CS156 (Introduction to AI), Spring 2021

List all your references and sources here. This includes all sites/discussion boards/blogs/posts/etc. where you grabbed some code examples. https://realpython.com/train-test-split-python-data/

Homework 2 submission Roster Name: Neeval Kumar

Preferred Name (if different): Chosen Name Student ID: 011877086

Email address: kumar.neeval@gmail.com

Any special notes or anything you would like to communicate to me about this homework submission goes in here.

References and sources

In [74]: import pandas as pd import numpy as np

import matplotlib.pyplot as plt

Solution Load libraries and set random number generator seed

from sklearn.model selection import train test split

pred_results.append(knn(test_df.loc[index], train_df,3))

if newObservation.size + 1 != referenceData.iloc[0].size:

Distances will be a dictionary to hold the distances

raise Exception ('newObservation dimensions are not the same as reference Datas

and it's output variable, label counts for how many neighbors are in majority

#Loop through reference data and find distance between every observation in refere

Sort the distances dictionary based on key (which is the distance), and then re-

np.random.seed(42)

import sys

from sklearn.metrics import accuracy score import collections

import math as math

Beginning with creating our 2D data

In [76]: n = 100X1 = np.random.normal(loc = -2, scale = 2, size = int(n/2))

Y1 = np.random.normal(loc = 0, scale = 1, size = int(n/2))

X2 = np.random.normal(loc = 2, scale = 2, size = int(n/2))Y2 = np.random.normal(loc = 0, scale = 1, size = int(n/2))X = np.concatenate((X1, X2), axis=0)Y = np.concatenate((Y1, Y2), axis=0)

dt = pd.DataFrame({'X': X, 'Y':Y}, columns=['X', 'Y']) dt.head() Χ **0** -1.006572 0.324084 **1** -2.276529 -0.385082 **2** -0.704623 -0.676922

1.046060 0.611676

4 -2.468307 1.031000 Divide dataset into training/test In [77]: 11 = [0] * int(n/2)12 = [1] * int(n/2)labels = 11 + 12 print(labels)

X train, X test, Y train, Y test = train test split(dt, labels, test size=0.2, random print(f"Length of training is {len(X train)} \nLength of test is {len(X test)}") Length of training is 80 Length of test is 20 Driver code, set up dataframe for reference data, create prediction results and

append the result of test dataframe's classification from the training data In [78]: train_df = pd.DataFrame(X train) train df["Output Variable"] = Y train test_df = pd.DataFrame(X test) pred_results = [] for index in test df.index:

In [79]: print(pred results) print(Y test) [0, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0][0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0]ACTUAL FUNCTION - New observation will represent our new data record to be classified as either 0 or 1. Reference data is our dataset and k = 3 will be the nearest

neighbor count to the new observation. In [80]: def knn(newObservation, referenceData, k=3):

distances = {} label0 count = 0label1_count = 0 point2 = []

for x in newObservation: point2.append(x)

point1 = []

if i == k: break

else:

return 1

return 0

plt.subplot(1, 2, 1)

plt.subplot(1, 2, 2)

plt.tight_layout()

Accuracy Score is: 0.9

plt.show()

2

1

0

 $^{-1}$

-2

plt.title("Actual labels")

plt.title("Predicted labels")

Actual labels

else:

Store the result in distances for index in referenceData.index: row = referenceData.loc[index]

> for x in range(row.size-1): point1.append(row[x])

point1 = np.array(point1) point2 = np.array(point2)

for value in sorted (distances):

if output variable == 0: label0 count += 1

label1 count += 1

if label0 count < label1 count:</pre>

output_variable = distances[value]

distance = np.linalg.norm(point1 - point2)

distances[distance] = referenceData['Output Variable'][index]

Print Accuracy score between y test and prediction results and then plot

print(f'Accuracy Score is: {float(accuracy score(pred results, Y test))}')

plt.scatter(X train.iloc[:,0],X train.iloc[:,1], s=25, c=Y train, marker=".") plt.scatter(X test.iloc[:,0],X test.iloc[:,1], s=50, c=Y test, marker="v")

plt.scatter(X_train.iloc[:,0],X_train.iloc[:,1], s=25, c=Y_train, marker=".") plt.scatter(X_test.iloc[:,0], X_test.iloc[:,1], s=50, c=pred_results, marker="v")

Predicted labels

2

1

0

 $^{-1}$

-2

X1 = np.random.normal(loc = 0, scale = 3, size = int(n/4))Y1 = np.random.normal(loc = -3, scale = 1, size = int(n/4))Z1 = np.random.normal(loc = -1, scale = 1, size = int(n/4))

X2 = np.random.normal(loc = 0, scale = 3, size = int(n/4))Y2 = np.random.normal(loc = 1, scale = 2, size = int(n/4))Z2 = np.random.normal(loc = 1, scale = 1, size = int(n/4))

X3 = np.random.normal(loc = 0, scale = 3, size = int(n/4))Y3 = np.random.normal(loc = 3, scale = 1, size = int(n/4))Z3 = np.random.normal(loc = 4, scale = 1, size = int(n/4))

X4 = np.random.normal(loc = 0, scale = 3, size = int(n/4))Y4 = np.random.normal(loc = 5, scale = 3, size = int(n/4))Z4 = np.random.normal(loc = -3, scale = 1, size = int(n/4))

three dt = $pd.DataFrame(\{'X': X, 'Y':Y, 'Z': Z\}, columns=['X', 'Y', 'Z'])$

pred results.append(knn(test df.loc[index], train df,3))

X train, X test, Y train, Y test = train test split(three dt, labels, test size=0.2, 1

[1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0,

0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0,

1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0] 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1,

X = np.concatenate((X1, X2, X3, X4), axis=0)Y = np.concatenate((Y1, Y2, Y3, Y4), axis=0)Z = np.concatenate((Z1, Z2, Z3, Z4), axis=0)

11 = [0] * int(n/2)12 = [1] * int(n/2)labels = 11+12

pred results = []

print(pred results)

print(Y test)

train df = pd.DataFrame(X train) train df["Output Variable"] = Y train

test df = pd.DataFrame(X test)

for index in test df.index:

5

Now onto 3D data, same process

which are the labels. Count the number of labels and return accordingly

1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0] In [84]: print(f'Accuracy Score is: {float(accuracy score(pred results, Y test))}') plt.subplot(1, 2, 1)plt.scatter(X train.iloc[:,0],X train.iloc[:,1], s=25, c=Y train, marker=".") plt.scatter(X test.iloc[:,0],X test.iloc[:,1], s=50, c=Y test, marker="v") plt.title("Actual labels") plt.subplot(1, 2, 2) plt.scatter(X train.iloc[:,0],X train.iloc[:,1], s=25, c=Y train, marker=".") plt.scatter(X test.iloc[:,0],X test.iloc[:,1], s=50, c=pred results, marker="v") plt.title("Predicted labels") Accuracy Score is: 0.945 Out[84]: Text(0.5, 1.0, 'Predicted labels') Actual labels Predicted labels 15 15 10 10 5 5 0 0 -10