Predictor #2: k-nearest neighbour predictor

The second predictor we trained was the k-nearest neighbour predictor. The original data have been splitted into training data and testing data for predictor training purposes. 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. Secondly, by observing the original data, we found that the range of each column in the input data are distinct, i.e., the range of tau1 (column 1) is around 0 to 10, and the range of g2 (column 10) is around 0 to 1, in this case, g2 may accidentally become more uninformative than tau1, 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 Euclidean

distance equation is: D(X, Y) = $\sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$, so we normalized all data by scaling them to the range of 0 to 1.

At this point, we trained our first demo predictor by based on the front 80% of the training data using the *knn* scikit tool (with default parameters) in python, and tested the predictor on the remaining 20% of the training data, the accuracy at this stage is 0.86625 and the F-*score* is 0.7834, which is not high enough. So we tried tuning the parameters listed in the knn predictor, and tested it on both the training / validation data by cross validation.

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 in knn, we only put three (underlined) of them into our tuning process:

- <u>n_neighbors</u>: decides the number of neighbors will be used to calculate kneighbors, 3 to 20 will be tested.
- <u>weights</u>: 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 runtime 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.
- <u>p</u>: 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 can just accept the default algorithm.
- 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 on validation data set is { n_neighbors= 16, weights= distance, p = 1}, and the Accuracy is 0.878125, also to avoid overfitting, we train the model with cross-validation of n=10, and we got that: The confusion matrices with *F1-score* on training data see Figure 1.

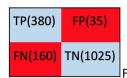


Figure 1:tn fp fn tp:1025 35 160 380,, F1-SCORE = 0.795812

Applying the model to predict the test data and we got Final Evaluation on testing data set with **Accuracy**: **0.861**, and the confusion matrices with *F1-score* on test data see Figure 2.

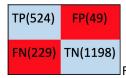


Figure 2: tn fp fn tp:1198 49 229 524, F1-SCORE = 0.790347

Reference: https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html