# CS156 (Introduction to AI), Spring 2021

## **Homework 2 submission**

Preferred Name (if different): Chosen Name

Roster Name: Neeval Kumar

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

References and sources 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/

#### import numpy as np import matplotlib.pyplot as plt import math as math

goes in here.

Solution

In [186... import pandas as pd

import sys

Out[188...

### from sklearn.metrics import accuracy score import collections from collections import Counter

```
In [187... np.random.seed(42)
Beginning with creating our 2D data
```

Load libraries and set random number generator seed

from sklearn.model selection import train test split

In [188... | n = 100

X1 = 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))

X = np.concatenate((X1, X2), axis=0)

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

Y = np.concatenate((Y1, Y2), axis=0)

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

dt.head()

Х Υ

**2** -0.704623 -0.676922 **3** 1.046060 0.611676

Divide dataset into training/test In [189...] 11 = [0]\*int(n/2)

**0** -1.006572 0.324084

**1** -2.276529 -0.385082

**4** -2.468307 1.031000

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

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

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

#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-

test\_df = pd.DataFrame(X test)

pred\_results = []

for index in test df.index: pred\_results.append(knn(test\_df.loc[index], train\_df,3))

In [191... 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.

def knn(newObservation, referenceData, k=3):

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

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

# Distances will be a dictionary to hold the distances # and it's output variable, label counts for how many neighbors are in majority

distances = {} point2 = []

for x in newObservation: point2.append(x)

point1 = []

**if** i == k: break

# Store the result in distances

row = referenceData.loc[index]

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

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

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

output variable = distances[value]

c.append(output\_variable)

**if** (val > (size/2)): return key

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

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

Predicted labels

2

1

0

 $^{-1}$ 

-2

Z2 = np.random.normal(loc = 1, scale = 1, size = int(n/4))

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

3, 1, 1,

1, 1,

In [199... | print(f'Accuracy Score is: {float(accuracy score(pred results, Y test))}')

3, 1, 1,

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

Ο, 3,

3, 3,

0.0

2, 0, 1, 2, 3, 3, 0, 3, 0, 2, 2,

3, 2, 1, 2, 2,

3, 1, 2, 1, 2, 0, 1, 2, 3, 3, 0, 3, 1, 2, 2, 3, 3, 0,

3, 2, 1, 2, 2, 2,

3, 3, 3, 1, 1, 3, 1, 2, 3, 2, 0, 3, 1, 1, 3, 1, 1,

3, 3,

1, 0,

2, 3, 2, 1, 0, 2,

3, 2,

3, 1, 2, 3, 2, 3, 3, 1, 1, 3, 1, 1,

3, 0,

2,

3, 0,

3, 1, 2, 1,

3, 3, 0, 3, 0,

3, 0,

3, 1, 1,

2,

plt.subplot(1, 2, 1)

1, 1, 2,

for index in referenceData.index:

c = []

for value in sorted (distances):

d = Counter(c) size = len(d)

for (key, val) in d.items():

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.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")

plt.tight\_layout()

plt.show() Accuracy Score is: 0.9

2 1

0

 $^{-1}$ 

-2

5

Now onto 3D data, same process

Actual labels

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))

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))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)

 $\label{three_dt} \mbox{three\_dt = pd.DataFrame} (\{'X': X, 'Y':Y, 'Z': Z\}, \mbox{columns=['X', 'Y', 'Z']})$ 11 = [0] \* int(n/4)12 = [1] \* int(n/4)

13 = [2] \* int(n/4)14 = [3] \* int(n/4)labels = 11+12+13+14X train, X test, Y train, Y test = train test split(three dt, labels, test size=0.2,

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:

print(pred results)

print(Y test) 3, 0, 3, 2, 2, 3, 1, 2, 0, 1, 3, 1, 0, 3, 2, 1, 1, 3, 3, 0, 1, 2, 1, 1, 1, 1, 3, 2, 2, 0, 1, 1, 2, 2, 1, 1, 1, 0, 1, 2, 1, 1, 2, 1, 2, 0, 3, 1, 1, 3, 2, 2, 1, 1, 1,

1, 1, 1, 0, 0, 2, 0, 1, 3, 1, 3, 0, 1, 3, 1, 3, 3, 1, 1, 3, 3, 2, 1, 3, 0, 3, 3, 3, 1, 3, 2, 2, 0, 2, 0] [3, 3, 1, 2, 2, 3, 0, 0, 1, 2, 1, 2, 3, 0, 2, 0, 3, 0, 0, 0, 0, 1, 1, 1, 3, 3, 0, 0, 3, 0, 3, 2, 2, 3, 1, 2, 3, 1, 3, 1, 0, 3, 2, 2, 1, 3, 3, 0, 1, 2, 1, 1, 1, 1, 1, 3, 1, 2, 2, 0, 0, 1, 0, 1, 2, 1, 2, 2, 1, 2, 0, 3, 1, 1, 3, 2, 2, 1, 1, 1, 3, 3, 2, 2, 1, 3, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 3, 1, 3, 0, 2, 3, 0, 3, 3, 1, 1, 3, 3, 2, 1, 3, 0, 3, 3, 1, 0, 3, 2, 2, 0, 2, 0]

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: 0.89 Out[199... Text(0.5, 1.0, 'Predicted labels') Actual labels Predicted labels 15.0 15.0 12.5 12.5 10.0 10.0 7.5 7.5 5.0 5.0 2.5

-2.52.5 -5.05.0-10

0.0