

Tutorial k-Nearest Neighbours Regression

Marcos del Cueto

2020
August

1 Introduction

In the last years, there has been an explosion on machine learning (ML) tools that make their use much more affordable to users. However, a risk of such advancements is that it is easier to treat the ML tools as a black box. I have seen several online tutorials dealing with k-nearest neighbours (k-NN) classification. However, one must remember that k-NN can also be used as a regression tool, to estimate values of a continuous target property. Given how extended the use of regression methods are in all branches of science, I am surprised there are no more tutorials of such a basic (yet useful) method as k-NN regression.

Goals tutorial. I have prepared this short tutorial as a direct way to:

- Provide simple example on how to prepare data for regression.
- Show how to perform k-NN regression and optimize k.
- Exemplify how k-NN regression works.
- Show how different weight functions can affect k-NN prediction.
- Discuss some limitations of k-NN regression.

Classification vs regression

Python and scikit-learn

2 Generate data

To simplify things here, we will consider here a one-dimensional dataset. This means that we have only one descriptor (aka feature) \mathbf{x} , and a target property $f(\mathbf{x})$. As a generic example, I am going to use in this tutorial a small database formed by ten points that follow the function:

$$f(x) = e^x \tag{1}$$

The points considered in this example are shown in Table 1, as well as in Figure 1.

Table 1: Caption

x_n	$f(x_n)$
5.00	148.41
5.20	181.27
5.40	221.41
5.60	270.43
5.80	330.30
6.00	403.43
6.20	492.75
6.40	601.85
6.60	735.10
6.80	897.85

I show now the code used to generate this dataset, as well as Figure 1.

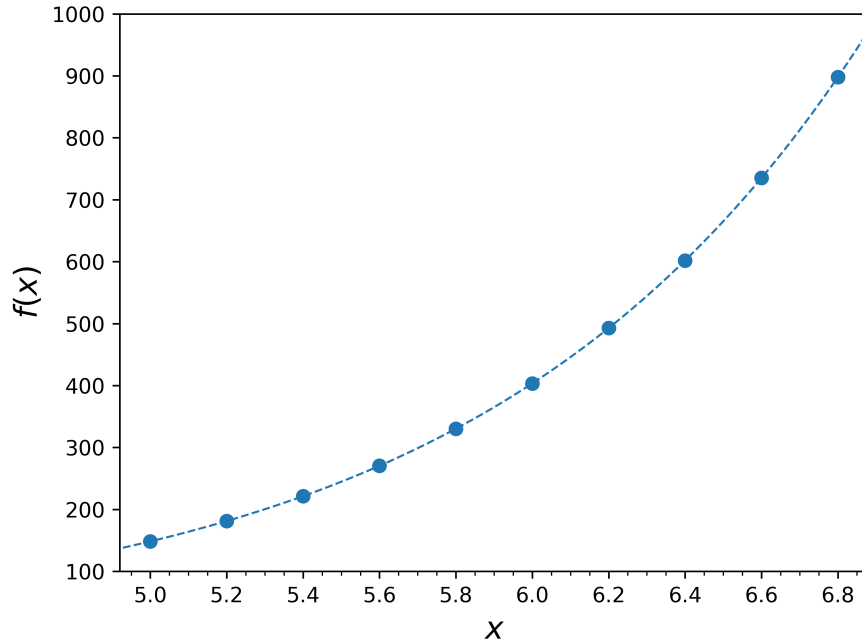


Figure 1: Caption.

Listing 1: Code1

```

1 import math
2 import matplotlib.pyplot as plt
3 from matplotlib.ticker import (MultipleLocator)
4 import numpy as np
5
6 ### 1) Generate data
7 list_x = []
8 list_y = []
9 for i in np.arange(5, 7, 0.2):
10     x = i
11     y = math.exp(x)
12     list_x.append(x)
13     list_y.append(y)
14     print("%.2f,%.6f" % (x, y))
15 list_x = np.array(list_x).reshape(-1, 1)
16 list_y = np.array(list_y)
17 basic_x = np.arange(4.9, 7.0, 0.01)
18 basic_y = [math.exp(x) for x in basic_x]
19 # Plot graph
20 plt.plot(basic_x, basic_y, color='C0', linestyle='dashed', linewidth=1)
21 plt.scatter(list_x, list_y, color='C0')
22 plt.xlabel('$x$', fontsize=15)
23 plt.ylabel('$f(x)$', fontsize=15)
24 plt.xticks(np.arange(5, 7, 0.2))
25 plt.xlim(4.92, 6.88)
26 plt.ylim(100, 1000)
27 axes = plt.gca()
28 axes.xaxis.set_minor_locator(MultipleLocator(0.05))
29 # Save plot into png
30 file_name='Fig1.png'
31 plt.savefig(file_name, format='png', dpi=600)
32 plt.close()

```

3 Data analysis

With k-NN regression, the y value at any given \mathbf{x} configuration is given as an average of nearby neighbours. When using this algorithm, one has to select two main parameters:

1. How many near neighbours are considered
2. How the average of the neighbours is calculated

3.1 Cross-validation to optimize k

Regarding the first of these points, a simple method to select the optimum number of k nearest neighbours to make accurate predictions is to use a grid search. A grid search involves trying different k values and finally choosing the one that minimizes the prediction error. As with any ML model, one is faced with the task of how to select the data that will be used to train the model and the data that will be used to test the model accuracy. Since we have a small amount of data here, we will use the leave-one-out (LOO) cross-validation (for more complex datasets, one may opt for a k -fold cross-validation).

I am not going to go into the details of LOO, but for a dataset with N samples, one trains the ML model with $N - 1$ points, and use the remaining point as a test. This procedure is repeated N times, so each point is used exactly once as a test (see Figure X). As a result, one ends with a predicted value for each point. Then, the predicted value of each point is compared with its actual value. As an error metric, we have used the root mean square error (rmse) here. Finally, we can compare the (rmse) obtained for different k values, and select as an optimum value the one that minimizes the rmse.

Here I provide the code necessary to do this:

Listing 2: Code2

```
1 import math
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn import preprocessing
5 from sklearn.model_selection import LeaveOneOut
6 from sklearn.metrics import mean_squared_error
7 from sklearn import neighbors
8 from matplotlib.ticker import MultipleLocator
9
10 ### 1) Generate data
11 list_x = []
12 list_y = []
13 for i in np.arange(5, 7, 0.2):
14     x = i
15     y = math.exp(x)
16     list_x.append(x)
17     list_y.append(y)
18 list_x = np.array(list_x).reshape(-1, 1)
19 list_y = np.array(list_y)
20 basic_x = np.arange(4.9, 7.0, 0.01)
21 basic_y = [math.exp(x) for x in basic_x]
22 ### 2) Leave one out
23 best_rmse = 0
24 best_k = 0
25 possible_k = [1, 2, 3, 4, 5, 6, 7, 8, 9]
26 for k in possible_k:
27     y_pred = []
28     x_pred = []
29     X = np.array(list_x).reshape(-1, 1)
30     validation = LeaveOneOut().split(X)
31     for train_index, test_index in validation:
32         X_train, X_test = list_x[train_index], list_x[test_index]
33         y_train, y_test = list_y[train_index], list_y[test_index]
34     # scale data
35     scaler = preprocessing.StandardScaler().fit(X_train)
36     X_train_scaled = scaler.transform(X_train)
```

```

37 X_test_scaled = scaler.transform(X_test)
38 knn = neighbors.KNeighborsRegressor(n_neighbors=k, weights='uniform')
39 pred = knn.fit(X_train_scaled, y_train).predict(X_test_scaled)
40 y_pred.append(pred)
41 x_pred.append(X_test)
42 mse = mean_squared_error(y_pred, list_y)
43 rmse = np.sqrt(mse)
44 print('k: %i, RMSE: %.2f' % (k, rmse))
45 if rmse < best_rmse or k==possible_k[0]:
46     best_rmse = rmse
47     best_k = k
48 print("Optimum: kNN, k=%i, RMSE: %.2f" % (best_k, best_rmse))

```

The output of the code above is:

```

k: 1, RMSE: 95.55
k: 2, RMSE: 74.87
k: 3, RMSE: 105.13
k: 4, RMSE: 122.74
k: 5, RMSE: 152.07
k: 6, RMSE: 173.34
k: 7, RMSE: 202.04
k: 8, RMSE: 230.51
k: 9, RMSE: 264.38
Optimum: kNN, k=2, RMSE: 74.87

```

We can use this optimum case, with $k=2$, to visualize how the k -NN regression algorithm is working. In Figure 2, we show how, when trying to predict the value of x_0 , the predicted value x'_0 is calculated as the average of its nearest neighbours, x_1 and x_2 ; x'_1 is predicted as the average of its two nearest neighbours, x_0 and x_2 and so on.

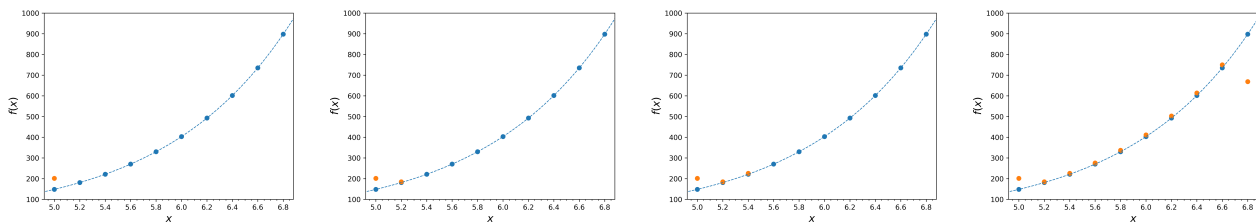


Figure 2: Caption.

In Figure 2, one can see how the predicted values behave reasonably well for the whole range, except on the extremes. This is because the prediction at the extremes will be biased due to the knowledge along the interval. There is also an additional issue that can't be observed in the previous Figure, since our \mathbf{x} values are equidistant. We can for example train the model with all our ten x values, and make predictions for all values inside the studied range. The code to make all these predictions is shown here, and the results are shown in Figure 3.

Listing 3: Code3

```

1 import math
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn import preprocessing
5 from sklearn.model_selection import LeaveOneOut
6 from sklearn.metrics import mean_squared_error
7 from sklearn import neighbors
8 from matplotlib.ticker import MultipleLocator
9
10 ### 1) Generate data
11 list_x = []
12 list_y = []
13 for i in np.arange(5, 7, 0.2):

```

```

14 x = i
15 y = math.exp(x)
16 list_x.append(x)
17 list_y.append(y)
18 list_x = np.array(list_x).reshape(-1, 1)
19 list_y = np.array(list_y)
20 basic_x = np.arange(4.9, 7.0, 0.01)
21 basic_y = [math.exp(x) for x in basic_x]
22
23 best_k = 2
24
25 #### 3) Do kNN regression
26 X_train = np.array(list_x).reshape(-1, 1)
27 y_train = np.array(list_y)
28 X_test = np.arange(5.0, 6.81, 0.01).reshape(-1, 1)
29 # scale data
30 scaler = preprocessing.StandardScaler().fit(X_train)
31 X_train_scaled = scaler.transform(X_train)
32 X_test_scaled = scaler.transform(X_test)
33 knn = neighbors.KNeighborsRegressor(n_neighbors=best_k, weights='uniform')
34 file_name='Fig3_1.png'
35 label_kNN = "k-NN_regression_(uniform)"
36 # kNN regression
37 y_test = knn.fit(X_train_scaled, y_train).predict(X_test_scaled)
38 # Plot graph
39 plt.plot(basic_x, basic_y, color='C0', linestyle='dashed', linewidth=1)
40 plt.scatter(list_x, list_y, color='C0', label='$x_n$')
41 plt.plot(X_test, y_test, color='C1', label=label_kNN)
42 plt.legend()
43 plt.xlabel('$x$', fontsize=15)
44 plt.ylabel('$f(x)$', fontsize=15)
45 plt.xticks(np.arange(5, 7, 0.2))
46 plt.xlim(4.92, 6.88)
47 plt.ylim(100, 1000)
48 axes = plt.gca()
49 axes.xaxis.set_minor_locator(MultipleLocator(0.05))
50 # Save plot into png
51 plt.savefig(file_name, format='png', dpi=600)
52 plt.close()

```

In Figure 3 we can see how the predicted values in between the training points x_n is constant, and it corresponds exactly to the average of the k nearest neighbours. This means that all k nearest neighbours have the same weight, which we specified with **weights='uniform'** (default option). However, one can clearly see how, in a case like this, it would be more accurate to give a larger weight to the neighbours that are closer.

3.2 Different weights for nearest neighbours

When using k -NN regression, one has to select a sensible weight function. In our case, one can see how it would be better to give a larger weight to points that are closer. This can be achieved by simply using the **weights='distance'** options, that calculates weights as the inverse of the distance between the studied i^{th} point and the neighbour n in the training set:

$$w_n = \frac{1}{d_n} \quad (2)$$

where d_n corresponds to the Euclidean distance between x_i and x_n . Note that one could use any other non-Euclidean metric, which might be more beneficial in other specific cases. I plan on writing a post specifically about this soon.

The result of using these new weights is shown in Figure 4 and we can observe how it behaves much better for the whole studied range.

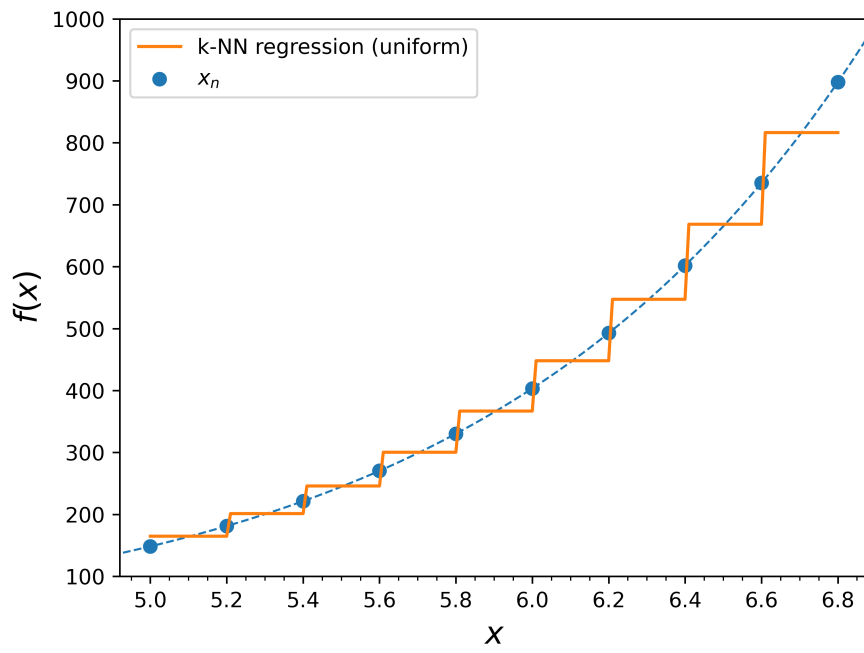


Figure 3: Caption.

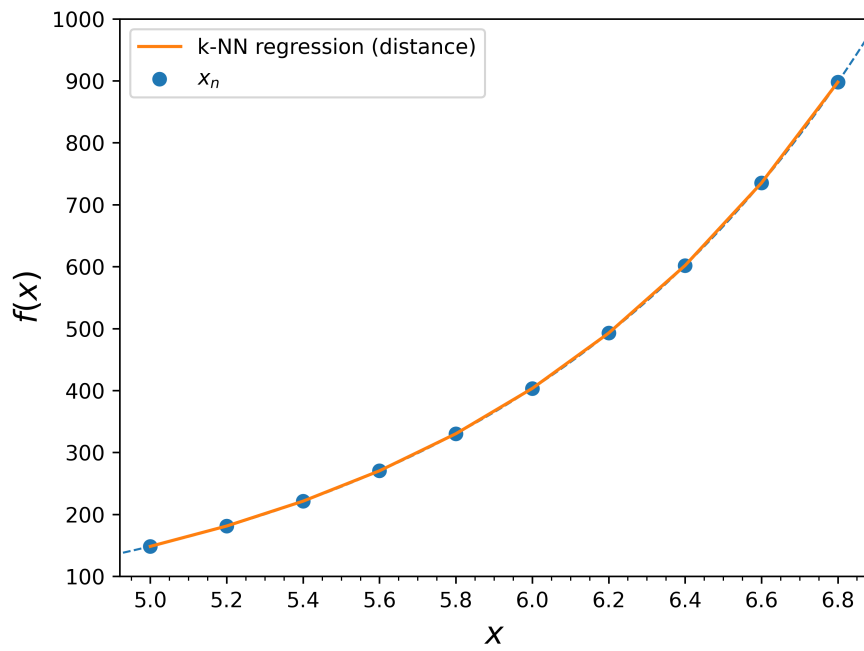


Figure 4: Caption.

4 Limitations k-NN regression

The example that we have seen here is extremely simple, and can be seen in Figure 4 how k-NN is able to make very accurate predictions within the studied range. However, one should note that even with such a simple algorithm as k-NN, one can make very accurate predictions (similar to much more sophisticated regression algorithms) in real cases in different branches of science.

However, it is worth mentioning a type of problem where k-NN performs particularly bad. This problematic task is extrapolation (i.e. predicting values that are outside the range of the training set). In general, this kind of extrapolation task is a challenge for all ML algorithms, but k-NN does particularly bad, since it purely relies in the information of neighbours in the training set, but we are in a region outside the training set.

We show in Figure 5 how k-NN regression (with different weights) performs for the range outside the training range. We can see how k-NN totally fails to predict the exponential behaviour expected outside the training range and it simply returns the weighted average value of the two nearest neighbours.

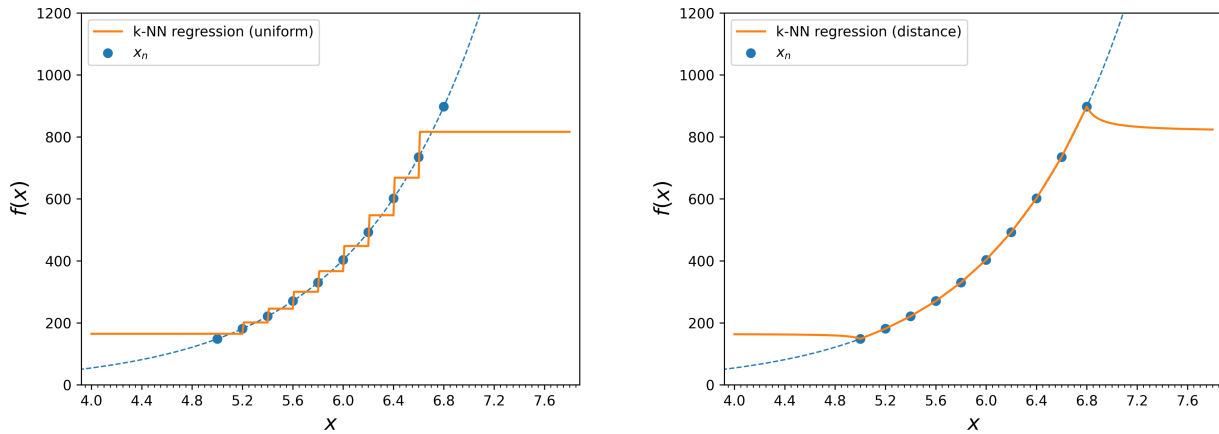


Figure 5: Caption.

The code used to generate Fig 5 is shown below.

Listing 4: Code3

```
1 import math
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn import preprocessing
5 from sklearn.model_selection import LeaveOneOut
6 from sklearn.metrics import mean_squared_error
7 from sklearn import neighbors
8 from matplotlib.ticker import MultipleLocator
9
10 ### 1) Generate data
11 list_x = []
12 list_y = []
13 for i in np.arange(5, 7, 0.2):
14     x = i
15     y = math.exp(x)
16     list_x.append(x)
17     list_y.append(y)
18 list_x = np.array(list_x).reshape(-1, 1)
19 list_y = np.array(list_y)
20 basic_x = np.arange(3.9, 8.0, 0.01)
21 basic_y = [math.exp(x) for x in basic_x]
22
23 best_k = 2
24
25 ### 4) Do kNN regression extrapolation
26 X_train = np.array(list_x).reshape(-1, 1)
```

```

27 y_train = np.array(list_y)
28 X_test = np.arange(4.0, 7.81, 0.01).reshape(-1, 1)
29 # scale data
30 scaler = preprocessing.StandardScaler().fit(X_train)
31 X_train_scaled = scaler.transform(X_train)
32 X_test_scaled = scaler.transform(X_test)
33 #####
34 ##### UNIFORM #####
35 #knn = neighbors.KNeighborsRegressor(n_neighbors=best_k, weights='uniform')
36 #file_name='Fig4.1.png'
37 #label_kNN = "k-NN regression (uniform)"
38 ##### DISTANCE #####
39 knn = neighbors.KNeighborsRegressor(n_neighbors=best_k, weights='distance')
40 file_name='Fig4.2.png'
41 label_kNN = "k-NN regression (distance)"
42 #####
43 # kNN regression
44 y_test = knn.fit(X_train_scaled, y_train).predict(X_test_scaled)
45 # Plot graph
46 plt.plot(basic_x, basic_y, color='C0', linestyle='dashed', linewidth=1)
47 plt.scatter(list_x, list_y, color='C0', label='$x_n$')
48 plt.plot(X_test, y_test, color='C1', label=label_kNN)
49 plt.legend()
50 plt.xlabel('$x$', fontsize=15)
51 plt.ylabel('$f(x)$', fontsize=15)
52 plt.xticks(np.arange(4, 8, 0.4))
53 plt.xlim(3.92, 7.88)
54 plt.ylim(0, 1200)
55 axes = plt.gca()
56 axes.xaxis.set_minor_locator(MultipleLocator(0.05))
57 # Save plot into png
58 plt.savefig(file_name, format='png', dpi=600)
59 plt.close()

```