

# Assignment 1: Classification

---

## Deadlines

**Submission:** 11:59pm, 24 April, 2020 (Friday week 8, Sydney time)

## Late submissions policy

Late submissions are allowed for up to 3 days late. A penalty of 5% per day late will apply. Assignments more than 3 days late will not be accepted (i.e. will get 0 marks). The day cut-off time is 11:59pm.

## Marking

This assignment is worth 15 marks = 15% of your final mark. It must be completed individually. See the submission details sections for more information about submission in PASTA.

Your mark depends on the number of tests passed. 12 marks will be allocated to normal tests and 3 marks to hidden tests. For the normal tests, you will know the result (if you have passed or failed the test) after every submission while for the hidden tests, you will know the result after the due date.

## Task description

In this assignment you will investigate a real dataset by implementing multiple classification algorithms. You will first pre-process the dataset by replacing missing values and normalising the dataset with a min-max scaler. You will then evaluate multiple classification algorithms: K-Nearest Neighbour, Logistic Regression, Naïve Bayes, Decision Tree, Support Vector Machine and ensembles -Bagging, Boosting (AdaBoost, Gradient Boosting) and Random Forest. These will then be evaluated using the stratified 10-fold cross validation method. You will also apply a grid search to find the best parameters for some of the classifiers.

## Programming language

This assignment must be written in **Python**. It should use the classification algorithms from the **sklearn** (scikit-learn) library used in the tutorials. You can also use the **numpy** and **pandas** libraries.

On PASTA we have the following versions:

1. Python version 3.7.0
2. Numpy: 1.18.2
3. Pandas: 1.0.3
4. Scikit-learn: 0.22.2.post1 – identical contents to 0.22.2. See: [https://scikit-learn.org/stable/whats\\_new/v0.22.html](https://scikit-learn.org/stable/whats_new/v0.22.html)

Please install these versions on your laptop/desktop computers. You must ensure that your code works with these versions.

## Submission

This assignment will be submitted using PASTA.

## VPN

PASTA is located at <https://comp5318.it.usyd.edu.au/PASTA/>. In order to connect to the website, you'll need to be connected to the university VPN. The University provides a CISCO VPN client

<https://vpn.sydney.edu.au>; please read [this page](#) to find out how to connect to the VPN using the Cisco VPN client.

The students in China impacted by the travel ban should use the [new VPN \(FortiClient\)](#), which is the special VPN for accessing university resources from China. **Only students in China should use this VPN.** To access PASTA, it may also be necessary to run both the FortiClient and Cisco VPN; connect to FortiClient first and then to CISCO VPN.

If you are on-campus and connected to the University wireless network, you don't need a VPN to access PASTA.

### SSL warning

If you see a message like this in your browser for the PASTA's website: "Warning: Potential Security Risk Ahead", please give an exception for this site, e.g. in Firefox you need to click on "Advanced", then "Accept the Risk and Continue". You need to do this once and your browser will remember the exception. The reason for this is that PASTA's SSL certificate expired a few days ago; PASTA is very busy now; we will take it down and update the certificate after the assignment is finished.

### Multiple submissions

PASTA will allow you to make as many submissions as you wish, and each submission will provide you with feedback on each of the components of the assignment. You **last** submission before the assignment deadline will be marked, and the mark displayed on PASTA will be the final mark for your assignment.

## 1. Dataset preparation

### The data

The dataset for this assignment is the Breast Cancer Wisconsin. It contains 699 examples described by 9 numeric attributes. There are two classes – **class1** and **class2**. The features are computed from a digitized image of a fine needle aspirate of a breast mass of the subject. Benign breast cancer tumours correspond to **class1** and malignant breast cancer tumours correspond to **class2**.

The dataset should be downloaded from Canvas: **breast-cancer-wisconsin.csv**. This file includes the attribute (feature) headings and each row corresponds to one individual. Missing attributes in the dataset are recorded with a '?'.

### Data pre-processing

You will need to pre-process the dataset, before you can apply the classification algorithms. Three types of pre-processing are required:

5. The **missing attribute values** should be replaced with the mean value of the column using `sklearn.impute.SimpleImputer`.
6. **Normalisation** of each attribute should be performed using a min-max scaler to normalise the values between [0,1] with `sklearn.preprocessing.MinMaxScaler`.
7. The classes **class1** and **class2** should be changed to **0** and **1** respectively.
8. The value of each attribute should be formatted to 4 decimal places using `.4f`.

The correctness of your pre-processed data will be automatically tested by PASTA. Thus you need to make sure that you follow the input and output instructions carefully.

## 2. Classification algorithms with 10-fold cross-validation

You will now apply multiple classifiers to the pre-processed dataset, in particular: Nearest Neighbor, Logistic Regression, Naïve Bayes, Decision Tree, Bagging, Ada Boost and Gradient Boosting. All classifiers should use the **sklearn** modules from the tutorials. All random states in the classifiers should be set to **random\_state=0**.

You need to evaluate the performance of these classifiers using 10-fold cross validation from **sklearn.model\_selection.StratifiedKFold** with these options:

```
cvKFold=StratifiedKFold(n_splits=10, shuffle=True, random_state=0)
```

For each classifier, write a function that accepts the required input and that outputs the average cross-validation score:

```
def exampleClassifier(X, y, [options]):
```

```
    ...
```

```
    return scores, scores.mean()
```

where **X** contains the attribute values and **y** contains the class (as in the tutorial exercises).

More specifically, the headers of the functions for the classifiers are given below:

### K-Nearest Neighbour

```
def kNNClassifier(X, y, K)
```

It should use the **KNeighboursClassifier** from **sklearn.neighbours**.

### Logistic Regression

```
def logregClassifier(X, y)
```

It should use **LogisticRegression** from **sklearn.linear\_model**.

### Naïve Bayes

```
def nbClassifier(X, y)
```

It should use **GaussianNB** from **sklearn.naive\_bayes**

### Decision Tree

```
def dtClassifier(X, y)
```

It should use **DecisionTreeClassifier** from **sklearn.tree**, with information gain (the entropy criterion)

### Ensembles

```
def bagDTClassifier(X, y, n_estimators, max_samples, max_depth)
```

```
def adaDTClassifier(X, y, n_estimators, learning_rate, max_depth)
```

```
def gbClassifier(X, y, n_estimators, learning_rate)
```

These functions should implement Bagging, Ada Boost and Gradient Boosting using `BaggingClassifier`, `AdaBoostClassifier` and `GradientBoostingClassifier` from `sklearn.ensemble`. All should combine Decision Trees with information gain.

### 3. Parameter Tuning

For two other classifiers, Linear SVM and Random Forest, we would like to find the best parameters using grid search with 10-fold stratified cross validation (`GridSearchCV` in `sklearn`).

The split into training and test subsets should be done using `train_test_split` from `sklearn.model_selection` with stratification and `random_state=0` (as in the tutorials but with `random_state=0`).

Write the following functions:

#### Linear SVM

**def bestLinClassifier(X,y)**

It should use `SVC` from `sklearn.svm`.

The grid search should consider the following values for the parameters `C` and `gamma`:

`C = {0.001, 0.01, 0.1, 1, 10, 100}`

`gamma = {0.001, 0.01, 0.1, 1, 10, 100}`

The function should print the best parameters found, the best-cross validation accuracy score and the best test set accuracy score, see Section 4.

#### Random Forest

**def bestRFClassifier(X,y)**

It should use `RandomForestClassifier` from `sklearn.ensemble` with information gain

The grid search should consider the following values for the parameters `n_estimators`, `max_features` and `max_leaf_nodes`:

`n_estimators = {10, 20, 50, 100}`

`max_features = {'auto', 'sqrt', 'log2'}`

`max_leaf_nodes = {10, 20, 30}`

The function should print the best parameters found, best-cross validation accuracy score and best test set accuracy score, see Section 4.

### 4. Submission details - PASTA

This assignment is to be submitted electronically via the PASTA submission system.

#### How to submit to PASTA

Your program will need to be named `MyClassifier` and may only be written in Python 3.7. If you use Jupyter notebook to write your code, please do the following:

1. Ensure that you have a main method that is able to parse 3 command line arguments. PASTA will run your code with 3 different arguments. Details on the 3 arguments are given below.
2. Ensure that you export your Jupyter notebook as Python code and only submit the Python code to PASTA.

Before submitting to PASTA please make sure that all your Python programs (`MyClassifier.py` and all other supporting `.py` programs, if you have any) are zipped together. If you only have 1 Python program (`MyClassifier.py`), you do not need to zip it. Just submit that file to PASTA.

If you are zipping, please do **not** zip the folder containing your `MyClassifier.py` program (and your other supporting programs, if any). Zip only the `.py` files. To do this, select the `MyClassifier.py` file (and your other supporting scripts if any), right click, and compress/zip them. Then submit the resulting zip file to PASTA.

## Input

Your program should take 3 command line arguments:

1. The first argument is the path to the data file.
2. The second is the name of the algorithm to be executed or the option for print the pre-processed dataset:
  - a. **NN** for Nearest Neighbour.
  - b. **LR** for Logistic Regression.
  - c. **NB** for Naïve Bayes.
  - d. **DT** for Decision Tree.
  - e. **BAG** for Ensemble Bagging DT.
  - f. **ADA** for Ensemble ADA boosting DT.
  - g. **GB** for Ensemble Gradient Boosting.
  - h. **RF** for Random Forest.
  - i. **SVM** for Linear SVM.
  - j. **P** for printing the pre-processed dataset
3. The third argument is **optional**, and should **only be supplied** to algorithms which require parameters, namely NN, BAG, ADA and GB. It is the path to the file containing the parameter values for the algorithm. The file should be formatted as a csv file like in the following examples:
  - a. NN (note the capital K); an example for 5-Nearest Neighbour:
 

```
K
5
```
  - b. BAG:
 

```
n_estimators,max_samples,max_depth
100,100,2
```
  - c. ADA:
 

```
n_estimators,learning_rate,max_depth
100,0.2,3
```
  - d. GB:
 

```
n_estimators,learning_rate
100,0.2
```

For algorithms which do not require any parameters (LR, NB, DT, RF, SVM, P), the third argument should not be supplied.

The file paths given to you (as the first and third arguments) represent files that will be supplied to your program for reading. You can test your submission using any files you like, but PASTA will provide your submission with its own files for testing, so do not assume any specific filenames.

**Your program must be able to correctly infer X and y from the file.** As in the given dataset, there will be a header line and the last column will correspond to the class.

The following examples show how your program would be run:

1. Assuming we want to run the k-Nearest Neighbour classifier, the data is in a file called **breast-cancer-wisconsin-normalised.csv**, and the parameters are stored in a file called **param.csv**:

```
python MyClassifier.py breast-cancer-wisconsin-normalised.csv NN param.csv
```

2. Assuming we want to run the data pre-processing and the data is in a file called **breast-cancer-wisconsin.csv**:

```
python MyClassifier.py breast-cancer-wisconsin.csv P
```

3. Assuming we want to run Naïve Bayes and the data is in a file called **breast-cancer-wisconsin-normalised.csv**.

```
python MyClassifier.py breast-cancer-wisconsin-normalised.csv NB
```

## Output

Your program will output to standard output (i.e. “the console”).

## Pre-processing task

You will need to output your pre-processed data onto the console using print command. Let’s assume your normalised data looks like the following:

Clump Thickness	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Bland Chromatin	Normal Nucleoli	Mitoses	Class
0.1343	0.4333	0.5432	0.8589	0.3737	0.9485	0.4834	0.9456	0.4329	0
0.1345	0.4432	0.4567	0.4323	0.1111	0.3456	0.3213	0.8985	0.3456	1
0.4948	0.4798	0.2543	0.1876	0.9846	0.3345	0.4567	0.4983	0.2845	0

Your program output should look as follows:

```
0.1343,0.4333,0.5432,0.8589,0.3737,0.9485,0.4834,0.9456,0.4329,0
0.1345,0.4432,0.4567,0.4323,0.1111,0.3456,0.3213,0.8985,0.3456,1
0.4948,0.4798,0.2543,0.1876,0.9846,0.3345,0.4567,0.4983,0.2845,0
```

You must print the feature values to 4 decimal places as in the example above, e.g. 0.1 should be printed as 0.100. The class value is an integer.

Note also that there is no new line character after the last line. To print the last line, ensure you pass empty character as the end parameter for the print statement.

```
print(data, end='')
```

### Classifiers with 10-fold cross validation task

Your classifier should only print the mean accuracy score with precision of 4 decimal places using .4f. Note that .4f does rounding too. For example, if your mean accuracy is 0.987689, your program output should look as follows:

```
0.9877
```

Note there is no new line character after the accuracy. See above on how to print like such.

### Parameter tuning task

For Linear SVM, your program should output **exactly 4 lines**. The first line contains the optimal C value, the second line contains the optimal gamma value, and the third line contains the best cross-validation accuracy score formatted to 4 decimal places using .4f and the fourth line contains the test set accuracy score also formatted to 4 decimal places. For instance, if the optimal C and gamma values are 0.001 and 0.1 respectively, and the two accuracies are 0.999874 and 0.952512, your program output should look like the following:

```
0.001
0.1
0.9999
0.9525
```

For Random Forest, your program should output **exactly 5 lines**. The first line contains the optimal n\_estimators, the second line contains the optimal max\_features, the third line contains the optimal max\_leaf\_nodes, the fourth line contains the best cross validation accuracy score truncated to 4 decimal places using .4f and the fifth line contains the test set accuracy score also truncated to 4 decimal places. For instance, if the optimal n\_estimators, max\_features and max\_leaf\_nodes are 10, sqrt, 4 respectively, and the two accuracies are 0.997623 and 0.87243, your program output should look like the following:

```
10
sqrt
4
0.9976
0.8724
```