

Predictor #2: k-nearest neighbour predictor

The second predictor we trained was the k-nearest neighbour predictor. The original data have been spitted into training data and testing data for predictor training purpose. The job is done in two phases, data preparation, parameter tuning and predictor evaluation.

Phase 1: Data preparation

First of all, we labeled both training data and testing data by their last column, +1 represents stable, -1 represents unstable. At this point, we trained our first demo predictor using the training data based on the *knn scikit tool* (with default parameters) in python, and testing the predictor on testing data, the accuracy on testing data at this stage is 0.774, which is not acceptable, and then we tried a range of k , from 1 to 50, and the predictors are tested on both training data (figure 1) and testing data (figure 2) by cross validation, the result in show below. By observing the original data, we found that the range of each column in the input data are distinct, i.e., the range of τ_{a1} (column 1) is around 0 to 10, and the range of g_2 (column 10) is around 0 to 1, in this case, g_2 may accidentally become more uninformative than τ_{a1} , for example, the calculation of distance in *knn* is based on the input features, a larger range feature has a larger impact on the distance, i.e., the Euclidian distance equation is: $D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$, we normalized all data by scaling them to the range of 0 to 1.

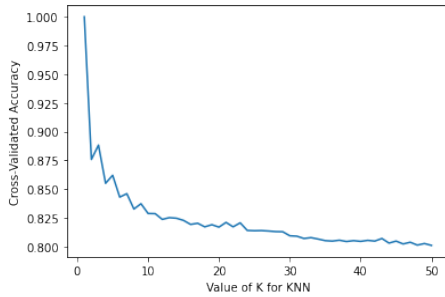


Figure 1 best $k = 1$, accuracy = 1 on training data

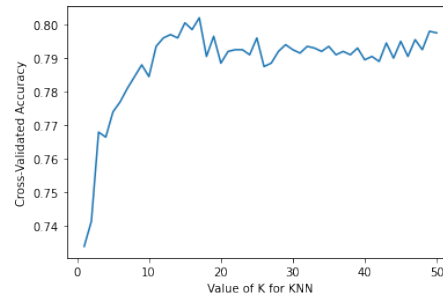


Figure 2 best $k = 17$, accuracy = 0.8020 on testing data

Phase 2: Parameter tuning and predictor evaluation

parameter tuning is one of the most important processes that dominate the behaviour of a predictor, and ultimately resulting in a most optimal combination of parameters. According to the python *sklearn*'s document, there are eight parameters, we only put four (underlined>) of them into our tuning process:

- $n_neighbors$: decides the number of neighbors will be used to calculate kneighbors, 1 to 50 will be tested.
- $weights$: there are two choices, uniform or distance, uniform means all points in each neighborhood are weighted equally, and distance means the closer point is weighted more, we should test both cases.
- $algorithm$: (auto, ball_tree, kd_tree or brute), since run time is not in our consideration, they will not be tested.
- $leaf_size$: only related to the construction time of the predictor and the memory cost of the query, so we can just use the default value, which is 30.
- p : determines the distance calculation procedure that the predictor will use, 1 (Manhattan distance), 2 (Euclidean distance) and arbitrary number (Minkowski distance).
- $metric$: the distance metric to use for the tree, we will only test Euclidean distance and Manhattan distance/
- $metric_params$: additional keyword arguments for the metric function, we have no additional arguments.
- n_jobs : determines the number of jobs will run in parallel, we will not test it.

Among all possible combinations, the one that has the best accuracy is $\{n_neighbors = 12, weights = distance, p = 1, metric = Manhattan\}$, the accuracy is 0.8785, the confusion matrices with *F1-score* see figure 3 and figure 4.

Confusion matrix

	Actual 1	Actual 2	sum_col
Predicted 1	5133 64.16%	5133 100%	10266
Predicted 2	2867 35.84%	2867 100%	5734
sum_row	8000	8000	16000

Actual

F1-score:

$$F1_{train} = \frac{TP}{TP + \frac{1}{2} \cdot (FP + FN)}$$

$$= \frac{2867}{2867 + \frac{1}{2} \cdot (0 + 0)}$$

$$= 1$$

Figure 3 Confusion matrix on training data

Confusion matrix

	Actual 1	Actual 2	sum_col
Predicted 1	1209 60.45%	205 10.25%	1414
Predicted 2	38 1.90%	548 27.40%	586
sum_row	1247	753	2000

Actual

F1-score:

$$F1_{test} = \frac{TP}{TP + \frac{1}{2} \cdot (FP + FN)}$$

$$= \frac{548}{548 + \frac{1}{2} \cdot (38 + 205)}$$

$$= 0.8185$$

Figure 3 Confusion matrix on testing data

Even though our predictor's *F1-score* in training data is 1, we ran cross validation multiple times, there is no significant change between different runs, overfitting is not occurred in the final predictor #2.