

# Amin Fadaee

## Machine Learning HW2

Python 3.5, Libraries: Numpy, Pandas, Random, Documentation: HTML, Bootstrap 4, jupyterMath

## KNN Algorithm

### Problem 1

KNN Algorithm needs sufficient data in order to have a reasonably good result and in times when there are not enough data points, the k-nearest neighbor for a test data deviates tremendously with the true label. By adding new dimensions to the data, the need for new point increase exponentially and sparse high-dimension data would make KNN a really poor choice for a classifier.

### Problem 2

The bias and variance of KNN algorithm is highly dependent on the **K** parameter. This relationship can be stated in the following way:

- For  $k = 1$  We have the maximum bias due to the fact that every test data would be classified as the class with highest prior probability. In other hand this would make the variance of the model as low as possible and if we suppose that the prior probability is fixed the variance would be actually zero for different sets of data points.
- In the case of  $k = 1$  we have the maximum variance, this is due to the fact that in this case the decision boundaries create a Voronoi cell network which is highly dependent on the positions of the data points and a new set of data would completely change the boundary.

### Problem 3

Based on Murphy's definition in his book **Machine Learning, a Probabilistic Approach**:

Does the model have a fixed number of parameters, or does the number of parameters grow with the amount of training data? The former is called a parametric model, and the latter is called a nonparametric model. Parametric models have the advantage of often being faster to use, but the disadvantage of making stronger assumptions about the nature of the data distributions. Nonparametric models are more flexible, but often computationally intractable for large datasets.

This definition is exactly correspondent to what KNN does and how it behaves.

### Problem 4

a. Using Euclidean Distance, 10-Fold cross validation has been done to find the best K for classifying the **Seeds Data** (<https://archive.ics.uci.edu/ml/datasets/seeds>) among options of K: 1, 3, 5, 7, 10. After performing the cross validation the following results were obtained:

K	Accuracy
1	0.9142
3	0.9809
5	0.9857
7	0.9864
10	0.9900

and  $K=1$  is the best choice for further classifications. The following functions were used in obtaining these results in `sklearn.py`:

- `sklearn.neighbors.KNeighborsClassifier` : By performing 10-Fold Cross Validation, finds the optimal k for KNN
- `sklearn.neighbors.KNeighborsClassifier` : Computes the KNN classification on test with the provided distance function

b. 10-Fold cross validation was again used to assess the performance of KNN with  $K=5$  with the following distance measures:

- Euclidean
- Manhattan
- Minkowski(4)
- Minkowski(1/2)
- Cosine

And here is the results:

Distance Measure	Accuracy
Euclidean	0.9857
Manhattan	0.9900
Minkowski(4)	0.9900
Minkowski(1/2)	0.9900
Cosine	0.9895

As we can see, for  $K=5$  the Cosine Distance obtains the best accuracy on the data. The following functions were used in obtaining these results in `sklearn.py`:

- `sklearn.neighbors.DistanceMetric` : Computes the euclidean distance between p1 and p2
- `sklearn.neighbors.DistanceMetric` : Computes the Manhattan distance between p1 and p2
- `sklearn.neighbors.DistanceMetric` : Computes the cosine distance between p1 and p2
- `sklearn.neighbors.DistanceMetric` : Computes the minkowski distance between p1 and p2 with parameter p
- `sklearn.neighbors.KNeighborsClassifier` : Computes the KNN classification on test with the provided distance function
- `sklearn.cross_validation.cross_val_score` : By performing 10-Fold Cross Validation, finds the test error of KNN for different distance measures.

### Problem 5

a. Based on the paper KNN shines for its simplicity and scalability, however it does not perform well in on imbalanced data.

b. Some ways were mentioned in dealing with imbalanced data, such as under-sampling of the major class instances or over-sampling for the minority class or changing the probability density distribution by setting different weights for different classes.

c. The WDKNN improves better in some cases because it assigns different weights on different nearest neighbors in a way that the weight has a reverse linear relationship with the point distance with the test data. The WDKNN algorithm has been implemented in `sklearn.py` and the results will be discussed in part e.

d. The paper deals with the imbalanced data by increasing the weight for the unsafe minority instances. This increase is determined by the number of majority instances that exist in that point's K-nearest neighborhood.

e. First the 7 original data sets were undergone some basic pre-processing:

- Meta-data in the beginning of the files were removed
- A header line was added with the names of each column
- The last column name is Class
- Positive was changed to 1 and Negative to 0
- Files format was changed to .csv

The following functions were implemented and used:

- `sklearn.neighbors.KNeighborsClassifier` : Computes the euclidean distance between p1 and p2
- `sklearn.neighbors.KNeighborsClassifier` : Computes the k nearest neighbor of test.
- `sklearn.neighbors.KNeighborsClassifier` : Computes the WDKNN classification on test with the provided distance function
- `sklearn.neighbors.KNeighborsClassifier` : Computes the Minority Class KNN classification on test with the provided distance function
- `sklearn.cross_validation.cross_val_score` : Performs the 5-Fold Cross validation of the above algorithms and returns G-Mean and F-Measure of the classification
- `sklearn.metrics.precision` : TP, FP, FN
- `sklearn.metrics.recall` : TP, FN, FP, FN
- `sklearn.metrics.g_mean` : TP, FN, FP, FN
- `sklearn.metrics.f_measure` : TP, FN, FP, FN

After performing both algorithms on all the 7 data sets, the G-Mean and F-Measure of each set was obtained and can be seen in the following table:

Data	WDKNN		MKKNN	
	F-Measure	G-Mean	F-Measure	G-Mean
yeast	0.7587	0.8153	0.7333	0.8522
svm3	0.6470	0.7782	0.6665	0.6371
yeast-4_vs_4	0.7294	0.7771	0.7022	0.7852
yeast-0-3-8_vs_7-8	0.3733	0.6227	0.3838	0.5898
yeast-0-5-8_vs_3-7-8-9	0.6341	0.7185	0.6192	0.7667
yeast-0-2-7-8_vs_3-9	0.8148	0.8755	0.8042	0.8889
svm-0-2-6_vs_3-5	0.7084	0.8295	0.7084	0.8236

As it can be clear, the results are nearly correspondent with the results in the paper and the minor differences are due to the randomness in 5-Fold CV; however, the paper contains an error in defining the F-Measure which was implemented correctly in both this code and their own.

## Decision Tree

### Problem 1

In decision trees, over-fitting occurs when the tree is designed so as to perfectly fit all samples in the training data set. Thus it ends up with branches with strict rules of sparse data. Thus this effects the accuracy when predicting samples that are not part of the training set.

### Problem 2

Pruning is a method to deal with over-fitting in decision trees. We can use pre-pruning (stop growing the tree after a certain depth) or post-pruning which is trimming the tree after it is built. The second method is more desirable due to the fact that in a lot of situations we can't derive a good depth before training.

### Problem 3

Based on the book **Understanding Machine Learning**:

A random forest is a classifier consisting of a collection of decision trees, where each tree is constructed by applying an algorithm  $A$  on the training set  $S$  and an additional random vector  $\theta$ , where  $\theta$  is sampled i.i.d. from some distribution. The prediction of the random forest is obtained by a majority vote over the predictions of the individual trees.

To specify a particular random forest, we need to define the algorithm  $A$  and the distribution over  $\theta$ . There are many ways to do this and here we describe one particular option. We generate  $\theta$  as follows. First, we take a random subsample from  $S$  with replacements; namely, we sample a new training set  $S'$  of size  $n'$  using the uniform distribution over  $S$ . Second, we construct a sequence  $J_1, J_2, \dots$ , where each  $J_i$  is a subset of  $[d]$  of size  $k$ , which is generated by sampling uniformly at random elements from  $[d]$ . All these random variables form the vector  $\theta$ . Then, the algorithm  $A$  grows a decision tree (e.g., using the ID3 algorithm) based on the sample  $S'$ , where at each splitting stage of the algorithm, the algorithm is restricted to choosing a feature that maximizes Gain from the set  $\theta$ . Intuitively, if  $k$  is small, this restriction may prevent overfitting.

Here are some advantages of using random forests:

- Can perform both classification and regression tasks
- Handle the missing values and maintains accuracy for missing data
- Won't overfit the model
- Handle large data set with higher dimensionality

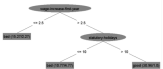
However random forest does have some draw backs and the most notable one is the little control we have on what the model does and its lack of interpretability.

### Problem 4

a. Accuracy= 73.68%

Confusion Matrix		
Confusion	Bad	Good
Bad	14	6
Good	9	20

- TP= 14
- FN= 6
- FP= 9
- TP= 28
- Precision=  $\frac{TP}{TP + FP} = 0.6125$
- Recall=  $\frac{TP}{TP + FN} = 0.7567$
- F1=  $2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 0.7087$

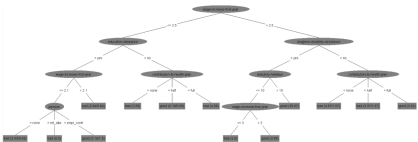


- The needed features for classification:
  - age -- increase --> bad -- yes --> 3. So we go RIGHT from ROOT
  - age --> bad -- yes --> 12. So we go RIGHT again.
  - The label is good

b. The `Pruned` parameter controls whether the tree should be pruned or not, unpruned trees are prone to overfitting the model.

Confusion Matrix		
Confusion	Bad	Good
Bad	14	6
Good	6	31

- TP= 14
- FN= 6
- FP= 6
- TP= 31
- Precision=  $\frac{TP}{TP + FP} = 0.6778$
- Recall=  $\frac{TP}{TP + FN} = 0.6778$
- F1=  $2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 0.6778$



- The needed features for classification:
  - age -- increase --> bad -- yes --> 3. So we go RIGHT from ROOT
  - age --> bad -- yes --> 12. So we go LEFT.
  - age --> bad -- yes --> 12. So we go RIGHT.
  - The label is good

c. As can be seen due to the lack of pruning in the algorithm this tree has a lot more leaves and branches than the pruned one.

Confusion Matrix		
Confusion	Bad	Good
Bad	15	4
Good	2	35

- TP= 15
- FN= 2
- FP= 4
- TP= 35
- Precision=  $\frac{TP}{TP + FP} = 0.8571$
- Recall=  $\frac{TP}{TP + FN} = 0.8571$
- F1=  $2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 0.8571$

d. As can be seen from the results, Random Forest obtained the best accuracy and in general, due to dealing with the variance part of the error it will achieve better performance.