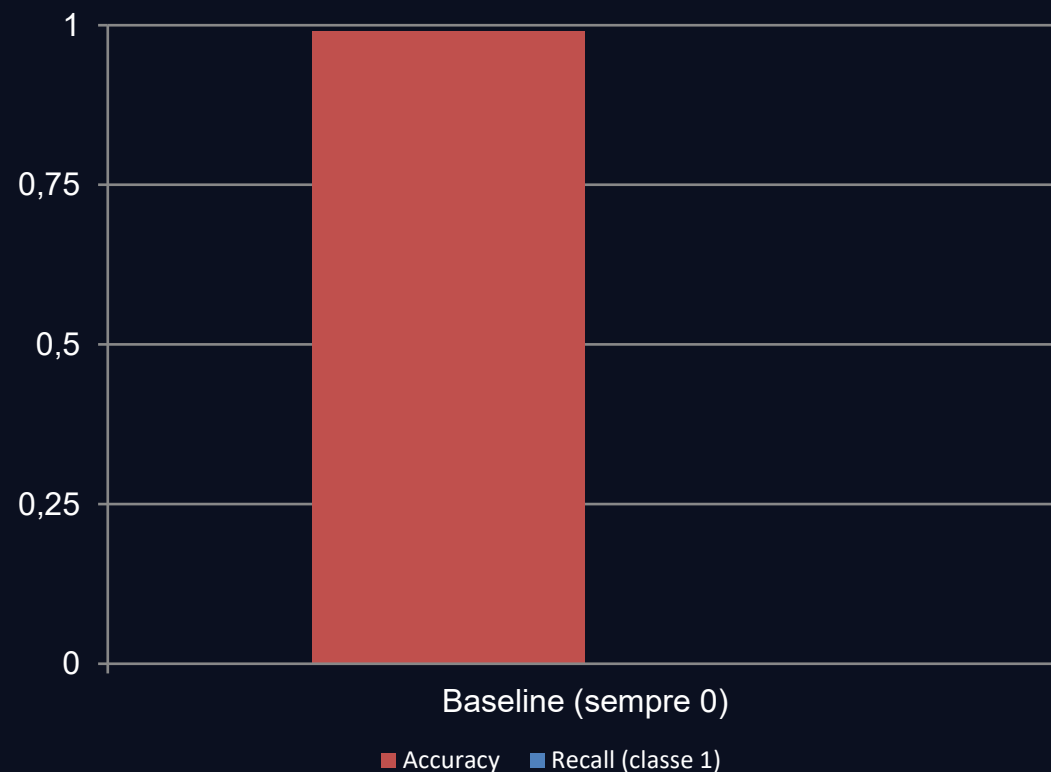
- Boa quando as classes são balanceadas
- Boa quando o custo de FP e FN é parecido

Armadilha clássica

- Dados desbalanceados: 99% da classe 0
- Modelo “sempre 0”: accuracy alta, mas recall da classe 1 = 0

⚠ Se o seu problema tem “raros importantes”, olhe além da accuracy.

Mini-exemplo (classe rara)



Accuracy é útil, mas não é onisciente. Combine com recall e F1.

Exercício Prático

Calcule a sua acurácia




Matriz de Confusão

O mapa dos acertos e erros por classe

Para binária (classe positiva = 1):

VN (Verdadeiro Negativo)	FP (Falso Positivo)
FN (Falso Negativo)	VP (Verdadeiro Positivo)

Linha = real | Coluna = previsto

 Cada célula é um tipo de erro. O que custa mais no seu caso: FP ou FN?

Como interpretar

- FP alto: muitos alarmes falsos (custo: tempo, atrito)
- FN alto: muitos casos perdidos (custo: risco, perda)
- Threshold muda o equilíbrio entre FP e FN


Onde isso aparece

- A matriz é base para Precision, Recall (sensibilidade) e F1
- Em multi-classe, vira uma matriz $N \times N$

Além da Acurácia

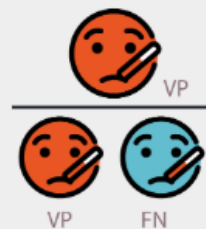
Precisão e Sensibilidade

Fórmulas (classe positiva)

	DOENTE	SAUDÁVEL
TESTE POSITIVO	(VP)  Verdadeiro positivo	(FP)  Falso positivo
TESTE NEGATIVO	(FN)  Falso negativo	(VN)  Verdadeiro negativo

Dos positivos reais, quantos eu consegui achar (RECALL) ?

SENSIBILIDADE



A **Sensibilidade** do teste é a proporção entre o número de doentes que o teste consegue detectar (VP) e o número total de doentes (ND). Em outras palavras, é a probabilidade de o teste ser positivo para uma pessoa doente: $P(\text{Teste} + \text{doente})$.

Dos positivos que eu marquei, quantos eram de verdade (PRECISION)?

PRECISÃO



A **Precisão** (ou valor preditivo positivo) é a relação entre a quantidade de pessoas doentes que testaram positivo (VP) e o número total de testes positivos (N+). Em outras palavras, é a probabilidade de você estar doente, dado que o teste deu positivo: $P(\text{Doente} / \text{Teste} +)$.

Exercício Prático

Calcule a sua precisão e sensibilidade



F1-Score

Média Harmônica entre precisão e sensibilidade

Fórmulas (classe positiva)

$Precision = TP / (TP + FP)$

$Recall = TP / (TP + FN)$

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

💡 Pune o modelo quando um dos dois (precision ou recall) é baixo.

Exemplos práticos

- Spam: priorize precision
- Diagnóstico: priorize recall
- Quando em dúvida: F1 (e olhe a matriz)

Intuição rápida

- F1: equilíbrio quando você precisa dos dois

É crucial para conjuntos de dados desbalanceados, equilibrando o custo de falsos positivos e falsos negativos. A pontuação varia de 0 (péssimo) a 1 (excelente).

Por que usar: Diferente da acurácia, o F1-Score não é enganado se uma classe dominar o conjunto de dados.

Aplicações: Ideal para diagnósticos médicos, recuperação de informações e classificação binária onde o equilíbrio é necessário.

Exercício Prático

Calcule o seu F1-score



Exemplo no Colab: treino e métricas

Dataset: Breast Cancer (scikit-learn)

O que faremos

- Carregar dados prontos (sem download)
- Treinar Logistic Regression com Pipeline
- Gerar previsões no teste

⚙ Rode em um Runtime CPU no Colab. É rápido.

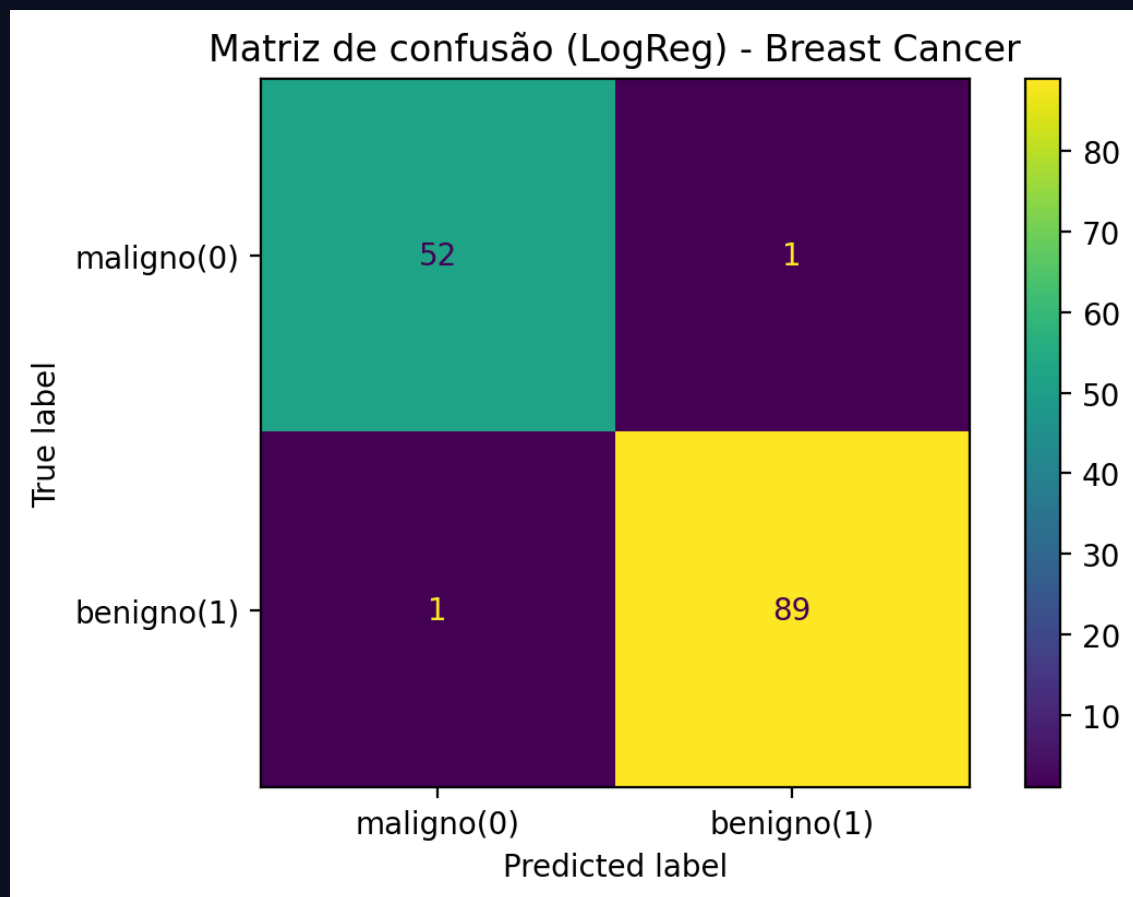
Código (Colab)

```
1 from sklearn.datasets import load_breast_cancer
2 from sklearn.model_selection import train_test_split
3 from sklearn.pipeline import make_pipeline
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.linear_model import LogisticRegression
6
7 X, y = load_breast_cancer(return_X_y=True)
8 X_tr, X_te, y_tr, y_te = train_test_split(
9     X, y, test_size=0.25, stratify=y, random_state=42
10 )
11
12 model = make_pipeline(StandardScaler(),
13     LogisticRegression(max_iter=2000))
14 model.fit(X_tr, y_tr)
15 pred = model.predict(X_te)
```

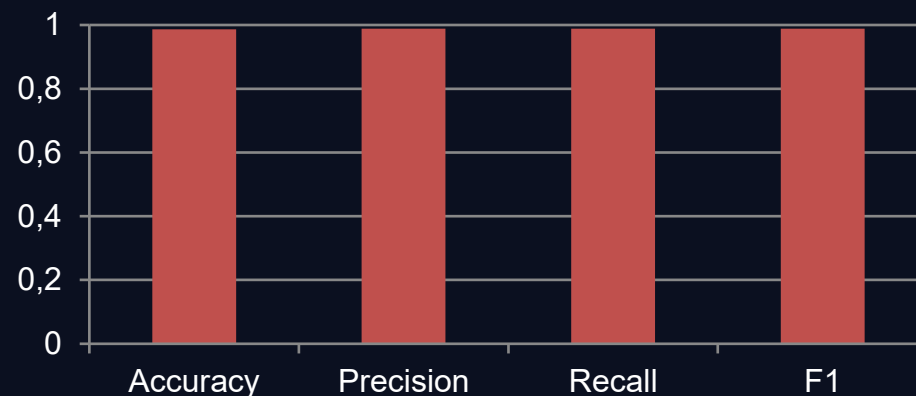

Exemplo: resultados e gráfico

Matriz de confusão + métricas no conjunto de teste

Matriz de confusão (teste)



Métricas (teste)



Código (Colab)


```
1 from sklearn.metrics import (accuracy_score,  
precision_score,  
2 recall_score, f1_score, confusion_matrix)  
3  
4 acc = accuracy_score(y_te, pred)  
5 prec = precision_score(y_te, pred)  
6 rec = recall_score(y_te, pred)  
7 f1 = f1_score(y_te, pred)  
8 cm = confusion_matrix(y_te, pred)  
9  
10 acc, prec, rec, f1, cm
```

Validação Cruzada (k-fold)

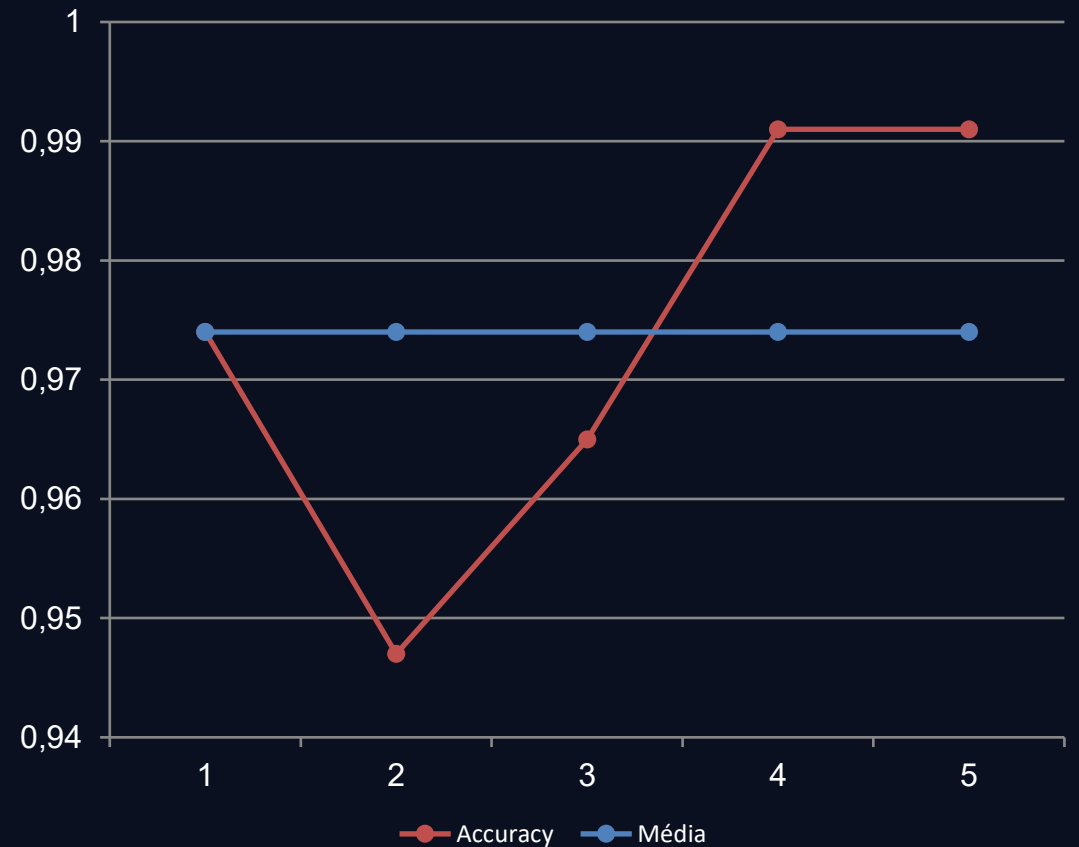
Estimando performance de generalização com menos “sorte”

Ideia

- Divida o dataset em k partes (folds)
- Treine em k-1 folds e teste no fold restante
- Repita k vezes e compute média e desvio padrão
- Use StratifiedKFold quando houver desbalanceamento

 CV é ótima para comparar modelos e hiperparâmetros. Guarde um teste final separado.

Exemplo: Accuracy por fold (k=5)



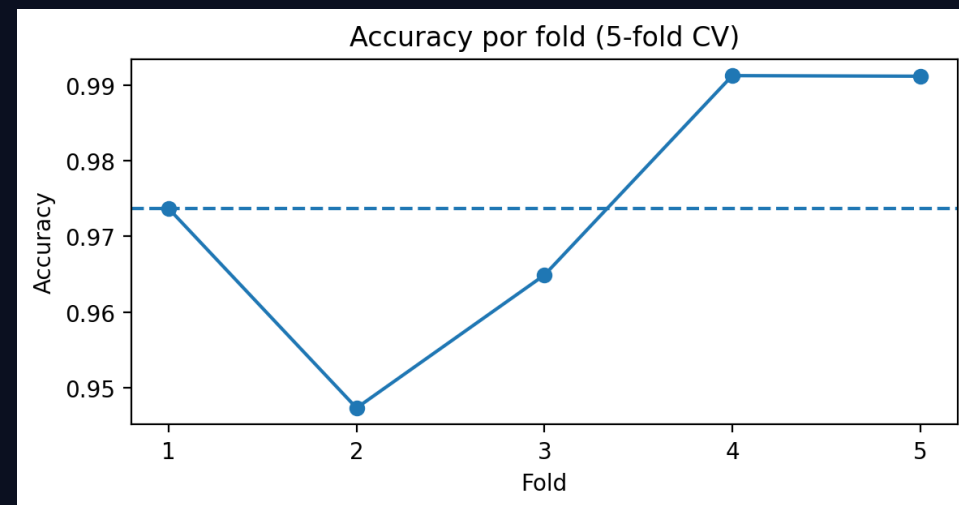
Exemplo no Colab: cross_validate

Calculando métricas em k-fold e gerando gráficos

Código (Colab)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import StratifiedKFold, cross_validate
4
5 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
6 res = cross_validate(model, X, y, cv=cv,
7                      scoring=['accuracy', 'precision', 'recall', 'f1'])
8
9 accs = res['test_accuracy']
10 print('accuracy média:', accs.mean(), '±', accs.std())
11
12 plt.plot(range(1, len(accs)+1), accs, marker='o')
13 plt.axhline(accs.mean(), linestyle='--')
14 plt.xticks(range(1, len(accs)+1))
15 plt.xlabel('Fold')
16 plt.ylabel('Accuracy')
17 plt.title('Accuracy por fold (CV)')
18 plt.show()
```

Gráfico (exemplo)



Exercício (mão na massa)

Dataset: Wine (3 classes) para vocês criarem seus próprios gráficos

Tarefas

- Carregue o dataset Wine e faça um train/test split estratificado
- Treine um modelo (LogReg, RandomForest, SVM...)
- Calcule accuracy e matriz de confusão (3×3) e plote o heatmap
- Calcule F1 macro e F1 weighted e compare
- Faça validação cruzada (k=5) para accuracy e F1 macro
- Gere 2 gráficos: matriz de confusão e performance por fold

Starter code (Colab)

```
1 from sklearn.datasets import load_wine
2 from sklearn.model_selection import train_test_split,
  StratifiedKFold, cross_validate
3 from sklearn.metrics import ConfusionMatrixDisplay, f1_score
4 from sklearn.ensemble import RandomForestClassifier
5 import matplotlib.pyplot as plt
6
7 X, y = load_wine(return_X_y=True)
8 X_tr, X_te, y_tr, y_te = train_test_split(
9     X, y, test_size=0.25, stratify=y, random_state=42
10 )
11
12 clf = RandomForestClassifier(random_state=42)
13 clf.fit(X_tr, y_tr)
14 pred = clf.predict(X_te)
15
16 # 1) Matriz de confusão
17 ConfusionMatrixDisplay.from_predictions(y_te, pred)
18 plt.title('Wine: matriz de confusão')
19 plt.show()
20
21 # 2) F1 macro vs weighted
22 print('F1 macro   :', f1_score(y_te, pred, average='macro'))
23 print('F1 weighted:', f1_score(y_te, pred, average='weighted'))
24
25 # 3) CV + gráfico por fold
26 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
27 res = cross_validate(clf, X, y, cv=cv,
28 scoring=['accuracy', 'f1_macro'])
29 plt.plot(res['test_accuracy'], marker='o', label='accuracy')
30 plt.plot(res['test_f1_macro'], marker='o', label='f1_macro')
31 plt.legend(); plt.title('Performance por fold')
32 plt.xlabel('Fold'); plt.show()
```

Boas práticas (checklist rápido)

Para não cair em “números bonitos” 🙄

- Separe um conjunto de teste final e não “olhe” para ele durante ajustes
- Use Pipeline para evitar vazamento de dados (scaler só com treino)
- Sempre inspecione a matriz de confusão (por classe)
- Em desbalanceamento: prefira F1, precision/recall e métricas macro
- Em CV: reporte média e desvio padrão, não só a média
- Escolha a métrica alinhada ao custo real do erro (FP vs FN)

Resumo em 4 frases

- Accuracy é um resumo, não um veredito.
- A matriz de confusão mostra onde o modelo erra.
- F1 combina precision e recall quando ambos importam.
- Validação cruzada reduz a dependência de um único split.