

I. Pen-and-paper

- 1) [3.0v] Compute the F1-measure of a NN with $k = 5$ and Hamming distance using a leave-one-out evaluation schema. Show all calculus.

1) Distâncias para x_i distância de Hamming 2 se 0 se 2
parâmetros iguais 1 se 1 parâmetro igual
2 se 0 parâmetros iguais

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
x_1		2	1	0	1	1	1	2
x_2			1	2	1	1	1	0
x_3				1	2	2	0	1
x_4					1	1	1	2
x_5						0	2	1
x_6							2	1
x_7								1
x_8								

vizinhos de x_1 : x_4 ; x_3 ; x_5 ; x_6 ; x_7
 x_2 : x_8 ; x_3 ; x_5 ; x_6 ; x_7
 x_3 : x_7 ; x_1 ; x_2 ; x_4 ; x_8
 x_4 : x_1 ; x_3 ; x_5 ; x_6 ; x_7
 x_5 : x_1 ; x_2 ; x_4 ; x_6 ; x_8
 x_6 : x_1 ; x_2 ; x_4 ; x_5 ; x_8
 x_7 : x_1 ; x_2 ; x_3 ; x_4 ; x_8
 x_8 : x_1 ; x_3 ; x_5 ; x_6 ; x_7

x_i classificação para cada x_i de acordo com os
seus vizinhos

x_1 2P 3N \rightarrow N	x_5 3P 2N \rightarrow P
x_2 1P 4N \rightarrow N	x_6 3P 2N \rightarrow P
x_3 3P 2N \rightarrow P	x_7 4P 1N \rightarrow P
x_4 2P 3N \rightarrow N	x_8 2P 3N \rightarrow N

matriz de confusão
Real

	P	N
P	1	3
N	3	1

$$P = \frac{TP}{TP+FP}$$

$$R = \frac{TP}{TP+FN}$$

$$F_{score} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

$$P = \frac{1}{1+3} = 0.25 \quad R = \frac{1}{1+3} = 0.25$$

$$F_1 = \frac{2}{\frac{1}{0.25} + \frac{1}{0.25}} = \frac{2}{\frac{4}{0.25}} = 0.25$$

- 2) [2.5v] Propose a new metric (distance and/or) that improves the latter's performance (i.e., k the F1-measure) by three fold.

2) para que o F1 seja aumente 3x
 a nossa matriz de confusão será do tipo

	Real	
	P	N
Predicted	P	3
	N	1

ficando o F1 score = $\frac{3}{4}$

decidi usar uma distancia ponderada do tipo

$$d(x_i, x_j) = w_1 (y_{1i} - y_{1j}) + w_2 (y_{2i} - y_{2j})$$

usando distancias de hamming para y_1 e y_2
 e os pesos $w_1 = 1.10$ e $w_2 = 0.10$

e $k = 3$

Usamos $k = 3$ porque existem 4 positivos e 4 negativos
 e quando estamos a verificar se uma instancia é
 de uma certa classe apenas sobre 3 instancias desta
 mesma classe logo não se podem ser afetados
 ter os mesmos resultados afetados por dados não
 pretendidos.

Estes pesos foram escolhidos para favorecer a escolha dos vizinhos que sejam da mesma classe.

- 3) [2.5v] Considering the nine training observations, learn a Bayesian classifier assuming:
i) y_1 and y_2 are dependent; ii) $\{y_1, y_2\}$ and $\{y_3\}$ variable sets are independent and equally important; and iii) y_3 is normally distributed. Show all parameters.

3

Posterior: $P(\text{out} = \text{class} / y_1 = x_1, y_2 = x_2, y_3 = x_3) \stackrel{\text{T. Bayes}}{=} =$

$$= \frac{P(\text{out} = \text{class}) P(y_1 = x_1, y_2 = x_2, y_3 = x_3 / \text{out} = \text{class})}{P(y_1 = x_1, y_2 = x_2, y_3 = x_3)}$$

①
②
3

$$= \frac{P(\text{out} = \text{class}) P(y_1 = x_1, y_2 = x_2 / \text{out} = \text{class}) P(y_3 = x_3 / \text{out} = \text{class})}{P(y_1 = x_1, y_2 = x_2) P(y_3 = x_3)}$$

Para uma nova observação (y_1, y_2, y_3) , nós queremos classificar a observação $P(\text{class} / O)$ para todas as classes $C \in N$ e escolher a com maior probabilidade

$$\hat{C} = \underset{C}{\operatorname{argmax}} (P(C / O)) = \underset{C}{\operatorname{argmax}} \frac{P(y_1, y_2 / C) P(y_3 / C)}{P(y_1, y_2) P(y_3)} =$$

$$= \underset{C}{\operatorname{argmax}} P(y_1, y_2 / C) P(y_3 / C) \quad \text{igual em todos os } C$$

① Priors: $P(\text{out} = P) = \frac{5}{9}$ $P(\text{out} = N) = \frac{4}{9}$

② $P(y_1 = A, y_2 = 0 / \text{out} = P) = \frac{2}{5}$ $P(y_1 = A, y_2 = 0 / \text{out} = N) = 0$
 $P(y_1 = A, y_2 = 1 / \text{out} = P) = \frac{1}{5}$ $P(y_1 = A, y_2 = 1 / \text{out} = N) = \frac{1}{4}$
 $P(y_1 = B, y_2 = 0 / \text{out} = P) = \frac{1}{5}$ $P(y_1 = B, y_2 = 0 / \text{out} = N) = \frac{2}{4}$
 $P(y_1 = B, y_2 = 1 / \text{out} = P) = \frac{1}{5}$ $P(y_1 = B, y_2 = 1 / \text{out} = N) = \frac{1}{4}$

③ $P(y_1 = A, y_2 = 0) = \frac{2}{9}$
 $P(y_1 = A, y_2 = 1) = \frac{2}{9}$

$P(y_1 = B, y_2 = 0) = \frac{3}{9}$
 $P(y_1 = B, y_2 = 1) = \frac{2}{9}$

③

$$\mu_{out=P} = \frac{1.1 + 0.8 + 0.5 + 0.9 + 0.8}{5} = 0.82$$

$$\mu_{out=N} = \frac{1 + 0.9 + 1.2 + 0.9}{4} = 1$$

$$\sigma_{out=P} = \sqrt{\frac{(1.1-0.82)^2 + (0.8-0.82)^2 + (0.5-0.82)^2 + (0.9-0.82)^2 + (0.8-0.82)^2}{5-1}}$$

$$\approx 0.2168$$

$$\sigma_{out=N} = \sqrt{\frac{(1-1)^2 + (0.9-1)^2 + (1.2-1)^2 + (0.9-1)^2}{4-1}}$$

$$\approx 0.1414$$

$$R(x|K, G^2) = \frac{1}{\sqrt{2\pi} \times 0.2168} \exp\left(-\frac{x^2}{2 \times (0.2168)^2}\right)$$

$$P(y_3 = x | y_{out} = class) = \frac{1}{\sqrt{2\pi} \times \sigma_{class}} \exp\left(-\frac{(x - \mu_{class})^2}{2 \times (\sigma_{class})^2}\right)$$

③'

$$\mu = \frac{1.1 + 0.8 + 0.5 + 0.9 + 0.8 + 1 + 0.9 + 1.2 + 0.9}{9} = 0.9$$

$$\sigma = \sqrt{\frac{(1.1-0.9)^2 + (0.8-0.9)^2 + (0.5-0.9)^2 + (0.9-0.9)^2 + (0.8-0.9)^2 + (1-0.9)^2 + (0.9-0.9)^2 + (1.2-0.9)^2 + (0.9-0.9)^2}{9-1}}$$

$$= 0.2$$

$$P(y_3 = x) = \frac{1}{\sqrt{2\pi} \times 0.2} \exp\left(-\frac{(x - \mu)^2}{2 \times 0.2^2}\right)$$

Agora já temos todos os parâmetros para aplicar o classificador

- 4) [2.5v] Under a MAP assumption, classify each testing observation showing all your calculus.

4

preciso
sabemos que as probabilidades vão ser da forma

$$P(\text{out} = \text{class} / y_1 = x_1, y_2 = x_2, y_3 = x_3) = \frac{P(\text{out} = \text{class})}{P(y_1 = x_1, y_2 = x_2) P(y_3 = x_3)}$$

sabemos que $P(\text{class} = P/x) + P(\text{class} = N/x) = 1$
~~é que o denominador não é igual~~
 e os denominadores são iguais para P e N
 assim sendo

para $x = (A, 1, 0.8)$
 temos $P(y_1 = A, y_2 = 1 / \text{out} = P) = \frac{1}{9}$

$$P(y_3 = 0.8 / \text{out} = P) = \frac{1}{\sqrt{2\pi} \times 0.2168} \times \exp\left(-\frac{(0.8 - 0.82)^2}{2 \times 0.2168}\right)$$

$$\approx 2.7932 \times 1.83$$

$$P(\text{out} = P) P(y_1 = A, y_2 = 1 / \text{out} = P) P(y_3 = 0.8 / \text{out} = P) \approx 0.20$$

para $\text{out} = N$
 $P(y_1 = A, y_2 = 1 / \text{out} = N) = \frac{1}{9}$

$$P(y_3 = 0.8 / \text{out} = N) = \frac{1}{\sqrt{2\pi} \times 0.1414} \times \exp\left(-\frac{(0.8 - 1)^2}{2 \times 0.1414}\right) \approx 1.0376$$

$$P(\text{out} = N) P(y_1 = A, y_2 = 1 / \text{out} = N) P(y_3 = 0.8 / \text{out} = N) \approx \frac{1}{9} \times \frac{1}{9} \times 1.0376 \approx 0.115$$

Como $0.2 > 0.115$ escolhemos o primeiro valor $(A, 1, 0.8)$ e escolhemos o positivo

para $(B, 1, 1)$

$$P(y_3 = 1 / \text{out} = P) = \frac{1}{\sqrt{2\pi} \times 0.2168} \times \exp\left(-\frac{(1 - 0.82)^2}{2 \times (0.2168)^2}\right)$$

$$\approx 1.304$$

$$P(\text{out} = P) P(y_1 = B, y_2 = 1 / \text{out} = P) P(y_3 = 1 / \text{out} = P)$$

$$= \frac{5}{9} \times \frac{1}{5} \times 1.304 \approx 0.1449$$

$$P(y_3 = 1 / \text{out} = N) = \frac{1}{\sqrt{2\pi} \times 0.1414} \times \exp\left(-\frac{(1 - 1)^2}{2 \times (0.1414)^2}\right)$$

$$= 2.821$$

$$P(\text{out} = N) P(y_1 = B, y_2 = 1 / \text{out} = N) P(y_3 = 1 / \text{out} = N)$$

$$= \frac{4}{9} \times \frac{1}{4} \times 2.821 = 0.3134$$

logo este valor é negativo

para $(B, 0, 0.9)$

$$P(y_3 = 0.9 / \text{out} = P) = \frac{1}{\sqrt{2\pi} \times 0.2168} \times \exp\left(-\frac{(0.9 - 0.82)^2}{2 \times (0.2168)^2}\right)$$

$$\approx 1.719$$

$$P(\text{out} = P) P(y_1 = B, y_2 = 0 / \text{out} = P) P(y_3 = 0.9 / \text{out} = P) \approx$$

$$\frac{5}{9} \times 1.719 \times \frac{1}{5} \approx 0.191$$

$$P(y_3 = 0.9 / \text{out} = N) = \frac{1}{\sqrt{2\pi} \times 0.1414} \times \exp\left(-\frac{(0.9 - 1)^2}{2 \times (0.1414)^2}\right)$$

$$\approx 2.197$$

$$P(\text{out} = N) P(y_1 = B, y_2 = 0 / \text{out} = N) P(y_3 = 0.9 / \text{out} = N) =$$

$$= 2.197 \times \frac{4}{9} \times \frac{2}{4} = 0.488$$

assim sendo
 $0.191 < 0.488$ logo

$(B, 0, 0.9)$ é negativo

- 5) [2.5v] Using a naïve Bayes under a ML assumption, classify the new sentence

"*I like to run*" . For the likelihoods calculation consider the following formula,

$$p(t_i|c) = (freq(t_i) + 1)/(N_c + V) ,$$

where t_i represents a certain term i , V the number of unique terms in the vocabulary, and N_c the total number of terms in class c . Show all calculus.

⑤ **Class Positivas** termos únicos:
Amazing run
I like it
NP = 5

Negativas termos únicos:
Too tired.
Bad run
NN = 4

total de termos = 8 = V

usar a fórmula para calcular a probabilidade de cada termo

$$P(t_i | c) = \frac{\text{freq}(t_i + 1)}{N_c + V}$$

para P

$$P("I" / \text{out} = P) = \frac{1+1}{5+8} = \frac{2}{13}$$

$$P("like" / \text{out} = P) = \frac{1+1}{5+8} = \frac{2}{13}$$

$$P("too" / \text{out} = P) = \frac{1+0}{5+8} = \frac{1}{13}$$

$$P("run" / \text{out} = P) = \frac{1+1}{5+8} = \frac{2}{13}$$

para N

$$P("I" / \text{out} = N) = \frac{0+1}{4+8} = \frac{1}{12}$$

$$P("like" / \text{out} = N) = \frac{0+1}{4+8} = \frac{1}{12}$$

$$P("too" / \text{out} = N) = \frac{0+1}{4+8} = \frac{1}{12}$$

$$P("run" / \text{out} = N) = \frac{1+1}{4+8} = \frac{2}{12}$$

$$P(\text{"I like to run"} / \text{out} = P) = P(\text{"I"} / \text{out} = P) \times P(\text{"like"} / \text{out} = P) \times \\ \times P(\text{"to"} / \text{out} = P) \times P(\text{"run"} / \text{out} = P) = \\ = \frac{2}{13} \times \frac{2}{13} \times \frac{1}{13} \times \frac{2}{13} = \frac{8}{28561} = 2.8 \times 10^{-4}$$

$$P(\text{"I like to run"} / \text{out} = N) = P(\text{"I"} / \text{out} = N) \times P(\text{"like"} / \text{out} = N) \times \\ \times P(\text{"to"} / \text{out} = N) \times P(\text{"run"} / \text{out} = N) = \frac{1}{12} \times \frac{1}{12} \times \frac{1}{12} \times \frac{2}{12} = \frac{2}{20736} \\ \approx 9.645 \times 10^{-5}$$

$P(\text{"I like to run"} / P) > P(\text{"I like to run"} / N)$ ~~Logo~~
a frase é classificada como Positiva
para (P)
k argmax_n $P(h) \prod_{i=1}^n \frac{P(y_i = x_i)}{2-1}$

Como ambos têm o mesmo Prior basta apenas
comparar as probabilidades obtidas anteriormente
e com isso podemos afirmar que a frase é
classificada como positiva $P(P) = P(N) = 0.5$

II. Programming and critical analysis

1) Compare the performance of a kNN $k = 5$ with and a naïve Bayes with Gaussian assumption (consider all remaining parameters as default):

a) [1.0v] Plot two boxplots with the fold accuracies for each classifier. Is there one more stable than the other regarding performance? Why do you think that is the case? Explain.

```
import pandas as pd
from sklearn.model_selection import StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn import metrics
import matplotlib.pyplot as plt

# Reading the CSV file
df = pd.read_csv('heart-disease.csv')

X = df.drop('target', axis=1)
y = df['target']

knn_predictor = KNeighborsClassifier(n_neighbors=5)
nb_predictor = GaussianNB()

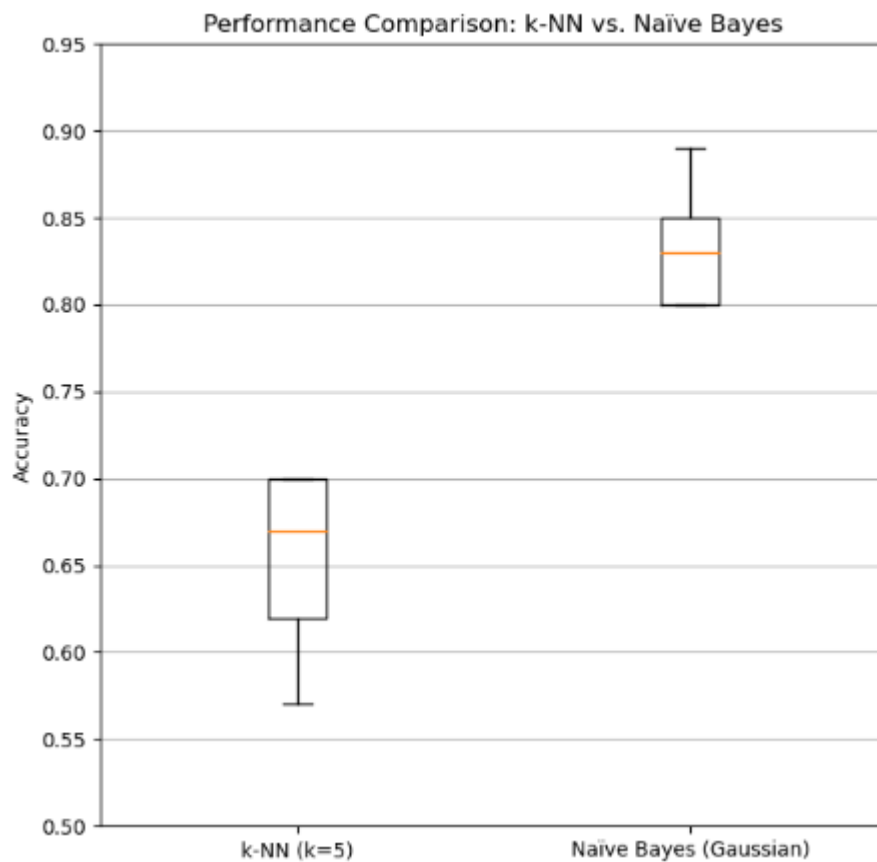
folds = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
knn_accuracies = []
nb_accuracies = []

# iterate per fold
for train_k, test_k in folds.split(X, y):
    X_train, X_test = X.iloc[train_k], X.iloc[test_k]
    y_train, y_test = y.iloc[train_k], y.iloc[test_k]

    # train and assess k-NN
    knn_predictor.fit(X_train, y_train)
    y_pred = knn_predictor.predict(X_test)
    knn_accuracies.append(round(metrics.accuracy_score(y_test, y_pred), 2))

    # train and assess Naïve Bayes with Gaussian assumption
    nb_predictor.fit(X_train, y_train)
    y_pred = nb_predictor.predict(X_test)
    nb_accuracies.append(round(metrics.accuracy_score(y_test, y_pred), 2))

# plots
plt.figure(figsize=(7, 7))
plt.boxplot([knn_accuracies, nb_accuracies], labels=['k-NN (k=5)', 'Naïve Bayes (Gaussian)'])
plt.title('Performance Comparison: k-NN vs. Naïve Bayes')
plt.ylabel('Accuracy')
plt.ylim(0.5, 0.95)
plt.grid(axis='y')
plt.show()
```



O modelo Naïve Bayes tem consistentemente maior precisão em comparação com o modelo k-NN. O modelo k-NN exibe maior variabilidade no desempenho, como indicado pela maior dispersão no boxplot. O modelo Naïve Bayes é mais estável em termos de desempenho, provavelmente devido à sua natureza probabilística, enquanto a precisão do k-NN depende mais da estrutura local dos dados e é mais sensível ao tamanho e distribuição dos vizinhos.

b) [1.0v] Report the accuracy of both models, this time scaling the data with a Min-Max scaler before training the models. Explain the impact that this preprocessing step has on the performance of each model, providing an explanation for the results.

```
import pandas as pd
from sklearn.model_selection import StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.preprocessing import MinMaxScaler
from sklearn import metrics
import matplotlib.pyplot as plt

# Reading the CSV file
df = pd.read_csv('heart-disease.csv')

X = df.drop('target', axis=1)
y = df['target']

# Initializing the models
knn_predictor = KNeighborsClassifier(n_neighbors=5)
nb_predictor = GaussianNB()

# Stratified k-folds
folds = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
knn_accuracies = []
nb_accuracies = []

# Min-Max Scaler
scaler = MinMaxScaler()

# iterate per fold
for train_k, test_k in folds.split(X, y):
    X_train, X_test = X.iloc[train_k], X.iloc[test_k]
    y_train, y_test = y.iloc[train_k], y.iloc[test_k]

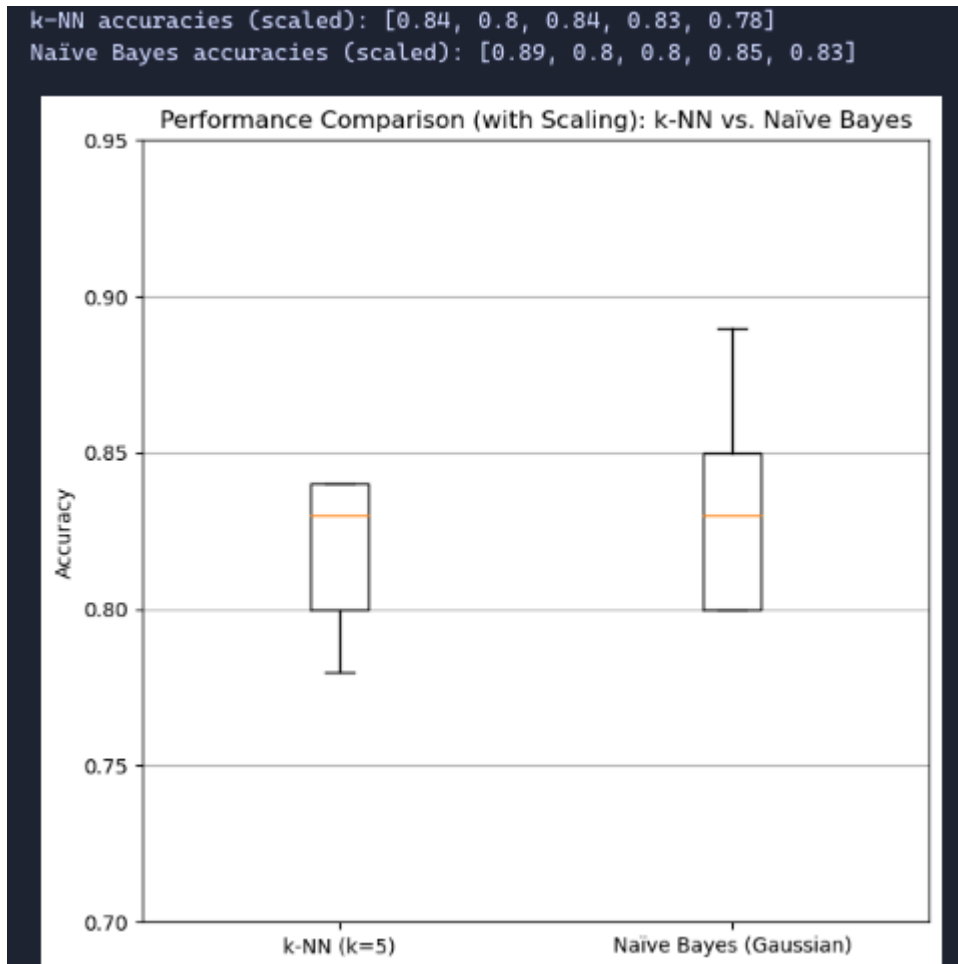
    # Scaling the data using Min-Max Scaler
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)

    # Train and assess k-NN
    knn_predictor.fit(X_train_scaled, y_train)
    y_pred_knn = knn_predictor.predict(X_test_scaled)
    knn_accuracies.append(round(metrics.accuracy_score(y_test, y_pred_knn), 2))

    # Train and assess Naïve Bayes with Gaussian assumption
    nb_predictor.fit(X_train_scaled, y_train)
    y_pred_nb = nb_predictor.predict(X_test_scaled)
    nb_accuracies.append(round(metrics.accuracy_score(y_test, y_pred_nb), 2))

# Report the accuracy
print("k-NN accuracies (scaled):", knn_accuracies)
print("Naïve Bayes accuracies (scaled):", nb_accuracies)

# Boxplot to compare
plt.figure(figsize=(7, 7))
plt.boxplot([knn_accuracies, nb_accuracies], labels=['k-NN (k=5)', 'Naïve Bayes (Gaussian)'])
plt.title('Performance Comparison (with Scaling): k-NN vs. Naïve Bayes')
plt.ylabel('Accuracy')
plt.ylim(0.7, 0.95)
plt.grid(axis='y')
plt.show()
```



O desempenho do k-NN melhorou significativamente após a escala, mostrando precisões mais altas e consistentes em todos os folds. O desempenho do Naïve Bayes permaneceu relativamente estável, sem mudanças perceptíveis após a escala. Isso indica que o k-NN é mais sensível à escala dos dados, provavelmente porque depende de métricas de distância, enquanto o Naïve Bayes não depende de distância e, portanto, é menos afetado pela escala das características.

c) [1.0v] Using scipy, test the hypothesis “the kNN model is statistically superior to naïve Bayes regarding accuracy”, asserting whether it is true

```
import pandas as pd
from sklearn.model_selection import StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn import metrics
from scipy import stats

# Reading the CSV file
df = pd.read_csv('heart-disease.csv')

X = df.drop('target', axis=1)
y = df['target']

knn_predictor = KNeighborsClassifier(n_neighbors=5)
nb_predictor = GaussianNB()

# 10-fold stratified cross-validator with shuffling
folds = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)

knn_accuracies = []
nb_accuracies = []

# iterate per fold
for train_k, test_k in folds.split(X, y):
    X_train, X_test = X.iloc[train_k], X.iloc[test_k]
    y_train, y_test = y.iloc[train_k], y.iloc[test_k]

    # train and assess k-MN
    knn_predictor.fit(X_train, y_train)
    y_pred = knn_predictor.predict(X_test)
    knn_accuracies.append(round(metrics.accuracy_score(y_test, y_pred), 2))

    # train and assess Naive Bayes with Gaussian assumption
    nb_predictor.fit(X_train, y_train)
    y_pred = nb_predictor.predict(X_test)
    nb_accuracies.append(round(metrics.accuracy_score(y_test, y_pred), 2))

# kNN is better than naïve Bayes?
res = stats.ttest_rel(knn_accuracies, nb_accuracies, alternative='greater')
print(f"KNN > naïve Bayes? pval=", res.pvalue)
```

✓ 0.1s
KNN > naïve Bayes? pval= 0.998415501126768

Python

KNN > naïve Bayes? pval= 0.998415501126768 como pval > valores de significancia usuais (ou seja, pval é superior a 0.01, 0.05 e 0.1), não é possível rejeitar a hipótese nula e afirmar que o KNN é estatisticamente superior ao Naive Bayes.

- 2) Using a 80-20 train-test split, vary the number of neighbors of a kNN classifier using weights $k = \{1, 5, 10, 20, 30\}$. Additionally, for each k , train one classifier using uniform and distance weights.
- a. [1.0v] Plot the train and test accuracy for each model.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics
import matplotlib.pyplot as plt

# Reading the CSV file
df = pd.read_csv('heart-disease.csv')

X = df.drop('target', axis=1)
y = df['target']

# Split the data into 80% train and 20% test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# List of k values to try
k_values = [1, 5, 10, 20, 30]

# Store accuracies
train_accuracies_uniform = []
test_accuracies_uniform = []
train_accuracies_distance = []
test_accuracies_distance = []

# Iterate over each k value
for k in k_values:
    # Uniform weights
    knn_uniform = KNeighborsClassifier(n_neighbors=k, weights='uniform')
    knn_uniform.fit(X_train, y_train)
    train_acc_uniform = metrics.accuracy_score(y_train, knn_uniform.predict(X_train))
    test_acc_uniform = metrics.accuracy_score(y_test, knn_uniform.predict(X_test))
    train_accuracies_uniform.append(train_acc_uniform)
    test_accuracies_uniform.append(test_acc_uniform)

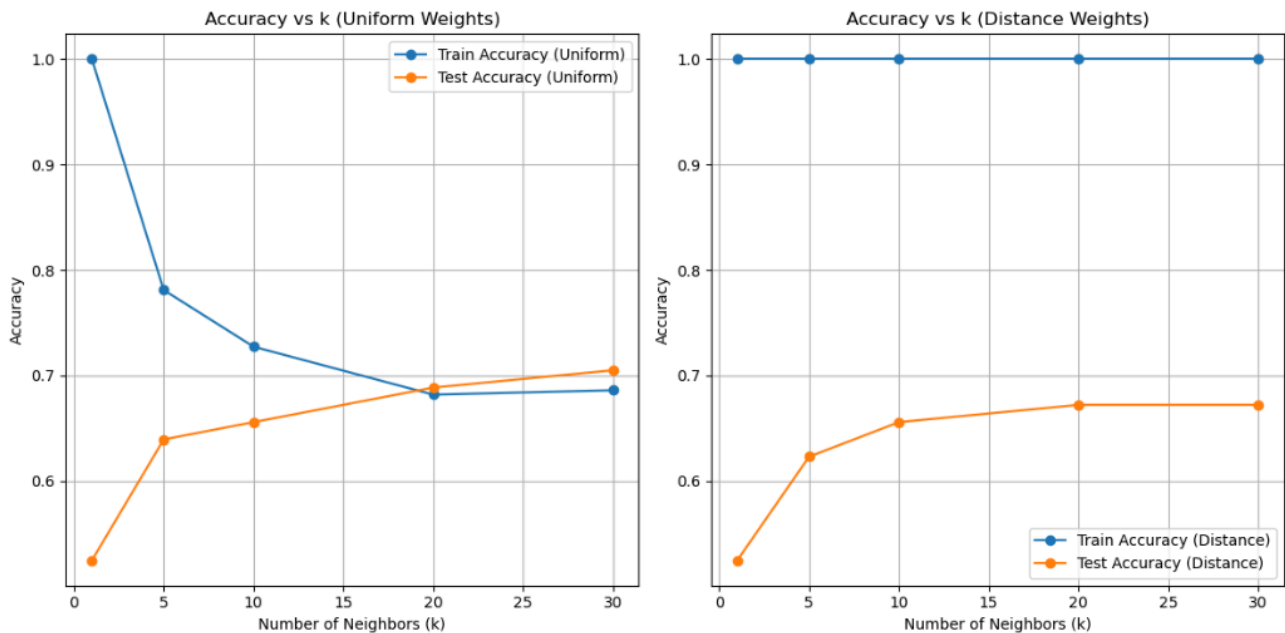
    # Distance weights
    knn_distance = KNeighborsClassifier(n_neighbors=k, weights='distance')
    knn_distance.fit(X_train, y_train)
    train_acc_distance = metrics.accuracy_score(y_train, knn_distance.predict(X_train))
    test_acc_distance = metrics.accuracy_score(y_test, knn_distance.predict(X_test))
    train_accuracies_distance.append(train_acc_distance)
    test_accuracies_distance.append(test_acc_distance)

# Plot the results
plt.figure(figsize=(12, 6))

# Uniform weights plot
plt.subplot(1, 2, 1)
plt.plot(k_values, train_accuracies_uniform, marker='o', label='Train Accuracy (Uniform)')
plt.plot(k_values, test_accuracies_uniform, marker='o', label='Test Accuracy (Uniform)')
plt.title('Accuracy vs k (Uniform Weights)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.legend()
plt.grid(True)

# Distance weights plot
plt.subplot(1, 2, 2)
plt.plot(k_values, train_accuracies_distance, marker='o', label='Train Accuracy (Distance)')
plt.plot(k_values, test_accuracies_distance, marker='o', label='Test Accuracy (Distance)')
plt.title('Accuracy vs k (Distance Weights)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()
```

b. [1.5v] Explain the impact of increasing the neighbors on the generalization ability of the models.

Pesos Uniformes (Gráfico 1 à Esquerda):

$k=1$: O modelo apresenta uma acurácia de treino perfeita, mas a acurácia de teste é baixa. Este é um claro indicativo de overfitting, onde o modelo está a memorizar os dados de treino, mas falha em generalizar para dados de teste.

$k=5$ a $k=10$: Observa-se uma melhoria na acurácia de teste.

$k>10$: À medida que k aumenta, a acurácia de treino vai aumentando ultrapassando a acurácia de treino em $k=20$. Esta diminuição da acurácia de treino e aumento lento da acurácia de teste pode sugerir uma leve tendência de underfitting, onde, com muitos vizinhos, o modelo perde a capacidade de capturar variações locais nos dados.

Pesos Baseados na Distância (Gráfico 2 à Direita):

$k=1$: O modelo novamente apresenta acurácia de treino perfeita, mas a acurácia de teste permanece baixa, semelhante ao caso anterior.

$k=5$ a $k=10$: Para valores intermediários de k , especialmente até de $k=10$, a acurácia de teste é semelhante ao modelo com pesos uniformes.

$k>10$: Com k maior que 10, a acurácia de teste estabiliza. Embora a acurácia de treino continue constante em 1.0, o impacto de atribuir pesos mais altos aos vizinhos mais próximos ajuda a mitigar a queda na acurácia de teste, mantendo um grau de localidade nas decisões.

Impacto no Poder de Generalização: Aumentar o número de vizinhos tende a suavizar o modelo, reduzindo o overfitting e melhorando a generalização até certo ponto. Contudo, quando k é excessivamente grande, pode resultar em underfitting, levando a uma queda na acurácia tanto de treino quanto de teste (neste caso não acontece redução da acurácia de treino).

O uso de pesos baseados na distância atenua o impacto negativo de um aumento excessivo de vizinhos, já que o modelo ainda prioriza os vizinhos mais próximos, preservando um nível de localidade nas suas decisões.

Resumindo, aumentar k ajuda a reduzir o overfitting, mas um valor muito elevado pode resultar em underfitting. O uso de pesos por distância geralmente melhora a generalização para valores maiores de k .

- 3) [1.5v] Considering the unique properties of the heart-disease.csv dataset, identify two possible difficulties of the naïve Bayes model used in the previous exercises when learning from the given dataset.

O Naïve Bayes supõe que todas as características (features) são independentes entre si, para uma dada classe alvo. Porém, no conjunto de dados de doenças cardíacas, é provável que muitas das características sejam correlacionadas. Por exemplo, variáveis como nível de colesterol, pressão arterial e idade podem ter dependências entre si, o que viola essa suposição de independência. Isto pode penalizar o desempenho do modelo Naïve Bayes, que não consegue modelar essas correlações entre os features.

Este modelo assume que os dados numéricos seguem uma distribuição normal (gaussiana). Se as features numéricas do conjunto de dados de doenças cardíacas não seguirem a distribuição normal (forem muito assimétricos, ou tiverem outras distribuições), o Naïve Bayes pode não se ajustar bem aos dados. Dados como níveis de colesterol ou idade podem não seguir uma distribuição normal, podendo diminuir a capacidade do modelo de fazer previsões precisas.

END