

```
from google.colab import drive
drive.mount('/content/drive')
```

☞ Drive already mounted at /content/drive; to attempt to forcibly remount, call dr:

```
from google.colab import files
from IPython.display import HTML, display
```

```
import numpy as np
import io
import re
from copy import deepcopy
```

```
# REQUIRED
testFileName = 'test_3.txt'
trainFileName = 'train_3.txt'
classAttributeIndex = 14
attributesIgnore = [2, 4, 10, 11]
```

```
# PARAMETERS
dataSplitRatio = 0
```

```
# Function to read a file
def readFile( fileName ):
    with open( fileName, 'r' ) as f:
        lines = f.read().split( '\n' )
        return lines
```

```
print("#### FILE DATA ####")
trainData = readFile( trainFileName )
testData = readFile( testFileName )
for line in testData:
    print( line )
```

☞

```
#### FILE DATA ####
```

```
66 Self-emp-not-inc 240294 Some-college 10 Married-civ-spouse Transport-moving Hu
40 Self-emp-not-inc 170214 Bachelors 13 Married-civ-spouse Transport-moving Husb
35 Private 194404 Bachelors 13 Married-civ-spouse Prof-specialty Husband White M
44 Private 96321 Some-college 10 Married-civ-spouse Prof-specialty Husband White
23 Private 202373 HS-grad 9 Never-married Sales Own-child Black Male 0 0 20 Unite
77 State-gov 267799 Doctorate 16 Married-spouse-absent Prof-specialty Not-in-fam
32 Private 64658 HS-grad 9 Divorced Craft-repair Not-in-family White Male 0 0 40
37 Private 248919 1st-4th 2 Married-civ-spouse Adm-clerical Husband White Male 0
67 Self-emp-inc 22313 Some-college 10 Married-civ-spouse Sales Husband White Male
22 Private 103805 Some-college 10 Never-married Other-service Own-child White Ma
49 Private 27802 Some-college 10 Married-civ-spouse Craft-repair Husband White M
39 Self-emp-not-inc 169542 Some-college 10 Married-civ-spouse Craft-repair Husba
43 Private 143809 HS-grad 9 Married-civ-spouse Tech-support Husband White Male 0
39 Self-emp-not-inc 33001 HS-grad 9 Married-civ-spouse Craft-repair Husband White
23 Private 194630 HS-grad 9 Separated Machine-op-inspct Own-child White Male 0 0
31 State-gov 203572 HS-grad 9 Divorced Adm-clerical Unmarried White Female 0 0 30
61 Private 95680 Some-college 10 Married-civ-spouse Machine-op-inspct Husband As
46 Private 212162 5th-6th 3 Married-civ-spouse Other-service Husband Black Male (
36 Private 115336 Assoc-voc 11 Married-civ-spouse Prof-specialty Wife White Fema
22 Private 316304 Some-college 10 Never-married Farming-fishing Own-child White I
```

```
# Converting the file data into a 2D array
```

```
def tabulateData( data, delimiter = ' ', hasHeader = True ):
```

```
    X = []
```

```
    for line in data:
```

```
        words = line.split(delimiter)
```

```
        X.append(words)
```

```
    return X
```

```
print("#### TABULATED DATA ####")
```

```
trainTabulatedData = tabulateData( trainData )
```

```
testTabulatedData = tabulateData( testData )
```

```
display(HTML(
```

```
    '<table><tr>{}</tr></table>'.format(
```

```
        '</tr><tr>'.join(
```

```
            '<td>{}</td>'.format('</td><td>'.join(str(_) for _ in row)) for row in tes
```

```
        )
```

```
))
```



#### TABULATED DATA ####

Self- 66 emp- not-inc	240294	Some- college	10	Married- civ- spouse	Transport- moving	Husband	Black	Male	0	0	60	United- States
Self- 40 emp- not-inc	170214	Bachelors	13	Married- civ- spouse	Transport- moving	Husband	White	Male	0	0	70	Iran
35 Private	194404	Bachelors	13	Married- civ- spouse	Prof- specialty	Husband	White	Male	0	0	40	United- States
44 Private	96321	Some- college	10	Married- civ- spouse	Prof- specialty	Husband	White	Male	0	0	40	United- States
23 Private	202373	HS-grad	9	Never- married	Sales	Own-child	Black	Male	0	0	20	United- States
77 State- gov	267799	Doctorate	16	Married- spouse- absent	Prof- specialty	Not-in- family	White	Male	0	0	4	United- States
32 Private	64658	HS-grad	9	Divorced	Craft- repair	Not-in- family	White	Male	0	0	40	United- States
37 Private	248919	1st-4th	2	Married- civ- spouse	Adm- clerical	Husband	White	Male	0	0	66	Mexico
Self- 67 emp- inc	22313	Some- college	10	Married- civ- spouse	Sales	Husband	White	Male	20051	0	40	United- States
22 Private	103805	Some- college	10	Never- married	Other- service	Own-child	White	Male	0	0	36	United- States
49 Private	27802	Some- college	10	Married- civ- spouse	Craft- repair	Husband	White	Male	0	1887	40	United- States
Self- 39 emp- not-inc	169542	Some- college	10	Married- civ- spouse	Craft- repair	Husband	White	Male	5178	0	40	United- States
43 Private	143809	HS-grad	9	Married- civ- spouse	Tech- support	Husband	White	Male	0	0	40	United- States
Self- 39 emp- not-inc	33001	HS-grad	9	Married- civ- spouse	Craft- repair	Husband	White	Male	0	0	40	United- States
23 Private	194630	HS-grad	9	Separated	Machine- op-inspct	Own-child	White	Male	0	0	53	United- States

```
# Removing data points which consists of null values
def preprocessData( tabulatedData, classAttributeIndex, train = True ):
    X = []
    Y_train = [ ]
    requiredLength = len( tabulatedData[0] )
    for dataPoint in tabulatedData:
        if( len(dataPoint) < requiredLength ):
            continue
        # if "none" in dataPoint:
```

```

# If none in dataframe.
# continue
X.append( dataPoint[ :requiredLength ] )

X = np.asanyarray(X)
if(train is True):
    Y_train = X[:, classAttributeIndex]
    X = np.delete(X, classAttributeIndex, axis = 1)
return X, Y_train

print("#### PREPROCESSED DATA ####")
X_train, Y_train = preprocessData( trainTabulatedData, classAttributeIndex = classAtt
X_test, Y_test = preprocessData( testTabulatedData, classAttributeIndex = classAttrib
print(X_train[10,:])
# print(Y_train)

[ ] #### PREPROCESSED DATA ####
['43' 'Private' '128354' 'HS-grad' '9' 'Married-civ-spouse' 'Adm-clerical'
 'Wife' 'White' 'Female' '0' '0' '30' 'United-States']

def categorical_distance(ptA, ptB):
    diff = ( ptA == ptB )
    return np.size(diff) - np.sum(diff)

def euclidean_distance(ptA, ptB):
    a = ptA.astype(np.float)
    b = ptB.astype(np.float)
    return np.sum((a - b)**2)**0.5

def distance(ptA, ptB, numeric_attributes, categorical_attributes):
    dist = 0
    dist += euclidean_distance(ptA[numeric_attributes], ptB[numeric_attributes])
    dist += categorical_distance(ptA[categorical_attributes], ptB[categorical_attribute
    return dist

def findAttributeTypes(X):
    N, M = np.shape(X)
    i = 0
    dataSet = X[0,:]
    while('? ' in dataSet):
        i += 1
        dataSet = X[i, :]
    categorical_attributes = []
    numeric_attributes = []
    array = dataSet

    for i in range(len(array)):
        regex_output = None

```

```

x = re.search('^[A-Za-z]+[-]*', array[i])
if x is not None:
    categorical_attributes.append(i)
    continue
x = re.search('^[0-9]+[.]*[0-9]+$', array[i])
if x is not None:
    numeric_attributes.append(i)
    continue
else:
    categorical_attributes.append(i)

return numeric_attributes, categorical_attributes

# Function to process Data that is removing the columns
def processData( data, removeColumns ):
    data = np.delete( data, removeColumns, axis = 1 )
    numeric_attributes, categorical_attributes = findAttributeTypes( data )
    return data, numeric_attributes, categorical_attributes

X_train, numeric_attributes, categorical_attributes = processData( X_train, attribute
X_test, numeric_attributes2, categorical_attributes2 = processData( X_test, attribute

print(numeric_attributes)
print(categorical_attributes)
print(X_train[0,:])

[0, 8]
[1, 2, 3, 4, 5, 6, 7, 9]
['42' 'Federal-gov' 'Bachelors' 'Never-married' 'Prof-specialty'
 'Not-in-family' 'White' 'Male' '40' '?']

def getMissedDataPoint( data ):
    N, M = np.shape(data)
    for i in range(N):
        point = data[i, :]
        if '?' in point:
            return i
    return -1

def getKNNNeighbours(X_train, Y_train, testPoint, k, numeric_attributes, categorical_a
    dist = np.empty((1,3))

    # Finding distance with all the training nodes and storing in dist matrix
    N_train, M = np.shape(X_train)

    for j in range( N_train ):
        trainPoint = X_train[j, :]
        if not '?' in trainPoint:
            temp = np.array([[ j, distance( trainPoint, testPoint, numeric_attributes, ca
            dist = np.append( dist, temp, axis = 0 )
        dist = np.delete(dist, 0, axis = 0)

```

```

# Sorting the distances
dist = dist[dist[:, 1].argsort()]

# Selecting top K elements as our neighbours
neighbours = dist[:k, :]

return neighbours

def fillTrainSet( data, Y, numeric_attributes, categorical_attributes ):
    N, M = np.shape( data )

    # i = getMissedDataPoint( data )
    for i in range(10):
        myPoint = data[i, :]
        missingAttributesIndex = []
        print("-----")
        print("Data point with missing value:")
        print(data[i, :])
        for j in range( len(myPoint) ):
            if '?' in myPoint[j]:
                missingAttributesIndex.append(j)
        copy_myPoint = deepcopy(myPoint)
        copy_data = deepcopy(data)

        copy_myPoint = np.delete(copy_myPoint, missingAttributesIndex )
        copy_data = np.delete(copy_data, missingAttributesIndex, axis = 1)

        numeric_attributes, categorical_attributes = findAttributeTypes( copy_data )

        myNeighbour = getKNNNeighbours( copy_data[10:,:], Y, copy_myPoint, 1, numeric_attr

        for j in range( len(myPoint) ):
            if '?' in myPoint[j]:
                # Adding 10 as we started with 10 training instances
                data[i, j] = data[ int(myNeighbour[0][0]) + 10, j ]
        print("My nearest equivalent:")
        print((data[int(myNeighbour[0][0])+10, :]) )
        print("My updated value:")
        print(data[i, :])
    return data

X_train = fillTrainSet(X_train, Y_train, numeric_attributes, categorical_attributes )

```



-----  
Data point with missing value:

```
['42' 'Federal-gov' 'Bachelors' 'Never-married' 'Prof-specialty'
 'Not-in-family' 'White' 'Male' '40' '?']
```

My nearest equivalent:

```
['37' 'Local-gov' 'Assoc-acdm' 'Married-civ-spouse' 'Craft-repair'
 'Husband' 'White' 'Male' '40' 'United-States']
```

My updated value:

```
['42' 'Federal-gov' 'Bachelors' 'Never-married' 'Prof-specialty'
 'Not-in-family' 'White' 'Male' '40' 'United-States']
```

-----  
Data point with missing value:

```
['30' '?' '11th' 'Never-married' '?' 'Unmarried' 'Black' 'Female' '40'
 'United-States']
```

My nearest equivalent:

```
['25' 'Private' 'HS-grad' 'Married-civ-spouse' 'Craft-repair' 'Husband'
 'White' 'Male' '40' 'United-States']
```

My updated value:

```
['30' 'Private' '11th' 'Never-married' 'Craft-repair' 'Unmarried' 'Black'
 'Female' '40' 'United-States']
```

-----  
Data point with missing value:

```
['83' '?' '7th-8th' 'Married-civ-spouse' '?' 'Husband' 'White' 'Male' '20'
 'United-States']
```

My nearest equivalent:

```
['76' 'Local-gov' '10th' 'Divorced' 'Transport-moving' 'Not-in-family'
 'Black' 'Male' '20' 'United-States']
```

My updated value:

```
['83' 'Local-gov' '7th-8th' 'Married-civ-spouse' 'Transport-moving'
 'Husband' 'White' 'Male' '20' 'United-States']
```

-----  
Data point with missing value:

```
['24' 'Private' 'Some-college' 'Married-civ-spouse' 'Transport-moving'
 'Husband' 'White' 'Male' '40' '?']
```

My nearest equivalent:

```
['27' 'Local-gov' 'HS-grad' 'Married-spouse-absent' 'Sales' 'Own-child'
 'Black' 'Female' '40' 'United-States']
```

My updated value:

```
['24' 'Private' 'Some-college' 'Married-civ-spouse' 'Transport-moving'
 'Husband' 'White' 'Male' '40' 'United-States']
```

-----  
Data point with missing value:

```
['60' '?' '7th-8th' 'Married-civ-spouse' '?' 'Husband' 'White' 'Male' '40'
 'United-States']
```

My nearest equivalent:

```
['60' 'Private' 'Bachelors' 'Married-civ-spouse' 'Exec-managerial'
 'Husband' 'White' 'Male' '40' 'United-States']
```

My updated value:

```
['60' 'Private' '7th-8th' 'Married-civ-spouse' 'Exec-managerial' 'Husband'
 'White' 'Male' '40' 'United-States']
```

-----  
Data point with missing value:

```
['22' '?' 'HS-grad' 'Never-married' '?' 'Own-child' 'White' 'Male' '40'
 'United-States']
```

My nearest equivalent:

```
['21' 'Private' 'HS-grad' 'Never-married' 'Other-service' 'Own-child'
 'White' 'Male' '40' 'United-States']
```

```

My updated value:
['22' 'Private' 'HS-grad' 'Never-married' 'Other-service' 'Own-child'
 'White' 'Male' '40' 'United-States']
-----
Data point with missing value:
['23' 'Private' 'HS-grad' 'Never-married' 'Sales' 'Not-in-family' 'Black'
 'Female' '27' '?']
My nearest equivalent:
['23' 'Private' 'Bachelors' 'Never-married' 'Prof-specialty'
 'Not-in-family' 'White' 'Female' '20' 'United-States']
My updated value:
['23' 'Private' 'HS-grad' 'Never-married' 'Sales' 'Not-in-family' 'Black'
 'Female' '27' 'United-States']
-----
Data point with missing value:
['19' '?' 'Some-college' 'Never-married' '?' 'Not-in-family' 'White'
 'Female' '30' 'United-States']
My nearest equivalent:
['19' 'Private' 'Some-college' 'Never-married' 'Handlers-cleaners'
 'Not-in-family' 'White' 'Male' '30' 'United-States']
My updated value:
['19' 'Private' 'Some-college' 'Never-married' 'Handlers-cleaners'
 'Not-in-family' 'White' 'Female' '30' 'United-States']
-----
Data point with missing value:
['51' 'Self-emp-not-inc' 'Masters' 'Married-civ-spouse' 'Prof-specialty'
 'Wife' 'White' 'Female' '55' '?']
My nearest equivalent:
['56' 'Private' '11th' 'Married-civ-spouse' 'Sales' 'Husband' 'White'
 'Male' '55' 'United-States']
My updated value:
['51' 'Self-emp-not-inc' 'Masters' 'Married-civ-spouse' 'Prof-specialty'
 'Wife' 'White' 'Female' '55' 'United-States']
-----
Data point with missing value:
['21' 'Private' 'HS-grad' 'Never-married' 'Transport-moving'
 'Other-relative' 'White' 'Male' '60' '?']
My nearest equivalent:
['28' 'Private' 'Assoc-voc' 'Never-married' 'Exec-managerial'
 'Not-in-family' 'White' 'Male' '60' 'United-States']
My updated value:
['21' 'Private' 'HS-grad' 'Never-married' 'Transport-moving'
 'Other-relative' 'White' 'Male' '60' 'United-States']

```

```
numeric_attributes, categorical_attributes = findAttributeTypes(X_train[10:, :])
```

```

def predict( neighbours ):
    # Finding count class
    countClass = dict()
    for neighbourClass in neighbours[:, 2]:
        if( neighbourClass in countClass ):
            countClass[neighbourClass] += 1
        else:
            countClass[neighbourClass] = 1

```

```
# Finding prediction
```



```
maxCount = -1
predClass = None
for countClassKey in countClass:
    if( maxCount < countClass[ countClassKey ] ):
        maxCount = countClass[ countClassKey ]
        predClass = countClassKey

return predClass

N_test, M = np.shape(X_test)
for i in range(N_test):
    testPoint = X_test[i]
    numeric_attributes, categorical_attributes = findAttributeTypes(X_train)
    neighbours = getKNNNeighbours(X_train, Y_train, testPoint, 3 , numeric_attributes, c
    print(testPoint)
    print("Predicted Class: " + str(predict(neighbours)))
    print("-----")
```



```

['66' 'Self-emp-not-inc' 'Some-college' 'Married-civ-spouse'
 'Transport-moving' 'Husband' 'Black' 'Male' '60' 'United-States']
Predicted Class: <=50K
-----
['40' 'Self-emp-not-inc' 'Bachelors' 'Married-civ-spouse'
 'Transport-moving' 'Husband' 'White' 'Male' '70' 'Iran']
Predicted Class: >50K
-----
['35' 'Private' 'Bachelors' 'Married-civ-spouse' 'Prof-specialty'
 'Husband' 'White' 'Male' '40' 'United-States']
Predicted Class: <=50K
-----
['44' 'Private' 'Some-college' 'Married-civ-spouse' 'Prof-specialty'
 'Husband' 'White' 'Male' '40' 'United-States']
Predicted Class: >50K
-----
['23' 'Private' 'HS-grad' 'Never-married' 'Sales' 'Own-child' 'Black'
 'Male' '20' 'United-States']
Predicted Class: <=50K
-----
['77' 'State-gov' 'Doctorate' 'Married-spouse-absent' 'Prof-specialty'
 'Not-in-family' 'White' 'Male' '4' 'United-States']
Predicted Class: >50K
-----
['32' 'Private' 'HS-grad' 'Divorced' 'Craft-repair' 'Not-in-family'
 'White' 'Male' '40' 'United-States']
Predicted Class: <=50K
-----
['37' 'Private' '1st-4th' 'Married-civ-spouse' 'Adm-clerical' 'Husband'
 'White' 'Male' '66' 'Mexico']
Predicted Class: >50K
-----
['67' 'Self-emp-inc' 'Some-college' 'Married-civ-spouse' 'Sales' 'Husband'
 'White' 'Male' '40' 'United-States']
Predicted Class: >50K
-----
['22' 'Private' 'Some-college' 'Never-married' 'Other-service' 'Own-child'
 'White' 'Male' '36' 'United-States']
Predicted Class: <=50K
-----
['49' 'Private' 'Some-college' 'Married-civ-spouse' 'Craft-repair'
 'Husband' 'White' 'Male' '40' 'United-States']
Predicted Class: <=50K
-----
['39' 'Self-emp-not-inc' 'Some-college' 'Married-civ-spouse'
 'Craft-repair' 'Husband' 'White' 'Male' '40' 'United-States']
Predicted Class: <=50K
-----
['43' 'Private' 'HS-grad' 'Married-civ-spouse' 'Tech-support' 'Husband'
 'White' 'Male' '40' 'United-States']
Predicted Class: <=50K
-----
['39' 'Self-emp-not-inc' 'HS-grad' 'Married-civ-spouse' 'Craft-repair'
 'Husband' 'White' 'Male' '40' 'United-States']
Predicted Class: >50K
-----
['23' 'Private' 'HS-grad' 'Separated' 'Machine-op-inspct' 'Own-child'

```