

# Module 2 Homework

ISE-529 Predictive Analytics

Solutions

*Note: In answering the following questions, you may only use functionality from Base Python, NumPy, Pandas, or Seaborn*

## 1. Evaluating regression functions

1A) The file "HW 2 Problem 1 Data.xlsx" contains two datasets - training and test in two different spreadsheet tabs. Read the tables into two dataframes (training\_data and test\_data) and display the first 10 rows of each dataframe. Hint - look up the Pandas function for reading in Excel files.

```
In [11]: import numpy as np
import pandas as pd
test_data = pd.read_excel('HW 2 Problem 1 Data.xlsx', sheet_name = 'Test Data')
training_data = pd.read_excel('HW 2 Problem 1 Data.xlsx', sheet_name = 'Training Data')
training_data.head(10)
```

```
Out[11]:
```

	X	Y
0	61	17661.067682
1	87	15482.455058
2	38	17444.767982
3	6	-4270.225550
4	54	8075.733045
5	29	16820.129406
6	77	23921.367914
7	70	16541.631267
8	21	4425.240897
9	76	30681.044509

```
In [12]: test_data.head(10)
```

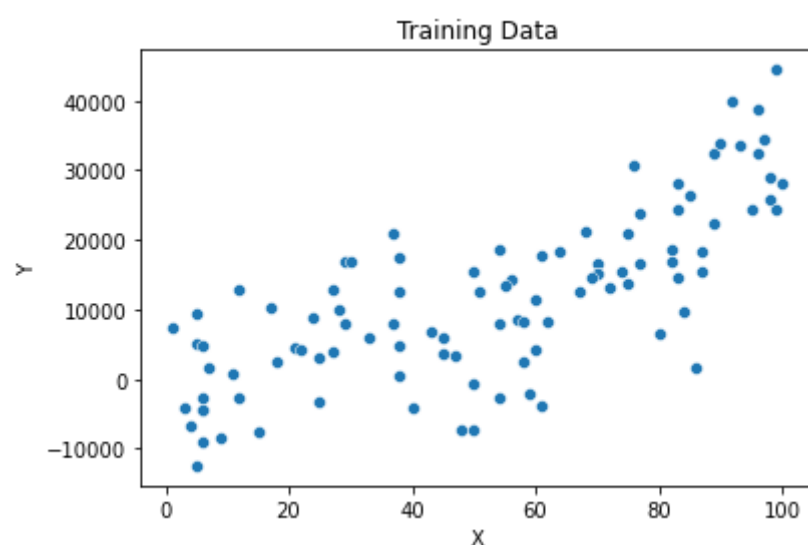
```
Out[12]:
```

	X	Y
0	36	15986.579536
1	64	-2761.487999
2	66	9752.039830
3	88	20038.697197
4	4	-13421.192312
5	80	6644.121448
6	22	-10124.906855
7	85	21502.561217
8	63	12105.189261
9	8	-12161.603294

1b) Using Seaborn, create scatterplots of the two dataframes. For full credit, include a title ("Training Data" and "Test Data") for the two scatterplots.

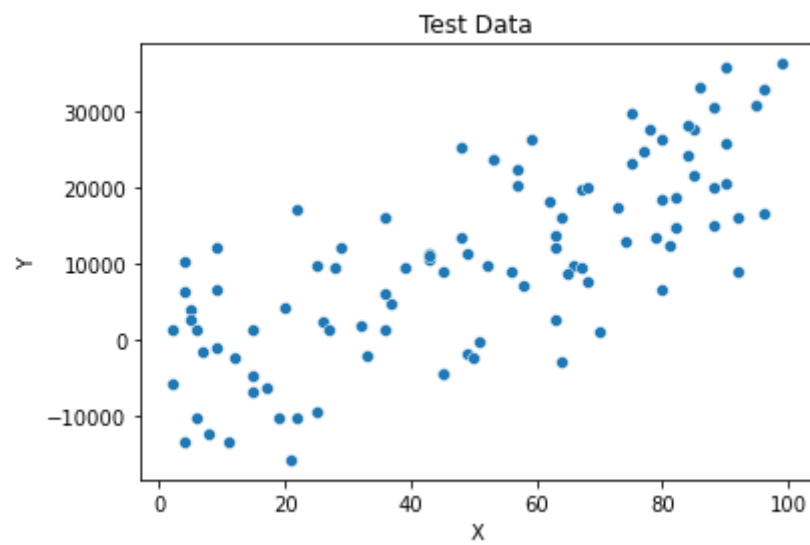
```
In [17]: import seaborn as sns
sns.scatterplot(data = training_data, x = 'X', y='Y').set(title = "Training Data")
```

```
Out[17]: [Text(0.5, 1.0, 'Training Data')]
```



```
In [18]: sns.scatterplot(data = test_data, x = 'X', y='Y').set(title = "Test Data")
```

```
Out[18]: [Text(0.5, 1.0, 'Test Data')]
```



1c) Now, we are going to calculate the test and training MSEs for three different candidate regression models:

- F1:  $300x - 4280$
- F2:  $3.37x^2 - 36.8x + 1220$
- F3:  $0.1x^3 - 12x^2 + 590x - 3600$

Create three functions for these three candidate models - f1(), f2(), and f3(). They should take as an input  $x$  and return the predicted value of  $y$  corresponding to that input:

```
In [34]: def f1(x):  
         return 300*x - 4280
```

```
In [64]: def f2(x):  
         return 3.37*x**2 - 36.8*x + 1220
```

```
In [58]: def f3(x):  
         return 0.1*x**3 - 12*x**2 + 590*x - 3600
```

1d) Now, write a function calc\_mse() that takes three parameter inputs:

- $x$  - an array (or Pandas series) of predictor  $X$  values
- $y$  - an array (or Pandas series) of response  $Y$  values
- $f$  - a function to be called for each value of the  $x$  and  $y$  arrays

The function should calculate the MSE for the function  $f$  using the  $x$  and  $y$  data arrays

```
In [49]: def calc_mse(x,y,f):  
         total_se = 0  
         for i in range(x.size):  
             total_se = total_se + (f(x[i])-y[i])**2  
         return total_se/x.size
```

1e) Call this calc\_mse function six times to calculate the training and test MSE for each of the three models:

```
In [94]: calc_mse(training_data['X'], training_data['Y'], f1)
```

```
Out[94]: 65116445.47500964
```

```
In [69]: calc_mse(training_data['X'], training_data['Y'], f2)
```

```
Out[69]: 57922250.15746113
```

```
In [70]: calc_mse(training_data['X'], training_data['Y'], f3)
```

```
Out[70]: 54103936.344325066
```

```
In [65]: calc_mse(test_data['X'], test_data['Y'], f1)
```

```
Out[65]: 67828868.62264524
```

```
In [66]: calc_mse(test_data['X'], test_data['Y'], f2)
```

```
Out[66]: 67805434.17472021
```

```
In [67]:
```

```
calc_mse(test_data['X'], test_data['Y'], f3)
```

Out[67]: 71191752.34205785

1f) Which of the three models would you select for use and why?

Model F2 because it has the lowest MSE on the test data.

1g) Instead of writing functions, write a single line of Python code to calculate the test MSE for function F1. Hint: use the Python map and lambda functions.

```
In [95]: (list(map(lambda x:300*x - 4280, test_data['X'])) - test_data['Y']).pow(2).mean()
```

Out[95]: 67828868.62264524

## 2. KNN and Calculate Misclassification Rates

2a) Read the file "HW 2 Problem 2 Data.csv" into a dataframe called knn\_data and display its first 10 rows

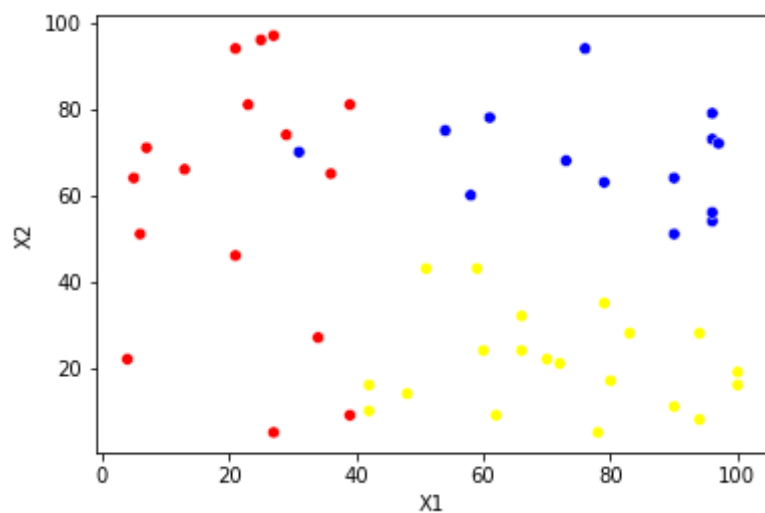
```
In [407... knn_data = pd.read_csv('HW 2 Problem 2 Data.csv')
knn_data.head(10)
```

```
Out[407...   X1  X2  Category
0   59  43    Yellow
1   60  24    Yellow
2   54  75     Blue
3   27   5     Red
4   96  73     Blue
5   51  43    Yellow
6   21  94     Red
7  100  19    Yellow
8   73  68     Blue
9   62   9    Yellow
```

2b) Using Seaborn, create a color-coded scatterplot of the data (where each point is colored with its category color). For full credit, be sure to color the points correctly with yellow, blue, and red colors.

```
In [408... sns.scatterplot(data = knn_data, x = 'X1', y = 'X2', hue = "Category", legend = False,
                      palette=dict(Yellow = "yellow", Red = 'red', Blue = 'blue'))
```

Out[408... <AxesSubplot:xlabel='X1', ylabel='X2'>



Now, we are going to create and assess predictions for the category of each observation using KNN-1 and KNN-3 algorithms. This means for each observation we are going to determine a prediction for its category color as if we didn't know what it was and then we will be checking to see if the prediction matches its actual category color.

The problem sub-parts below will step you through the process of doing this.

2c) Create a two-dimensional numpy array of size 50x50 that has the distance of each observation to every other observation. Use a Euclidean formula to calculate the distances. Display the first row of the distance table (this row will show the distance of each observation from the first observation)

```
In [410... num_obs = knn_data.shape[0]
distance_table = np.zeros((num_obs,num_obs))
X1 = knn_data['X1']
X2 = knn_data['X2']
for i in range (num_obs):
    for j in range (num_obs):
        distance_table[i,j] = np.sqrt((X1[i] - X1[j])**2 + (X2[i] - X2[j])**2)

distance_table[0,]
```

```
Out[410... array([ 0.          , 19.02629759, 32.38826948, 49.67896939, 47.63402146,
        8.          , 63.60031446, 47.50789408, 28.65309756, 34.13209633,
        32.01562119, 62.76941931, 57.93962375, 17.02938637, 38.07886553,
        33.42154993, 37.44329045, 39.44616585, 51.623638   , 28.3019434  ,
        42.48529157, 58.87274412, 38.60051813, 28.28427125, 31.82766093,
        43.13930922, 21.54065923, 44.55333882, 23.70653918, 38.11823711,
        31.01612484, 59.05929224, 38.89730068, 35.05709629, 29.68164416,
        13.03840481, 20.24845673, 42.94182111, 49.09175083, 39.2173431  ,
        51.42956348, 53.60037313, 62.96824597, 52.34500931, 31.90611227,
        25.55386468, 49.49747468, 37.12142239, 47.80167361, 53.75872022])
```

2d) Now create a dataframe with 50 rows (one per observation) and 3 columns labeled 'nn1\_cat', 'nn2\_cat', and 'nn3\_cat'. Populate the dataframe with the category color for the first, second, and third nearest neighbor for each observation.

Display the first ten rows of this dataframe

```
In [412... nn_cats = pd.DataFrame(np.empty((num_obs, 3), dtype = 'str'))
nn_cats.columns = ['nn1_cat', 'nn2_cat', 'nn3_cat']
for i in range(num_obs):
    nn_cats['nn1_cat'][i] = knn_data.iloc[np.argsort(distance_table[i,])[1]]['Category']
    nn_cats['nn2_cat'][i] = knn_data.iloc[np.argsort(distance_table[i,])[2]]['Category']
    nn_cats['nn3_cat'][i] = knn_data.iloc[np.argsort(distance_table[i,])[3]]['Category']

nn_cats.head(10)
```

```
Out[412...
```

	nn1_cat	nn2_cat	nn3_cat
0	Yellow	Yellow	Blue
1	Yellow	Yellow	Yellow
2	Blue	Blue	Red
3	Red	Yellow	Yellow
4	Blue	Blue	Blue
5	Yellow	Blue	Yellow
6	Red	Red	Red
7	Yellow	Yellow	Yellow
8	Blue	Blue	Blue
9	Yellow	Yellow	Yellow

2e) Create a Pandas series nn1\_preds that has the prediction for each observation using a KNN1 algorithm. Display the first 10 rows of the series.

```
In [413... nn1_preds = nn_cats['nn1_cat']
nn1_preds.head(10)
```

```
Out[413...
```

0	Yellow
1	Yellow
2	Blue
3	Red
4	Blue
5	Yellow
6	Red
7	Yellow
8	Blue
9	Yellow

Name: nn1\_cat, dtype: object

2f) Calculate the misclassification rate for the KNN1 algorithm on this dataset

```
In [402... nn1_misclass_rate = (num_obs - sum(nn1_preds == knn_data['Category']))/nn1_preds.size
nn1_misclass_rate
```

Out[402... 0.12

2g) Create a Pandas series nn3\_preds that has the prediction for each observation using a KNN3 algorithm. Display the first 10 rows of the series. (If the three nearest neighbors all have different color categories, use the first nearest neighbor category as the prediction)

```
In [414... nn3_preds = pd.Series(np.empty(num_obs), dtype = 'str')
for i in range(num_obs):
    if sum(nn_cats.iloc[i] == "Yellow") >= 2:
        nn3_preds[i] = 'Yellow'
    elif sum(nn_cats.iloc[i] == "Blue") >= 2:
        nn3_preds[i] = 'Blue'
    elif sum(nn_cats.iloc[i] == "Red") >= 2:
        nn3_preds[i] = 'Red'
    else:
        nn3_preds[i] = nn_cats["Category"][i]
```

```
In [404... nn3_preds.head(10)
```

```
Out[404...
```

0	Yellow
1	Yellow
2	Blue

```
3    Yellow
4     Blue
5    Yellow
6      Red
7    Yellow
8     Blue
9    Yellow
dtype: object
```

2h) Calculate the misclassification rate for the KNN3 algorithm on this dataset

In [415...

```
nn3_misclass_rate = (num_obs - sum(nn3_preds == knn_data['Category']))/nn3_preds.size
nn3_misclass_rate
```

Out[415... 0.08