

Introduction:

This project performs the K-Nearest-Neighbor Regression on a 100 generated data samples and obtains the closest neighbor information for various K values. The project is divided into 7 cases.

- Case1: K=1, the K neighbors contribute equally
- Case2: K=3, the K neighbors contribute equally
- Case3: K=5, the K neighbors contribute equally
- Case4: K=1, each of the K neighbors has an influence that is inversely proportional to the distance from the point
- Case5: K=3, each of the K neighbors has an influence that is inversely proportional to the distance from the point
- Case6: K=5, each of the K neighbors has an influence that is inversely proportional to the distance from the point
- Case7: K=N, each of the K neighbors has an influence that is inversely proportional to the distance from the point

Steps involved:

- Step1: Data Preparation of X
- Step2: Data Preparation of y
- Step3: KNN Regression modeling for each of the above seven cases
- Step4: Obtain predictions for the test set for each of the above seven cases
- Step5: Plots for each of the above seven cases

Step1: Data Preparation of X

Consider a single dimension (variable). Obtain N = 100 id samples x1, x2, ... of N uniformly randomly between 1 and 10, and then obtain the corresponding y values as the natural logarithm of x plus a Gaussian noise (mean 0, standard deviation 0.1), with different points having different amounts of noise.

```
In [3]: import pandas as pd
import numpy as np

In [30]: np.random.seed(44)
N = 100
x = np.random.uniform(low=1, high=10, size=N)
X

Out [30]: array([[8.51357934, 1.84316494, 7.76176434, 4.24456763, 4.23379754,
0.48314543, 4.54401596, 4.68165249, 9.63802189, 7.39133394,
9.64473693, 5.19958998, 4.84886369, 2.82117331, 2.96108982,
5.12321386, 9.49015648, 6.86188453, 8.81769696, 2.92442334,
6.73148812, 2.25231627, 5.12833663, 8.86476867, 3.32605275,
0.98595969, 8.76486902, 2.33963257, 0.80654962, 2.4329737,
5.9569454, 1.8362965, 8.8543834, 5.99676367, 8.35117425,
9.91749413, 8.24015483, 4.39673885, 5.64216363, 1.53809629,
7.39861118, 1.62527632, 8.94304624, 7.53452037, 9.58059591,
9.73027615, 7.2768926, 0.37892523, 8.9365474, 1.85474785,
5.10840179, 5.4345151, 1.98074554, 7.53452037, 9.58059591,
3.44484138, 9.67436608, 2.47811181, 2.19212465, 3.85618862,
3.76880383, 4.78948608, 3.97189752, 6.11267836, 1.85756195,
9.7819739, 3.41790753, 9.26712422, 9.13059582, 1.1427292,
8.68614021, 7.88945556, 1.33912263, 3.7890309, 6.07482543,
4.78486987, 7.12772067, 6.69188453, 5.48332341, 7.63407523,
7.70369599, 1.07997036, 5.44271745, 8.80075932, 1.16429597,
9.16213249, 2.0451139, 6.97967893, 1.43687471, 5.45486395,
5.52780892, 7.88264422, 4.12731829, 9.38277039, 3.69282354,
5.96915355, 1.51797934, 5.15767122, 6.78749758, 1.87174229])
```

Step2: Data Preparation of y

Gaussian noise (mean 0, standard deviation 0.1) with different points having different amounts of noise. Note: Gaussian is nothing but normal distribution, therefore using numpy.random.normal gives 100 noises (size=100)

```
In [7]: np.random.seed(44)
noise_list = np.random.normal(loc=0, scale=0.1, size=100)
noise_list

Out [7]: array([-7.58614717e-02, 1.31635732e-01, 1.24614800e-01, -1.68491574e-01,
-1.46814368e-01, -1.15978406e-01, 1.85873696e-01, 0.75877636e-01,
5.23221394e-01, 5.55471537e-02, 9.63403694e-02, 1.88321456e-02,
-1.18349104e-01, 6.05449224e-02, -9.51659553e-02, 3.68858066e-02,
1.91742862e-01, 1.12727170e-02, 8.75690334e-02, 1.28931205e-01,
-1.19959362e-01, 2.18871821e-02, 2.12931830e-02, 1.41148914e-01,
1.60961898e-02, 1.12073208e-01, 7.65993036e-02, 2.12615547e-01,
-3.91121558e-02, 1.03139124e-02, 3.49290932e-02, 5.91188382e-02,
5.24049870e-02, 8.25222896e-02, 4.26205070e-02, 1.95044026e-02,
5.14112178e-02, 3.07260882e-01, 3.85103078e-02, 7.85399184e-02,
7.93297152e-02, 1.73739208e-01, 7.58889259e-02, 1.86395144e-01,
9.43319252e-02, 6.61893984e-01, 6.07926410e-02, 6.13787638e-02,
2.08892219e-01, 1.37177891e-01, 1.05859789e-01, 1.44533756e-01,
4.46232084e-02, 0.93263204e-01, 4.92656139e-02, 3.06907736e-02,
1.11297983e-01, -1.93150426e-03, -0.07598399e-05, -0.38536234e-02,
1.05464837e-01, 1.86128686e-03, 5.61480670e-02, -1.98478822e-02,
4.62071922e-02, 1.43803984e-01, 1.23924819e-02, 6.13787638e-02,
1.38817959e-01, 7.75837130e-02, -2.88170397e-02, -7.23937713e-02,
9.95518933e-02, 6.61893984e-01, 6.07926410e-02, 6.13787638e-02,
-1.05172261e-01, -6.01513575e-02, -1.01089971e-01, -1.45889714e-01,
9.95518933e-02, 6.61893984e-01, 6.07926410e-02, 6.13787638e-02,
-1.88928831e-01, -8.64595954e-02, -5.71974903e-03, 1.23101785e-01,
5.33551496e-02, -3.25974324e-02, -1.05572381e-01, 1.21803879e-01,
2.87674775e-01, 2.28137836e-02, 5.51734929e-02, 2.81891408e-02,
5.98645881e-02, 1.02717438e-01, 6.94197135e-04, -1.75339973e-01])
```

```
In [8]: y = list()
for k_i, noise_i in zip(X, noise_list):
    natural_log_of_the_x_is = np.log where as log10 is np.log10
    y_i = np.log10(x_i) + y.append(y_i)
    print(y_i)

[2.666699891899353, 0.7595937882645583, 2.166684399873824, 1.28513422620288346, 1.2962849866692785, 1.6976987519753934, 1.6996895628426843, 1.552418154718667, 1.71589
328792944, 7.058851788981, 1.78787807856565, 1.6138787617888915, 1.466483471074546, 8.764227781435642, 0.989392136544697, 2.0984638238784384, 2.5831617174152136,
2.178418670256676, 2.09287932282761, 0.952115946887709, 1.78774887586914388, 0.8339463211880806, 1.6134894504636388, 2.840943296851077, 1.153488139825443, 2.0637493256
9454, 1.1089730964374504, 0.90829651538399, 0.9733985971538399, 0.7198495511146, 1.493189583787935, 1.7111162395858266, 1.407823468643388, 2.313852099728425, 1.699653455614377, 0.345786822873842, 1.931410888445642, 0.328592576191199, 2.2669831
38429, 0.34899898080933, 2.1818267456268, 1.978075252828095, 1.027311449549966, 2.3248412834866, 2.388844937028947, 0.754925047489924, 0.784978823819384, 1.88854184424829,
893778088322, 0.88826371955263, 0.877849119232946, 2.249314867152353, 1.268711671489424, 2.18715160236779, 8.8656398159284, 0.784978823819384, 1.88854184424829,
9.1.485694754939817, 1.5761639760148238, 1.437163234856286, 1.789517351156743, 0.6655856057176135, 1.928905487769829, 1.242328951598002, 2.157782673519933, 2.34081
6823085, 0.7108951544674674, 1.333661935101571, 8.84444734861052, 0.107858156437855, 1.3860312687768775, 1.758438892512222, 0.818118740276885, 1.85871613078808, 1.85871613078808,
0.1.83478955455818, 2.144747231535882, 1.886733828276285, 2.880756013115857, 0.682874535953942, 1.818045613228856, 1.278014183225527, 0.03757624899357808, 0.8.6823085, 0.6944826919399238, 2.8808740919595106, 1.676215885825488, 1.5542191213416828, 2.1851450987269, 1.1931861648458163, 2.268916354
144785, 1.3898822897635441, 1.8232166744289538, 0.4775964923791013, 1.9523937423825743, 1.9157617602150214, 1.8628925937980772]
```

Step3: KNN Regression Modeling

Now use K-NN regression to obtain y values (= estimates of y) at x-values of 1, 3, 5, 7 and 9 for each of the following three schemes:

case1.2.3: the K neighbors contribute equally (separately for K = 1, 3, 5)

case1.5.6: each of the K neighbors has an influence that is inversely proportional to the distance from the point (separately for K = 1, 3, 5)

case7: all the N points contribute, with each contribution proportional to e-12d2, where d represents distance.

```
In [21]: from sklearn.neighbors import KNeighborsRegressor

In [22]: X_test = [1,3,5,7,9] # Should we do x.delete? or change this if required

In [23]: # convert shapes and to numpy to give as model input
X_numpy = np.asarray(X).reshape(-1,1) # use np.delete if required
y_numpy = np.asarray(y).reshape(-1,1)

In [24]: # case1: the K neighbors contribute equally
uniform_k_1_nn_reg_model_class = KNeighborsRegressor(n_neighbors=1, weights='uniform')
uniform_k_1_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [24]: KNeighborsRegressor(n_neighbors=1)

In [25]: # case2: the K neighbors contribute equally
uniform_k_3_nn_reg_model_class = KNeighborsRegressor(n_neighbors=3, weights='uniform')
uniform_k_3_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [25]: KNeighborsRegressor(n_neighbors=3)

In [26]: # case3: the K neighbors contribute equally
uniform_k_50_nn_reg_model_class = KNeighborsRegressor(n_neighbors=50, weights='uniform')
uniform_k_50_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [26]: KNeighborsRegressor(n_neighbors=50)

In [27]: # case4: each of the K neighbors has an influence that is inversely proportional to the distance from the point
distance_k_1_nn_reg_model_class = KNeighborsRegressor(n_neighbors=1, weights='distance')
distance_k_1_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [27]: KNeighborsRegressor(n_neighbors=1, weights='distance')

In [28]: # case5: each of the K neighbors has an influence that is inversely proportional to the distance from the point
distance_k_3_nn_reg_model_class = KNeighborsRegressor(n_neighbors=3, weights='distance')
distance_k_3_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [28]: KNeighborsRegressor(n_neighbors=3, weights='distance')

In [29]: # case6: each of the K neighbors has an influence that is inversely proportional to the distance from the point
distance_k_50_nn_reg_model_class = KNeighborsRegressor(n_neighbors=50, weights='distance')
distance_k_50_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [29]: KNeighborsRegressor(n_neighbors=50, weights='distance')

In [31]: # case7: all the N points contribute, with each contribution proportional to e-12d2, where d represents distance.
def custom_function(array_distances:list)->list:
    """
    Returns:
    (Array of weights)

    output_weights = []

    for dist in array_distances:
        power_value = -1 * 8 * dist * dist
        output_weights.append(np.exp(power_value))
    return output_weights

# We 100 is defined in the data preparation
k_all_n_nn_reg_model_class = KNeighborsRegressor(n_neighbors=N, weights=custom_function)
k_all_n_nn_reg_model_class.fit(X_numpy,y_numpy)

Out [31]: KNeighborsRegressor(n_neighbors=100,
weights=function custom_function at 0x7f56bc859d0=)
```

Step4 and Step5: Obtain predictions for the test set for each of the above seven cases and Plot the closest neighbor

Next three cells are helper functions that are reused

```
In [33]: def get_and_print_predictions(model_class_obj, x_test_values):
# Print the numerical values of the (x, y) pairs
y_hat_output = [model_class_obj.predict([xi]) for xi in x_test_values]
predictions_x_y_hat = []
for x_pred, y_pred in zip(x_test_values, y_hat_output):
    (x, y)
    current_prediction = (x_pred, y_pred[0][0])
    print (current_prediction)
    predictions.append(current_prediction)
    return prediction_x_y_hat

In [40]: #Also, plot the (x',y') points for each of these seven cases, where x' is the point
# out of the 100 sample points) closest to x and y' is the y-value of x'
# Reference: https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html#sklearn.neighbors.KNeighborsRegressor.kneighbors

def get_X_y_prime_for_plot(model_class_obj, preds_x_y_hat):
    x_values = []
    y_hat_values = []
    x_prime_values = []
    y_prime_values = []

    for x, y_hat in preds_x_y_hat:
        y_hat("Extracting closest neighbor for x=", x, ", y_hat=", y_hat)

        # n_neighbors = 1, because we need the closest neighbor
        closes_neighbor_info = model_class_obj.kneighbors(X=[x]),n_neighbors=1, return_distance=True)

        closes_neighbor_distance = closes_neighbor_info[0][0][0]
        closes_neighbor_index_in_X = closes_neighbor_info[1][0][0]

        # closes_neighbor is x_prime and simultaneously y_prime
        x_prime = X[closes_neighbor_index_in_X]
        y_prime = y[closes_neighbor_index_in_X]
        print("Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime)", (x_prime, y_prime))

        # gather the below values for scatter plot
        x_values.append(x)
        y_hat_values.append(y_hat)
        x_prime_values.append(x_prime)
        y_prime_values.append(y_prime)

    return x_values, y_hat_values, x_prime_values, y_prime_values

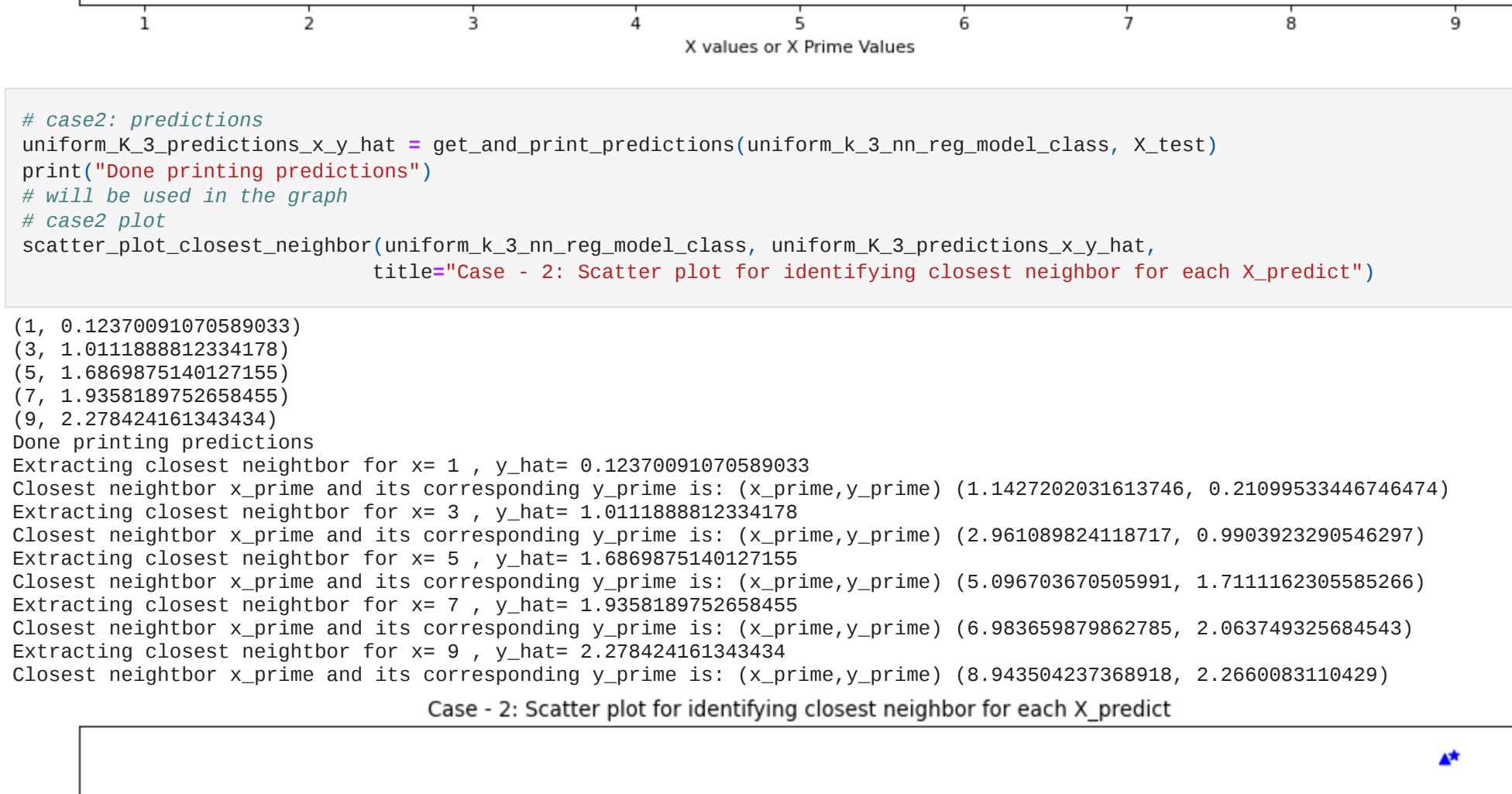
import matplotlib.pyplot as plt

def scatter_plot_closest_neighbor(model_class_obj, preds_x_y_hat, title: str):
    x_values, y_hat_values, x_prime_values, y_prime_values = get_X_y_prime_for_plot(model_class_obj, preds_x_y_hat)
    plt.figure(figsize=(15, 7), dpi=80)
    plt.scatter(x_values, y_hat_values, marker='x', color=['red', 'green', 'black', 'orange', 'blue'])
    plt.scatter(x_prime_values, y_prime_values, marker='o', color=['red', 'green', 'black', 'orange', 'blue'])
    plt.xlabel("x' values or X Prime Values")
    plt.ylabel("y' values or Y Prime Values")
    plt.title(title)
    plt.show()
```

```
In [51]: # case1: predictions
uniform_k_1_predictions_x_y_hat = get_and_print_predictions(uniform_k_1_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case1 plot
scatter_plot_closest_neighbor(uniform_k_1_nn_reg_model_class, uniform_k_1_predictions_x_y_hat,
title="Case - 1: Scatter plot for identifying closest neighbor for each X_predict")

(1, 0.21999533446746474)
(3, 0.996932299546297)
(5, 0.996932299546297)
(7, 2.063749325684543)
(9, 2.266983118429)
Done printing predictions
Extracting closest neighbor for x= 1, y_hat= 0.21999533446746474
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (1.1427282831613746, 0.21999533446746474)
Extracting closest neighbor for x= 3, y_hat= 0.996932299546297
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (2.96189824118717, 0.996932299546297)
Extracting closest neighbor for x= 5, y_hat= 1.711162395858266
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (5.096783678656991, 1.711162395858266)
Extracting closest neighbor for x= 7, y_hat= 1.869887518612155
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (6.986783678656991, 1.711162395858266)
Extracting closest neighbor for x= 9, y_hat= 2.266983118429
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)

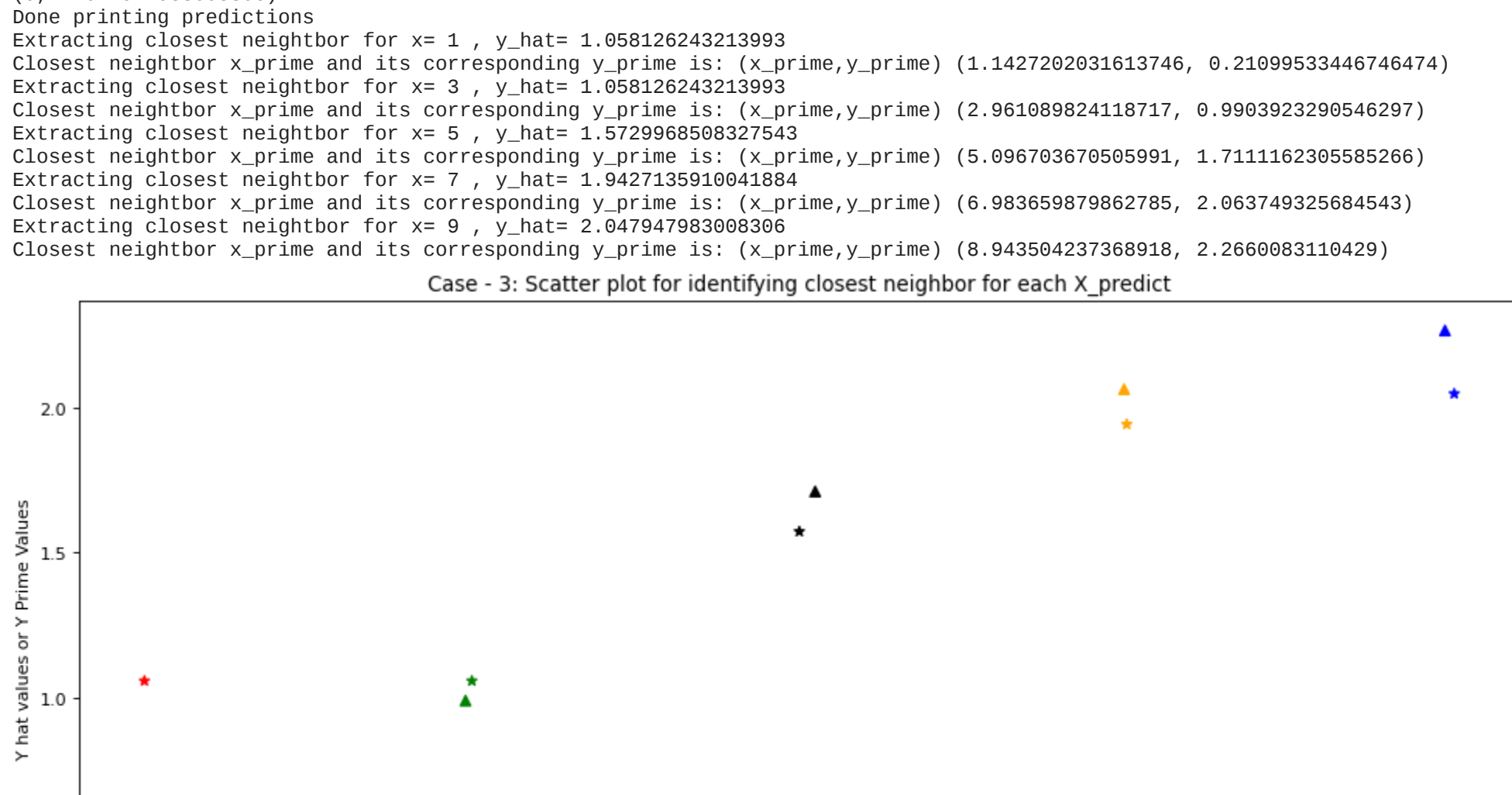
Case - 1: Scatter plot for identifying closest neighbor for each X_predict
```



```
In [50]: # case2: predictions
uniform_k_3_predictions_x_y_hat = get_and_print_predictions(uniform_k_3_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case2 plot
scatter_plot_closest_neighbor(uniform_k_3_nn_reg_model_class, uniform_k_3_predictions_x_y_hat,
title="Case - 2: Scatter plot for identifying closest neighbor for each X_predict")

(1, 0.1237089197058933)
(3, 1.011888812334178)
(5, 1.5729985808327543)
(7, 1.935818975268455)
(9, 2.27842161343434)
Done printing predictions
Extracting closest neighbor for x= 1, y_hat= 0.1237089197058933
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (1.1427282831613746, 0.21999533446746474)
Extracting closest neighbor for x= 3, y_hat= 1.011888812334178
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (2.96189824118717, 0.996932299546297)
Extracting closest neighbor for x= 5, y_hat= 1.5729985808327543
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (5.096783678656991, 1.711162395858266)
Extracting closest neighbor for x= 7, y_hat= 1.935818975268455
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (6.986783678656991, 1.711162395858266)
Extracting closest neighbor for x= 9, y_hat= 2.27842161343434
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)

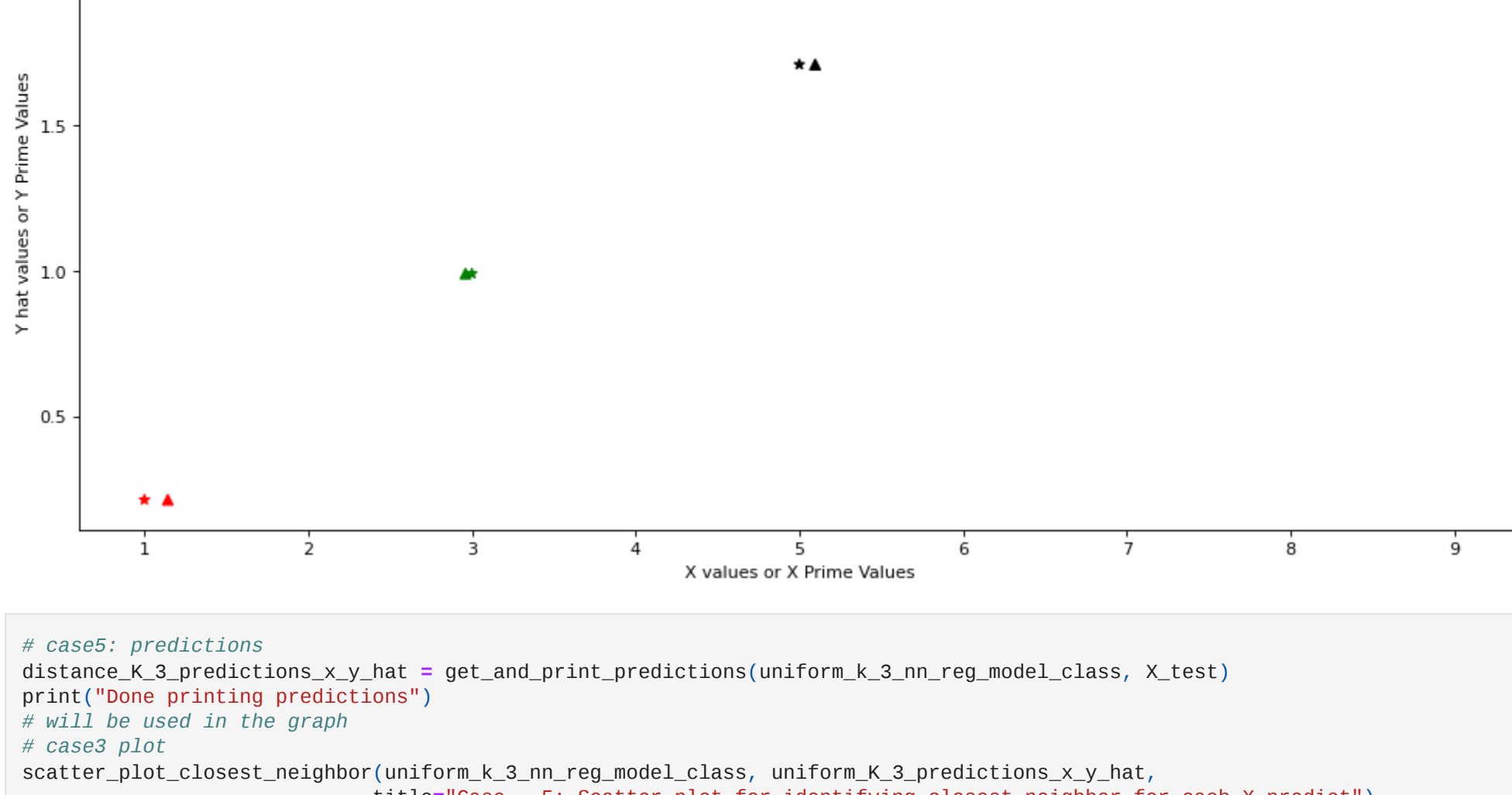
Case - 2: Scatter plot for identifying closest neighbor for each X_predict
```



```
In [52]: # case3: predictions
uniform_k_50_predictions_x_y_hat = get_and_print_predictions(uniform_k_50_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case3 plot
scatter_plot_closest_neighbor(uniform_k_50_nn_reg_model_class, uniform_k_50_predictions_x_y_hat,
title="Case - 3: Scatter plot for identifying closest neighbor for each X_predict")

(1, 1.058126423213993)
(3, 1.058126423213993)
(5, 1.5729985808327543)
(7, 1.9427135910041884)
(9, 2.9427135910041884)
Done printing predictions
Extracting closest neighbor for x= 1, y_hat= 1.058126423213993
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (1.1427282831613746, 0.21999533446746474)
Extracting closest neighbor for x= 3, y_hat= 1.058126423213993
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (2.96189824118717, 0.996932299546297)
Extracting closest neighbor for x= 5, y_hat= 1.5729985808327543
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (5.096783678656991, 1.711162395858266)
Extracting closest neighbor for x= 7, y_hat= 1.9427135910041884
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (6.986783678656991, 1.711162395858266)
Extracting closest neighbor for x= 9, y_hat= 2.9427135910041884
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)

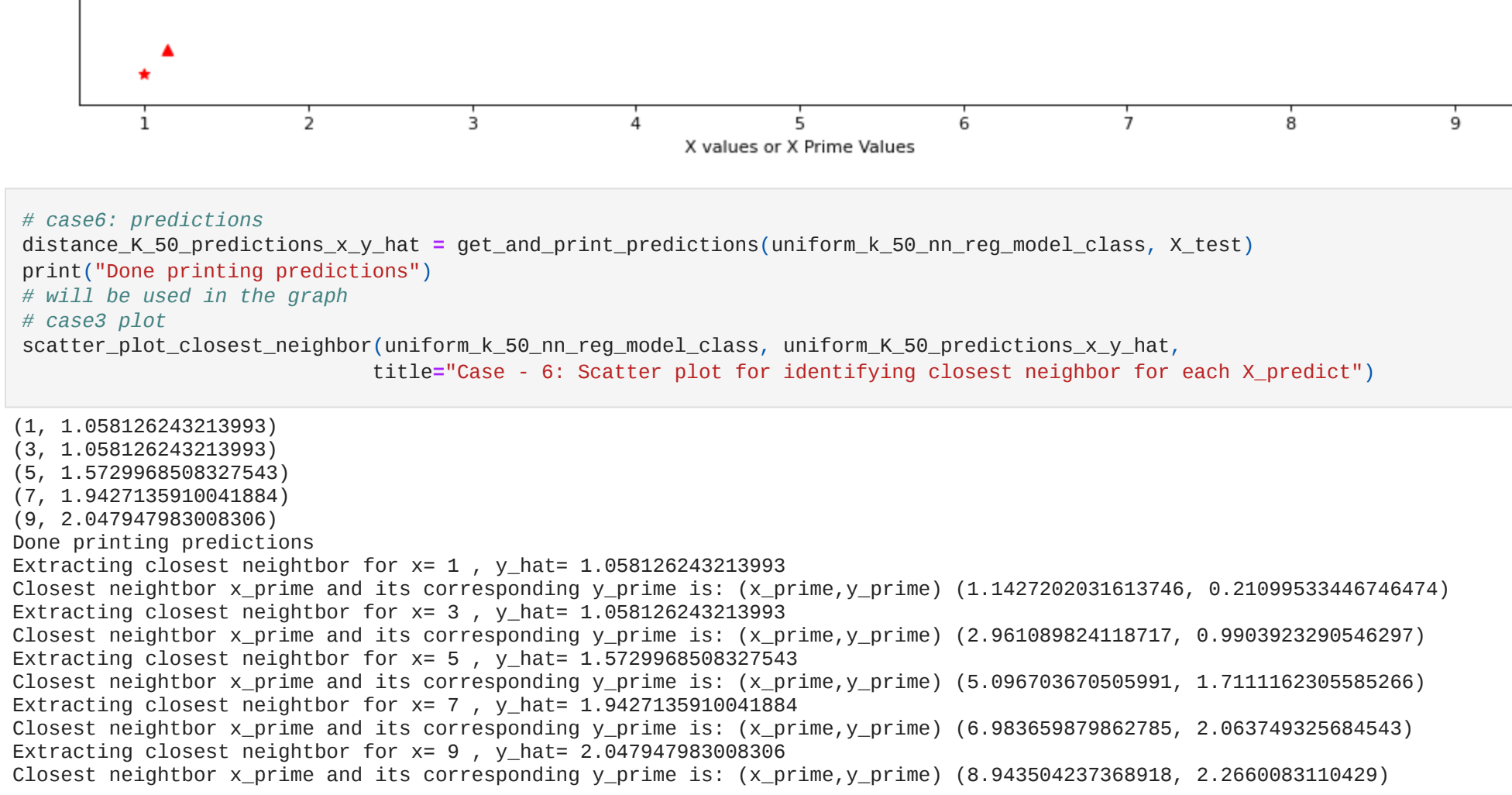
Case - 3: Scatter plot for identifying closest neighbor for each X_predict
```



```
In [54]: # case4: predictions
distance_k_1_predictions_x_y_hat = get_and_print_predictions(uniform_k_1_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case4 plot
scatter_plot_closest_neighbor(uniform_k_1_nn_reg_model_class, distance_k_1_predictions_x_y_hat,
title="Case - 4: Scatter plot for identifying closest neighbor for each X_predict")

(1, 0.21999533446746474)
(3, 0.996932299546297)
(5, 1.5729985808327543)
(7, 1.9427135910041884)
(9, 2.9427135910041884)
Done printing predictions
Extracting closest neighbor for x= 1, y_hat= 0.21999533446746474
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (1.1427282831613746, 0.21999533446746474)
Extracting closest neighbor for x= 3, y_hat= 0.996932299546297
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (2.96189824118717, 0.996932299546297)
Extracting closest neighbor for x= 5, y_hat= 1.5729985808327543
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (5.096783678656991, 1.711162395858266)
Extracting closest neighbor for x= 7, y_hat= 1.9427135910041884
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (6.986783678656991, 1.711162395858266)
Extracting closest neighbor for x= 9, y_hat= 2.9427135910041884
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)

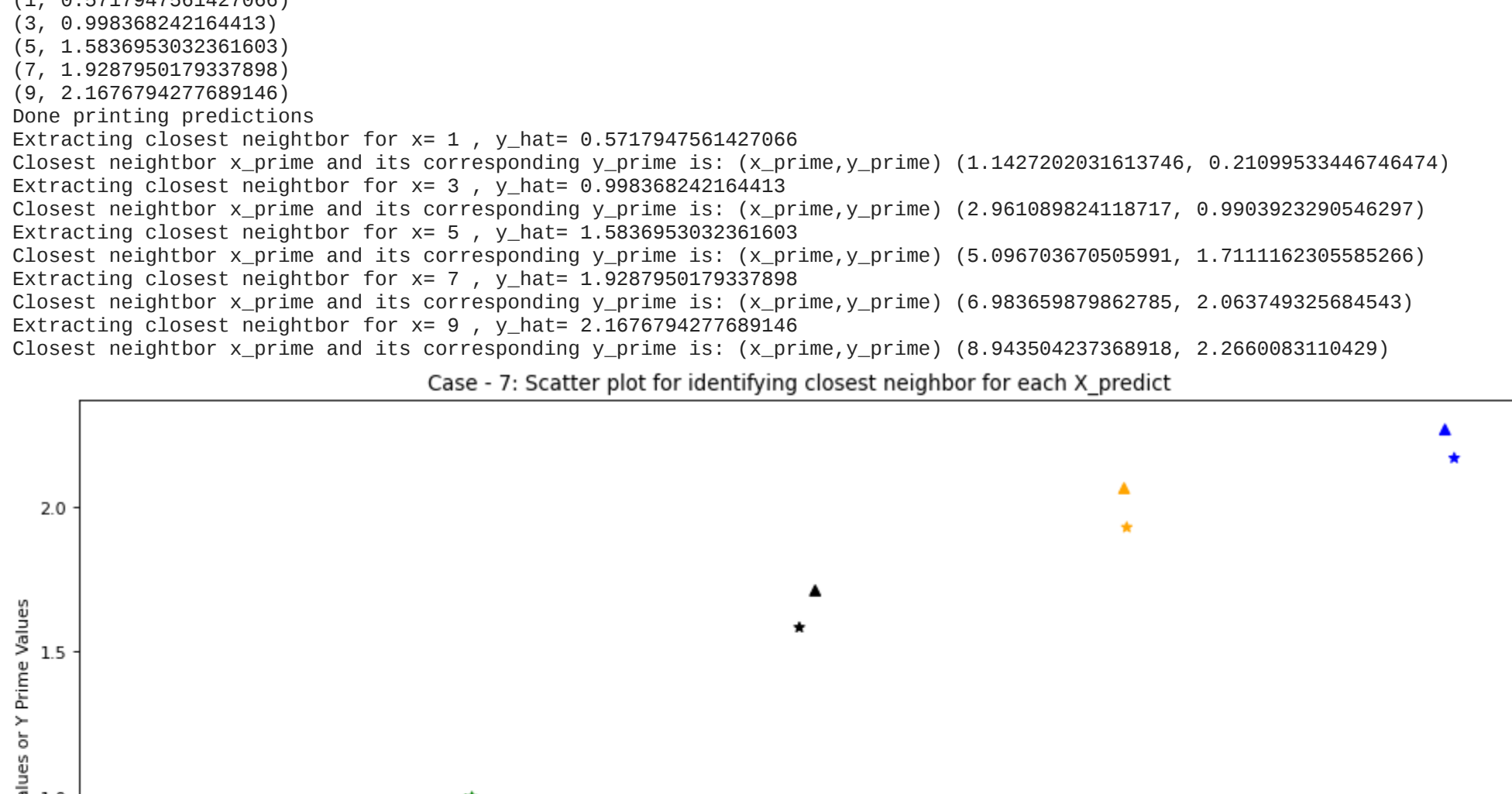
Case - 4: Scatter plot for identifying closest neighbor for each X_predict
```



```
In [55]: # case5: predictions
distance_k_3_predictions_x_y_hat = get_and_print_predictions(uniform_k_3_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case5 plot
scatter_plot_closest_neighbor(uniform_k_3_nn_reg_model_class, distance_k_3_predictions_x_y_hat,
title="Case - 5: Scatter plot for identifying closest neighbor for each X_predict")

(1, 0.1237089197058933)
(3, 1.011888812334178)
(5, 1.5729985808327543)
(7, 1.935818975268455)
(9, 2.27842161343434)
Done printing predictions
Extracting closest neighbor for x= 1, y_hat= 0.1237089197058933
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (1.1427282831613746, 0.21999533446746474)
Extracting closest neighbor for x= 3, y_hat= 1.011888812334178
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (2.96189824118717, 0.996932299546297)
Extracting closest neighbor for x= 5, y_hat= 1.5729985808327543
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (5.096783678656991, 1.711162395858266)
Extracting closest neighbor for x= 7, y_hat= 1.935818975268455
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (6.986783678656991, 1.711162395858266)
Extracting closest neighbor for x= 9, y_hat= 2.27842161343434
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)
Closest neighbor x_prime and its corresponding y_prime is: (x_prime,y_prime) (8.943564237868918, 2.266983118429)

Case - 5: Scatter plot for identifying closest neighbor for each X_predict
```



```
In [58]: # case6: predictions
distance_k_50_predictions_x_y_hat = get_and_print_predictions(k_all_n_nn_reg_model_class, X_test)
print("Done printing predictions")
# will be used in the graph
# case6 plot
scatter_plot_closest_neighbor(k_all_n_nn_reg_model_class, k_all_n_predictions_x_y_hat,
title="Case - 6: Scatter plot for identifying closest neighbor for each X_predict")

(1, 0.571947561427866)
(3, 1.58959530923616
```