

Hong Kong University of Science and Technology
COMP 4211: Machine Learning
Spring 2021

Programming Assignment 1
Due: 8 March 2021, Monday, 11:59pm

1 Objectives

The objectives of this programming assignment are:

- To practise some data importing and preprocessing skills by using the **pandas** library in Python.
- To acquire a better understanding of supervised learning methods by using a public-domain software package called **scikit-learn**.
- To evaluate the performance of several supervised learning methods by conducting empirical study on a real-world dataset.

2 Dataset

You will use a wine quality dataset provided in the form of a ZIP file (**data.zip**). There are two **csv** data files. The following table shows the attributes of the data in each **csv** file.

File	# of records	Has label?	# of columns
train.csv	1,620	yes	11
test.csv	400	yes	11

The first 10 columns are features and the last column named ‘label’ indicates whether the wine is of high quality (1) or low quality (0).

3 Major Tasks

The assignment consists of four parts and a written report:

PART 1: Use **pandas** for data importing and preprocessing.

PART 2: Use the linear regression model for regression.

PART 3: Use the logistic regression model and single-hidden-layer neural network model for classification.

PART 4: Use **scikit-learn** to tune the hyperparameters and deal with the imbalanced dataset problem.

WRITTEN REPORT: Report the results and answer some questions.

More details will be provided in the following sections. Note that **[Q n]** refers to a specific question (the n th question) that you need to answer in the written report. All the experiments are expected to be done with Python 3.

4 Part 1: Data Preprocessing

In this part, you are required to preprocess the data and visualize the basic properties of the dataset. To be specific, you need to remove the duplicates and fill in the missing values with the median value. After you have finished handling the above cases, visualize the correlation between every two of the 10 features with a heatmap.

[Q1] Report the number of remaining records after duplicate removal and paste the screenshot of the heatmap.

5 Part 2: Regression

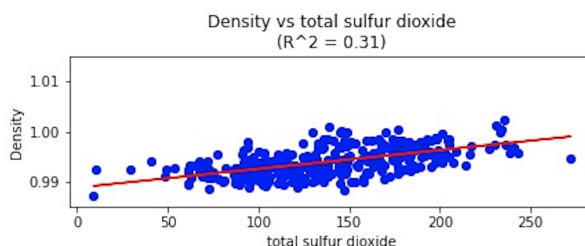
Linear regression is a basic model for regression which is expressed in the form $f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + \dots + w_dx_d$, where \mathbf{w} denotes the parameters to be learned from the data. Note that this basic model has no hyperparameters to set.

In this task, you will build six linear regression models in the first step, where each model uses one feature to find whether it is correlated with the feature ‘density’. The six features are ‘fixed acidity’, ‘residual sugar’, ‘chlorides’, ‘free sulfur dioxide’, ‘total sulfur dioxide’, and ‘alcohol’. Then, in the second step, you will build another linear regression model that explores the relationship between the linear combination of the six features and the feature ‘density’.

You are required to use the `train_test_split` submodule in `scikit-learn` to split the data in `train.csv`, with 80% for training and 20% for validation. You should set `random_state = 4211` for reproducibility.

[Q2] Report the validation R^2 score of each model to evaluate the relationship between different features and ‘density’.

[Q3] After training the models with the training set, use them to make prediction on the validation set. Then, plot the regression line and the data points of the validation set for each of the first six models. For illustration, the figure below shows a plot of the feature ‘density’ versus the feature ‘total sulfur dioxide’ and the regression line.



6 Part 3: Classification

In this task, you will build a logistic regression model as well as neural network classifiers to predict whether a certain type of wine is of high quality or not. You have to select the features according to some statistics and then use the selected features to complete the classification task.

You are also required to use the `train_test_split` submodule in `scikit-learn` to split the data, with 80% for training and 20% for validation. As before, we ask that you set `random_state = 4211` for reproducibility.

6.1 Feature Selection

To reduce the computational cost and remove the unrelated features, we would like to choose a subset of features for classification. You can use the feature selection module in `scikit-learn`. Select chi-squared statistics as the score for feature selection and then drop the two least important features.

[Q4] Report the score for each of the 10 features.

6.2 Logistic Regression

Learning of the logistic regression model should use a gradient-descent algorithm by minimizing the cross-entropy loss. It requires that the step size parameter η be specified. Try out a few values (<1) and choose one that leads to stable convergence. You may also decrease η gradually during the learning process to enhance convergence. This can be done automatically in `scikit-learn` when set properly.

Use the features selected in Section 6.1 above to train the model. During training, record the training time for the logistic regression model. After training, you are required to evaluate your model using accuracy, the F1 score¹ and the ROC curve on the *validation set*. Remember to standardize the features before training and validation.

[Q5] Report the model setting, training time, and performance of the logistic regression model. Since the solution found may depend on the initial weight values, you are expected to repeat each setting three times and report the corresponding mean and standard deviation of the training time, accuracy, and F1 score for each setting.

[Q6] Plot the ROC curve calculated on the validation set with the last model in [Q5] and report the AUC value.² Give one advantage of the ROC curve for model evaluation.

¹The F1 score is the harmonic mean of precision and sensitivity. You can find this metric in `sklearn.metrics`.

²Details of the ROC curve and AUC can be found in Wikipedia. Hints: You need to predict the probability of being of high quality for each wine first using the last model in [Q5].

6.2.1 Single-hidden-layer Neural Networks

Neural network classifiers generalize logistic regression by introducing one or more hidden layers. The learning algorithm for them is similar to that for logistic regression as described above. Remember to standardize the features before training and validation.

For the single-hidden-layer neural network model, you need to try different number of hidden units $H \in \{1, 2, 4, 8, 16, 32, 64, 128\}$. The hyperparameter `max_iter` can be set to 500 (default is 200) and `early_stopping` can be set to 'True' to avoid overfitting (default is 'False'). The other hyperparameters may just take their default values. During training, you are expected to record the training time of the models. After training, evaluate your models using the accuracy and the F1 score on the *validation set*. You have to report the accuracy and the F1 score for *each value* of H by plotting them using `matplotlib`.

[Q7] Report the model setting, training time, and performance of the neural networks for each value of H . You are also expected to repeat each setting three times for the same hyperparameter setting and report the mean and standard deviation of the training time, accuracy, and F1 score for each setting.

[Q8] Compare the training time, accuracy and F1 score of the logistic regression model and the best neural network model.

[Q9] Plot the accuracy and the F1 score for different values of H . Suggest a possible reason for the gap between the accuracy and the F1 score.

[Q10] Do you notice any trend when you increase the hidden layer size from 1 to 128? If so, please describe what the trend is and suggest a reason for your observation.

7 Part 4: Performance Enhancement

7.1 Hyperparameter Tuning

In this task, you need to use grid search to tune a single-hidden-layer neural network model to predict whether a type of wine is of high quality or not. Use the features selected in Section 6.1 for training and testing.

This time, you are required to evaluate your model on the test set `test.csv` provided. You need to import `test.csv` as a data frame and standardize the features using the statistics you used for the training data. Remember to standardize the features in the training set as well. (We assume that the test set comes from the same distribution as the training set.)

You are required to use the `model_selection` submodule in `scikit-learn` to facilitate performing grid search cross validation for hyperparameter tuning. This is done by randomly sampling 80% of the training instances to train a classifier and then validating it on the remaining 20%. Five such random data splits are performed and the average over these five trials is used to estimate the generalization performance. You are expected to search at least 10 combinations of the hyperparameter setting. Set the `random_state` hyperparameter of the single-hidden-layer neural network to 4211 for reproducibility and `early_stopping` to 'True' to avoid overfitting.

[Q11] Report 10 combinations of the hyperparameter setting.

[Q12] Report the three best hyperparameter settings in terms of accuracy as well as the mean and standard deviation of the validation accuracy of the five random data splits for each hyperparameter setting.

[Q13] Use the best model in terms of accuracy to predict the instances in the test set (`test.csv`). Report the accuracy, F1 score and the confusion matrix of the predictions on the test set.

7.2 Oversampling

By counting the number of records in the high-quality and low-quality classes, you will notice that the training set is imbalanced with more low-quality examples than high-quality ones. In this task, you will apply the oversampling strategy to tackle this problem. To be specific, you have to randomly oversample from the minority class and add these examples to the training set so that the two classes will become balanced in size. Set the `random_state` hyperparameter of the single-hidden-layer neural network and oversampling to 4211 for reproducibility and `early_stopping` to 'True' to avoid overfitting.

Based on the resampled dataset, perform grid search again with the same hyperparameter settings used in Section 7.1. Remember to standardize the features.

[Q14] Name another two methods to deal with the imbalanced dataset problem.

[Q15] Use the best model in terms of accuracy to predict the instances in the test set (`test.csv`). Report the accuracy, F1 score and the confusion matrix of the predictions on the test set.

[Q16] Compare the accuracy, F1 score and the confusion matrix with those in Section 7.1.

8 Report Writing

Answer [Q1] to [Q16] in the report.

9 Some Programming Tips

As is always the case, good programming practices should be applied when coding your program. Below are some common ones but they are by no means complete:

- Using functions to structure program clearly
- Using meaningful variable and function names to improve readability
- Using consistent styles
- Including concise but informative comments

For `scikit-learn` in particular, you are recommended to take full advantage of the built-in classes which can keep your program both short and efficient. Proper use of implementation tricks often leads to speedup by orders of magnitude. Please be careful to choose the built-in

models that are suitable for your tasks, e.g., `sklearn.linear_model.LogisticRegression` is not a correct choice for our logistic regression model since it does not use gradient descent.

10 Assignment Submission

Assignment submission should only be done electronically using the Course Assignment Submission System (CASS):

<https://cssystem.cse.ust.hk/UGuides/cass/student.html>

There should be two files in your submission with the following naming convention required:

1. **Report** (with filename `<StudentID>_report`): preferably in PDF format.
2. **Source code and a README file** (with filename `<StudentID>_code`): all necessary code for running your program as well as a brief user guide for the TA to run the programs easily to verify your results, all compressed into a single ZIP or RAR file. The data should not be submitted to keep the file size small.

When multiple versions with the same filename are submitted, only the latest version according to the timestamp will be used for grading. Files not adhering to the naming convention above will be ignored.

11 Grading Scheme

This programming assignment will be counted towards 10% of your final course grade. Note that the plus sign (+) in the last column of the table below indicates that reporting without providing the corresponding code will get zero point. The maximum scores for different tasks are shown below:

Grading scheme	Code (60)	Report (+40)
Part 1		
- Remove the duplicates + [Q1]	1	+1
- Fill in the missing values	1	
- Visualize the correlation between features + [Q1]	1	+1
Part 2		
- Build the linear regression model	3	
- Compute the R^2 score of the 6 linear regression models + [Q2]	2	+3
- Plot the regression line and the data points for each of the 6 linear regression models + [Q3]	3	+3
Part 3		
- Select the features according to chi-squared scores + [Q4]	3	+2
- Build the logistic regression model by adopting the gradient descent optimization algorithm	6	
- Compute the training time, accuracy, and F1 score of the logistic regression model + [Q5]	3	+2
- Plot the ROC curve and report the AUC value + [Q6]	3	+3
- Build the single-hidden-layer neural network model	6	
- Compute the training time, accuracy, and F1 score for each value of H in the single-hidden-layer neural network model + [Q7]	3	+3
- [Q8]		2
- Plot the accuracy and F1 score with different values of H for the single-hidden-layer neural network model + [Q9]	3	+3
- [Q10]		2
Part 4		
- Grid search on the single-hidden-layer neural network model for at least 10 combinations + [Q11]	6	+2
- Report the 3 best hyperparameter settings and the validation accuracy (both mean and standard deviation) for each setting + [Q12]	6	+2
- Report the accuracy and F1 score on the test set and visualize the confusion matrix + [Q13]	4	+4
- Oversample the dataset	2	
- [Q14]		2
- Grid search on the single-hidden-layer neural network model and report the accuracy, F1 score as well as the visualization of the confusion matrix on the test set + [Q15]	4	+3
- [Q16]		2

Late submission will be accepted but with penalty. The late penalty is deduction of one point (out of a maximum of 100 points) for every minute late after 11:59pm. Being late for a fraction of

a minute is considered a full minute. For example, two points will be deducted if the submission time is 00:00:34. At most one NQA coupon may be used to entitle you to submit this project late for one day without grade penalty.

12 Academic Integrity

Please refer to the regulations for student conduct and academic integrity on this webpage: <https://acadreg.ust.hk/generalreg>.

While you may discuss with your classmates on general ideas about the assignment, your submission should be based on your own independent effort. In case you seek help from any person or reference source, you should state it clearly in your submission. Failure to do so is considered plagiarism which will lead to appropriate disciplinary actions.