

Assignment 1

1. Individual Assignment: Predicting Plant Types with k -Nearest Neighbours

Instructions: This exercise should be done “by hand”, that is, not using R or Python. All necessary calculations should be included in the submission, as well as brief explanations of what you do.

The data below is a small subset of the famous *Iris database*, first used by Sir R. A. Fisher (<http://archive.ics.uci.edu/ml/datasets/Iris>). This is perhaps the best known database in the pattern recognition literature. Fisher’s paper (R. A. Fisher (1936), “*The use of multiple measurements in taxonomic problems*”, Annals of Eugenics 7(2):179–188) is a classic in the field and is referenced frequently to this day. The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. There are four features: sepal length, sepal width, petal length and petal width. For the exercise below we only use sepal length and sepal width.

sepal length	sepal width	class
6.3	2.9	virginica
5.1	3.4	setosa
5.7	2.5	virginica
5.0	3.5	setosa
4.8	3.4	setosa
6.6	2.9	versicolor
6.3	3.3	versicolor
6.3	3.4	versicolor
6.0	3.4	versicolor
4.7	3.2	???
5.0	3.3	???
6.1	2.9	???
6.8	3.2	???
4.5	2.3	???
7.7	3.0	???

(a) Normalise the data using the Z-score normalisation.

Note: Since you have all data (training and input for prediction), you can (and should) use all data to compute the μ and σ required for the normalisation. Keep in mind,

however, that in reality you have to use the training data to compute μ and σ — the input for the future predictions will typically not yet be available to you!

- (b) Use the k -Nearest Neighbours method with $k = 3$ to predict the missing classes. Use the Euclidean norm in your distance calculations.

Note: When making predictions, you need to transform the input data for the predictions using the same μ and σ as for the transformation of the training data in (a).

2. Group Assignment: Predicting Wine Quality with k -Nearest Neighbours

Instructions: This exercise should be done using R and Python. The source codes as well as all relevant outputs should be included in the submission, as well as brief explanations of the code as well as what you see.

Download the data files “wine_quality-red.csv” and “wine_quality-white.csv” from <http://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/> and read the descriptions of the data sets at <http://archive.ics.uci.edu/ml/datasets/Wine+Quality>.

- (a) Build a k -Nearest Neighbours classifier in R for “wine_quality-white.csv” that:

1. loads the data file;
2. construct a new binary column “good wine” that indicates whether the wine is good (which we define as having a quality of 6 or higher) or not;
3. splits the data set into a training data set (~40%), a validation data set (~30%) and a test data set (~30%) — make sure you shuffle the record before the split;
4. normalises the data according to the Z-score transform;
5. loads and trains the k -Nearest Neighbours classifiers for $k = 1, \dots, 80$;
6. evaluates each classifier on the validation set and selects the best classifier;
7. predicts the generalisation error using the test data set, as well as outputs the result in a confusion matrix.

How do you judge whether the classifier is well-suited for the data set?

Note: You may want to refer to the two textbooks “Machine Learning with R” (Lantz) and “An Introduction to Statistical Learning” (James, Witten, Hastie, Tibshirani) for hints.

- (b) Build a k -Nearest Neighbours classifier in Python for “wine_quality-red.csv” that:

1. loads the data file;
2. construct a new binary column “good wine” that indicates whether the wine is good (which we define as having a quality of 6 or higher) or not;
3. splits the data set into a training data set (~50%) and a test data set (~50%) — make sure you shuffle the record before the split;
4. normalises the data according to the Z-score transform;
5. loads and trains the k -Nearest Neighbours classifiers for $k = 1, 6, 11, 16, \dots, 500$;
6. evaluates each classifier using 5-fold cross validation and selects the best classifier;
7. predicts the generalisation error using the test data set, as well as outputs the result in a confusion matrix.

How do you judge whether the classifier is well-suited for the data set?

Hints for Using Python

If you have not done so already, you may want to install Anaconda Python 2.7. You may want to use *Jupyter Notebook* for the exercises: it provides a more convenient interface than the standard Python interpreter.

For your Python implementation, you may want to rely on the Pandas, NumPy and SciKit-Learn libraries. Part of the group assessment challenge is to explore the libraries on your own!

1. Pandas is useful to import, explore and manipulate data. A cheat sheet from the company *enthought* is attached to this assignment, to get you started.
2. SciKit-Learn is the machine learning library. The SciKit-Learn documentation is excellent and available at <http://scikit-learn.org/stable/documentation.html>.
3. The Python tutorial is also excellent if you need a refresher: <https://docs.python.org/2.7/tutorial/>.

Here is a non-exclusive list of commands for Pandas that you might find useful:

```
import pandas as pd
df=pd.read_csv('filename',sep='\t',names=["col1name", "col2name"])
df['newcolumn']=df.col1name.apply(lambda x: x+2) # apply the function x -> x+2
df['newcolumn']=df.col1name.apply(f) # apply the function f defined elsewhere
df.describe()
df.ix[5]
df.head(3)
df[2:4]
```

For NumPy and SciKit-Learn you may want to look at the following commands:

```
from sklearn.utils import shuffle
X,y=shuffle(X,y)
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.cross_validation import corss_val_score
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
```