

## Project 1 - Part3. k-NN

```
In [1]: 1 # Load Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 from sklearn import datasets
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import train_test_split
9 from sklearn.pipeline import Pipeline, FeatureUnion
10 from sklearn.model_selection import GridSearchCV
11
12 from sklearn import preprocessing
13 from sklearn.metrics import accuracy_score
14 from sklearn.neighbors import NearestNeighbors, KNeighborsClassifier
15
16 from dmbs import regressionSummary
```

```
In [2]: 1 # Load data
2
3 df = pd.read_csv('FinalNew_house.csv')
4 df.columns
```

```
Out[2]: Index(['price', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floor
s',
               'waterfront', 'view', 'condition', 'grade', 'sqft_above',
               'sqft_basement', 'zipcode', 'lat', 'long', 'sqft_living15',
               'sqft_lot15', 'age', 'renov_age'],
              dtype='object')
```

```
In [3]: 1 # Columns to be used were identified on previous section.
2
3 X = df[['grade', 'lat', 'age', 'sqft_living15']]
4 y = df['price']
```

```
In [4]: 1 # Split dataset into train and test
2
3 train_X, valid_X, train_y, valid_y = train_test_split(X,y, test_size=0.4,
```

```
In [5]: 1 # Transform pandas dataframe to numpy array.
2
3 train_X_array = train_X.to_numpy()
4 valid_X_array = valid_X.to_numpy()
5 train_y_array = train_y.to_numpy()
6 valid_y_array = valid_y.to_numpy()
```

```
In [6]: 1 # Create standardizer
2 standardizer = StandardScaler()
```

```
In [7]: 1 # Standardize features
        2 train_X_std = standardizer.fit_transform(train_X_array)
        3 valid_X_std = standardizer.fit_transform(valid_X_array)
```

```
In [8]: 1 # Train a KNN classifier with 5 neighbors
        2 knn5 = KNeighborsClassifier(n_neighbors=5, n_jobs=-1).fit(train_X_std, t
        3 knn5
```

```
Out[8]: KNeighborsClassifier
        KNeighborsClassifier(n_jobs=-1)
```

```
In [9]: 1 # Predict the class of two observations
        2 y_predict_5 = knn5.predict(valid_X_std)
        3 y_predict_5
```

```
Out[9]: array([257000., 130000., 315450., ..., 235000., 445000., 185000.])
```

## Identifying the Best Neighborhood Size

```
In [10]: 1 # Load Library
         2
         3 from sklearn.pipeline import Pipeline, FeatureUnion
         4 from sklearn.model_selection import GridSearchCV
```

```
In [11]: 1 # Create a KNN classifier
         2 knn = KNeighborsClassifier(n_neighbors=5, n_jobs=-1)
```

```
In [12]: 1 # Create a pipeline
         2 pipe = Pipeline([("standardizer", standardizer), ("knn", knn)])
```

```
In [13]: 1 # Create space of candidate values
         2 search_space = [{"knn__n_neighbors": [ 1,2,3, 4, 5, 6, 7, 8, 9, 10,11,12,
```

```
In [14]: 1 # Create grid search
         2 classifier = GridSearchCV(pipe, search_space, cv=5, verbose=0).fit(train
```

C:\Users\daarv\anaconda3\lib\site-packages\sklearn\model\_selection\\_split.py:684: UserWarning: The least populated class in y has only 1 members, which is less than n\_splits=5.  
warnings.warn(

```
In [15]: 1 # Best neighborhood size (k)
         2 best_K = classifier.best_estimator_.get_params()["knn__n_neighbors"]
         3 print('The best K for this K-NN model is: ', best_K)
```

The best K for this K-NN model is: 2

## Model Statistics with k = 5

```
In [16]: 1 regressionSummary(valid_y_array, y_predict_5)
```

Regression statistics

```

                Mean Error (ME) : 84916.5426
      Root Mean Squared Error (RMSE) : 141722.7023
        Mean Absolute Error (MAE) : 102478.5349
        Mean Percentage Error (MPE) : 15.7483
Mean Absolute Percentage Error (MAPE) : 21.5093

```

## Model Statistics with k = 2

```
In [17]: 1 knn = KNeighborsClassifier(n_neighbors=2, n_jobs=-1).fit(train_X_std, train_y)
        2
        3 y_predicted = knn.predict(valid_X_std)
```

```
In [18]: 1 regressionSummary(valid_y_array, y_predicted)
```

Regression statistics

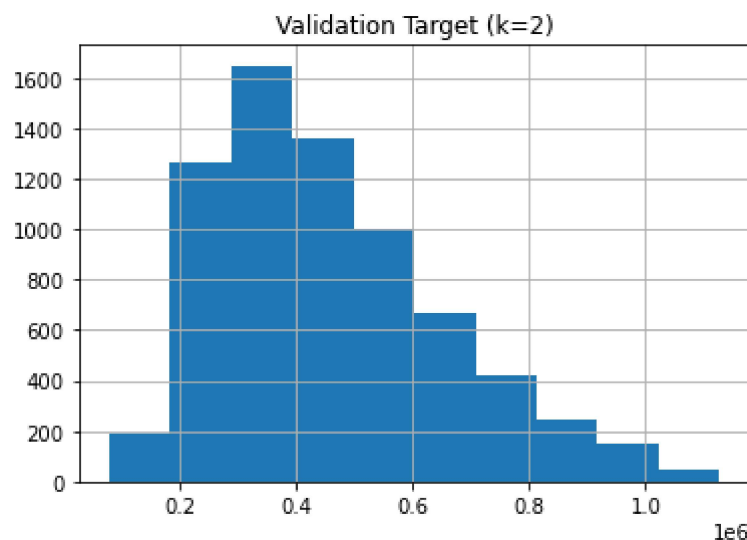
```

                Mean Error (ME) : 43693.3936
      Root Mean Squared Error (RMSE) : 123017.5274
        Mean Absolute Error (MAE) : 85991.5124
        Mean Percentage Error (MPE) : 6.3275
Mean Absolute Percentage Error (MAPE) : 18.7976

```

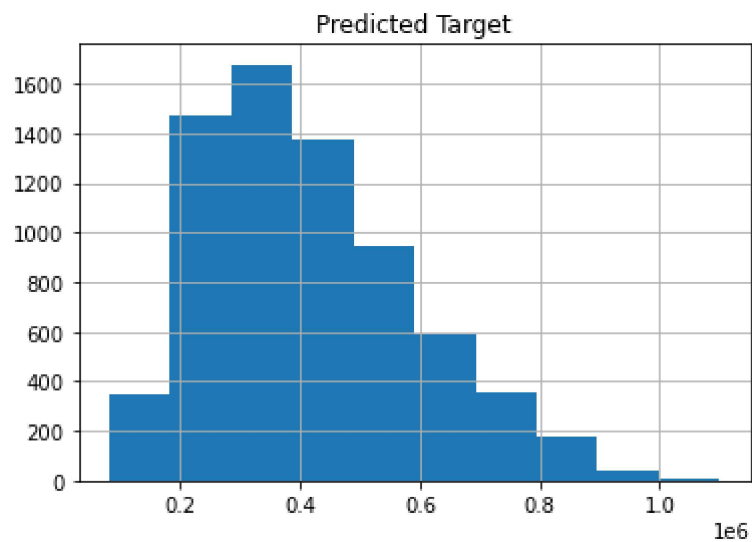
```
In [19]: 1 pd.DataFrame(valid_y_array).hist()
        2 plt.title('Validation Target (k=2)')
```

```
Out[19]: Text(0.5, 1.0, 'Validation Target (k=2)')
```



```
In [20]: 1 pd.DataFrame(y_predicted).hist()  
        2 plt.title('Predicted Target')
```

Out[20]: Text(0.5, 1.0, 'Predicted Target')



```

In [21]: ▶ 1 #fig, ((ax0, ax1), (ax2, ax3)) = plt.subplots(nrows=1, ncols=3)
2 fig, (ax1,ax2, ax3) = plt.subplots(nrows=1, ncols=3)
3
4 ax1.hist(y_predict_5, density=True, histtype='bar',color= 'red', stacked:
5 ax1.set_title('Predicted Target (k=5)')
6
7 ax2.hist(valid_y_array, histtype='bar', color= 'blue')
8 ax2.set_title('Validation Target')
9
10 ax3.hist(y_predicted, histtype='bar',color= 'tan')
11 ax3.set_title('Predicted Target (k=2)')
12
13 fig.tight_layout()
14 plt.show()

```

