

# MIE 451-1513 Decision Support Systems

## Assignment 2: Machine Learning (ML)

This assignment allows students to apply machine learning knowledge through document classification using supervised learning techniques and to perform performance analysis of the learned classifier. In addition, this assignment contains a competitive component to encourage students to explore feature sets, feature representations and hyperparameter tuning.

- Programming language: Python (Google Colab Environment)
- Due Date: Posted in Syllabus

**Marking scheme and requirements:** Full marks will be given for (1) working, readable, reasonably efficient, documented code that achieves the assignment goals, (2) for providing appropriate answers to the questions in Jupyter notebooks (`ml_assignment.ipynb` and `ml_assignment_Q7.ipynb`) committed to the student's github assignment repository, and (3) attendance at the code review session, running your solution notebook for instructors, and providing clear and succinct answers in response to instructor questions regarding your solution.

Please note the plagiarism policy in the syllabus. If you borrow or modify any *multiline* snippets of code from the web, you are required to cite the URL in a comment above the code that you used. You do not need to cite tutorials or reference content that demonstrate how to use packages – you should certainly be making use of such content.

### What/how to submit your work:

1. All your code should be included in the notebooks named `ml_assignment.ipynb` and `ml_assignment_Q7.ipynb` provided in the cloned assignment repository.
2. Commit and push your work to your github repository in order to submit it. Your last commit and push before the assignment deadline will be considered to be your submission. You can check your repository online to make sure that all required files have actually been committed and pushed to your repository.
3. A link to create a personal repository for this assignment is posted on Quercus.

**Credit:** This lab's notebook material has been prepared based on an Advanced Scikit-Learn tutorial provided by Data Scientist Workbench.

### Notes that you should pay attention to

1. The same auto-grade restrictions apply for this assignment as they did for the previous assignments. Please check the pinned posts about autograder on Piazza for more information.
2. The autograder will be making 2 separate commits. The first will run test cases for Q1-Q6. The second is for the competitive portion of Q7 where we'll test performance on our own train and test sets. The second commit may occur hours after the first. Please do not push any changes to your repository in between these two commits, or we will consider your submission as late.
3. DO NOT use global variables in your functions. You can only use variables passed into functions or defined within classes.
4. DO NOT change the names of functions and classes.
5. During the code review, we will ask questions about your free text answers and your results in general as well as some details about your code and possible alternative approaches you might have taken. If you made both an on-time and late submission, you will have to tell us which version we should use for code review with a 30% deduction applying to the code review if you wish us to use your late submission. Note: you will receive a zero for the code review if you did not make an on-time or late submission – in this case you are considered to not have submitted your assignment.

**This assignment has 7 points in total, allocated as follows:**

- Auto-grading points(5 points):
  - Q1: 1.5 points
  - Q2: 0.5 points
  - Q3: 0.5 points
  - Q4: 0.5 points
  - Q5: 0.5 points
  - Q6: 0.5 points
  - Q7: 1 point
- Code review (2 points)

Points will be deducted from the auto-grading points for missing or incomplete text answers.

# 1 Before the Introductory lab

In the lab you'll familiarize yourself with several python machine learning libraries;

- **scikit-learn** a powerful machine learning library for python.
- **pandas** a powerful python data analysis toolkit
- **numpy and scipy** powerful computing packages for python
- **matplotlib** a python 2D plotting library

The scikit-learn documentation can be found on its website at <http://scikit-learn.org/stable/>

The pandas documentation can be found on its website at <http://pandas.pydata.org/pandas-docs/stable/>

The numpy and scipy documentations can be found on its website at <http://docs.scipy.org/doc/>

The matplotlib documentation can be found on its website at <http://matplotlib.org/>

You just need to execute the import statements at the top of the lab ipynb notebook to import these standard Python libraries.

# 2 In the Introductory lab

The lab can be found within your assignment-ml repository, under the `lab_ml.ipynb` file; clone it and open the Jupyter Notebook in the Google Colab to follow along.

In the Introductory lab section, we will go through the process of loading a dataset called “20 Newsgroups” and run a baseline classification using logistic regression.

The data you need is distributed on the course website but should be identical to the data found at <http://qwone.com/~jason/20Newsgroups/> with the following description:

The 20 Newsgroups data set is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. To the best of my knowledge, it was originally collected by Ken Lang, probably for his Newsweeder: Learning to filter netnews paper, though he does not explicitly mention this collection. The 20 newsgroups collection has become a popular data set for experiments in text applications of machine learning techniques, such as text classification and text clustering.

# 3 Main Assignment

In the introductory lab, you have been exposed to the 20 newsgroup dataset. The lab included loading and basic preparation of the data and a baseline classification using logistic regression and covered the libraries used in this assignment.

In this assignment, we will analyze different choices when configuring your classifier:

**Feature Set** Feature selection may not only improve classification performance (especially when data is sparse), but it can also improve computational performance. For this assignment, you can experiment with any features and feature selection technique.

**Feature Encoding** As two simple variations of feature encoding, consider a boolean  $\{0, 1\}$  vector encoding (as produced in the lab) as well as a term frequency (TF) encoding (which you need to produce yourself).

**Amount of Data** While generally you should use all training data available, it is instructive to vary the amount of training data provided to an algorithm to assess the impact of the amount of data on learning performance.

**Hyperparameters** All algorithms typically have at least one hyperparameter that needs tuning. We should use cross-validation for tuning hyperparameters in practice, but in this assignment, we will simply analyze performance as a function of hyperparameters.

We provide a template notebook for the coding part that you must use in your submission. Autograding requires that all function names are not modified and that the entire notebook executes in one pass end-to-end in Colab environment (please verify before submitting). Please fill in the missing part of the functions in the template notebook to provide the experimental results. The function you need to implement for each part of the assignment is clearly noted in the question description and the return value is validated using asserts in the end of the functions in the template notebook.

Please note the following:

- For all questions except Q3, use a hyperparameter that you find performs well (can be the default value). In general one should use nested cross-validation (CV) for hyperparameter tuning, but we avoid that here due to the time-consuming nature of nested CV.
- Do NOT change the function and class names. We need those for grading.
- Do NOT use/refer to global variables in the functions. This will cause your code to fail during autograding because we execute functions independent of any content outside the function. However, you can (and should) call other functions when needed (of course, you can create new helper methods in classes).
- If you need to use **external libraries** that are not provided in lab, please ask for approval on Piazza.

Please answer the following questions:

### Q1. Binary Encoding

- (a) The class `ClassifierWithBinaryEncoding` implements the `get_dataset` method which takes in a list of files and returns the dataset with features  $X$  and labels  $y$ . Specifically, we use a binary encoding of the most common words as the feature set. Based on this description and the provided code, please describe in close detail how the feature set and the labels are constructed. Also, how large is the dataset?
- (b) The `Classifier` class has the method named `train_and_predict` which takes in the list of files (alternatively, a pre-processed dataset can also be passed), constructs the dataset (if not already provided with one), builds the train and test sets by *randomly* splitting the

entire dataset, trains a multi-class classifier model with the train set, evaluates the trained model with the test set, and finally returns the dataset and accuracy information. Complete the method (partial code is provided) and use it to compute the train/test accuracies of the logistic regression classifier that uses the baseline binary encoding defined in the `ClassifierWithBinaryEncoding` class. Note that you should pass the keyword arguments `'model_kwargs'` to the scikit-learn classification algorithm class you instantiate (i.e., `LogisticRegression`).

- (c) Try to improve the results of the baseline by improving (only) the feature set. You can use all the techniques covered in the IR lab to improve your features (e.g., stemming, lemmatization, lowercasing, stopwords; you can use NLTK for this purpose). Your code should be written in the provided method `get_dataset` of the `ClassifierWithImprovedBinaryEncoding` class (input and return values should be similar to the one returned by `ClassifierWithBinaryEncoding`).
- (d) Calculate the train and test accuracy of the model that uses the improved feature set. How did the results change?
- (e) Different train-test splits can lead to different results. In order to get a more robust estimation of the performance of your classifier, we want to calculate the mean and the 95% confidence interval on the accuracy of the classifier over a set of multiple runs with random splits. Notice that the function `train_test_split` takes an argument `random_state` that can be used to create different (random) splits by passing a random value to this argument. Please implement the method `get_performance_ci` that creates multiple random splits of your dataset (the argument `num_tests` will determine the number of splits to evaluate) and returns a tuple (`train_mean`, `train_ci_low`, `train_ci_high`, `test_mean`, `test_ci_low`, `test_ci_high`) that represents the mean and the low and high ends of the 95% confidence interval for both the training and test accuracy. We recommend you use `test_size=0.3` in `train_test_split`.

Note the following:

- To generate random numbers for the `random_state`, you can use, e.g., `random.randint(0,1000)` that generates a random integer in the range from 0 to 1000. Alternatively, you can generate multiple random integers at once by using `np.random.randint(0, 1000, size=num_tests)`. Yet another method is `np.random.permutation(range(1000))[:num_tests]`, which shuffles the integers from 0 to 999 and takes the first `num_tests` numbers. Make sure you don't use the same `random_state` more than once.
  - The code to calculate the mean and confidence interval given a list of accuracy results (the variables `train_results`, `test_results`) for different random splits is provided.
  - You should call `self.train_and_predict` (completed in Q1(b)) to get the accuracy results. Notice that you should not construct the dataset for every split since it will slow down your code; instead, just fix the dataset while varying how it is split.
- (f) Run the above function for 10 iterations (`num_tests=10`, see provided code). What do the average and 95% confidence intervals tell you? Are they more informative than a single trial? Yes or no, and why? [2 sentences.]
  - (g) Implement a function `create_cm` that produces a confusion matrix that is based on multiple random splits. Such matrix is created by summing the confusion matrices for the different

splits. Build a confusion matrix based on the results of 10 iterations (produced, as before, by calling `train_test_split` function with different random `random_state` values). Partial code is provided that includes the summation of the different confusion matrices.

Note also the following:

- Use the labels and predictions on the test set to compute the confusion matrices.
  - You should call `self.train_and_predict` (completed in Q1(b)) to get the prediction results. Notice that this method returns two dictionaries which contain the test predictions and labels.
- (h) Show the confusion matrix for 10 random splits (`num_tests=10`, see provided code). Are some classes more easily confused with others? Which ones and why? [2 sentences.]

## Q2. Number of Features

In this question, you will vary the number of words used as features and see how it affects the results. Please be careful to only use a single train/test split for this evaluation.

- (a) Calculate the train accuracy and the test accuracy when using the first  $p$  percent of the features (for a fixed ordering of features),  $p \in [10\%, 20\%, 40\%, 60\%, 80\%, 100\%]$ . The method `predict_with_varying_num_features` of the `ClassifierWithVaryingNumFeatures` class has partial code you need to complete. It returns a dataframe of the results.
- (b) Use the provided code to plot the results. Explain any trends you see (average over multiple trials if trends are not clear). [1 sentence.]

## Q3. Hyperparameter Tuning

- (a) Calculate the train and the test accuracy for different values of the regularization hyperparameter  $C$ :  $[10^{-3}, 10^{-2}, \dots, 10^0, \dots, 10^3]$ . The function `hyperparameter` has partial code you need to complete. It returns a dataframe of the results.

Note the following:

- A `Classifier` object (`classifier`) is passed to the `hyperparameter` function. Given a specific  $C$  value stored in `param` in the given code, figure out how you can use the `train_and_predict` method such that the logistic regression model uses the given hyperparameter for regularization.
- (b) Use the provided code to plot the results (we use a logarithmic  $x$  axis). Explain any trends you see (average over multiple trials if trends are not clear). [1 sentence.]

Note: In practice, you need to tune hyper-parameter on the validation set only, not the test set!

## Q4. Feature Encoding

In this question, you will evaluate the effect of using term-frequency (TF) encoding instead of a binary encoding on this dataset.

- (a) Implement a TF encoding in the method `get_dataset` of the `ClassifierWithTFEncoding` class. You should only use `num_words` number of terms as features. For better performance, use the same document processing method you used to build the improved binary feature set (as opposed to the baseline one we provided). Change the encoding from binary to TF.

- (b) Compare the two encodings by comparing their mean accuracies and 95% confidence intervals over 10 trials. Use the function `get_performance_ci` defined for the `Classifier` class. Which method performs better on this dataset? Why do you think this occurs? [1 sentence.]

### Q5. Comparison vs. Naive Bayes

In this question, you will compare the mean accuracy and 95% confidence interval of the logistic regression classifier to those of a naive Bayes (NB) classifier.

- (a) Implement a NB classifier evaluated over multiple random splits in the function `evaluate_nb`. You should instantiate an object of type `Classifier` that uses NB as the classification algorithm, instead of logistic regression. Then, you should use the `get_performance_ci` method to get the performance metrics. Use the encoding (binary or TF) you found to be better.
- (b) Compare the two classifiers by comparing their mean accuracies and 95% confidence intervals over 10 trials. Which method performs better on this dataset? Why do you think this occurs? [1 sentence.]

### Q6. Binary Logistic Regression

In this question you will build a binary logistic regression that is trained to classify the target `sci.med` vs. any other targets. Use the binary encoding of features for this question.

- (a) Implement the method `get_dataset` of the `BinaryLogisticClassifier` class. In this question, there are only two possible targets: 1 for `sci.med` and 0 for any other labels. Notice that you do not have to rebuild the feature set as long as it is already built and saved to the class variable `ClassifierWithImprovedBinaryEncoding._X`. Hence, this only leaves correctly computing the binary targets.
- (b) Using the method `get_performance_ci`, calculate the average accuracy and 95% confidence interval over 10 iterations (`num_tests=10`, see provided code). What do the average and 95% confidence intervals tell you? How do they compare to the multi-class logistic regression in Q1? [1 sentence.]

### Q7. Classification with Ranking

This question is the competitive portion of the assignment and should be answered in the `ml_assignment_Q7.ipynb` notebook. In this question you will again build a multi-class classification system. However, this time you will be interested in the ranked results of your system. You will be putting together a system to optimize the **average precision score**. The classification results will be ranked by the probability estimates from the classifier. To evaluate your system, use the provided `calculate_average_precision_at_k` function (where `k` is defaulted to number of testing documents). You may use any feature set and feature encoding you wish for this question, but you may only use the standard classifiers built into scikit-learn. Your model should be built and returned in the provided function `build_model_q7`. There is a runtime limit of 7 minutes for this question. Running the function `calculate_average_precision_at_k` with `build_model_q7()` as an input argument must run on Google Colab with a runtime less than 7 minutes. Models that exceed this limit will receive 0 marks for this question. For example usage, please see the Jupyter notebook.

- (a) Implement the `data_q7` function using your chosen feature set and feature encoding.

- (b) Implement the `build_model_q7` function with a machine learning model of your choice.
- (c) In markdown cells, provide:
- A clear and concise description of your chosen feature set and feature encoding
  - The name of the classifier you chose
  - Why you chose the feature set, feature encoding, and classifier you used
  - The final AP performance that your choices attained. We will verify this score by running the model returned by your `build_model_q7` function, and the data as returned by your `data_q7` function.

**Your score for Q4 will be calculated using the min-max normalization shown below:**

$$Q7\_score(AP_{yours}) = \max(0, \frac{AP_{yours} - 0.6}{AP_{max} - 0.6}) \quad (1)$$

Where  $AP_{yours}$  is your average precision score with a specific train test split and  $AP_{max}$  is 0.9 which we have chosen based on the results from previous years. You will not receive the credit for Q7 if your AP score is below 0.6 and you will receive the full credit(1 point) for Q7 if your score is greater than or equal to  $AP_{max} = 0.9$ .

**Note:** We provide a set of training and testing file names as an example of the final dataset your models will be tested against. However, the actual dataset we use in grading will be different. It is your responsibility to try to improve your classifier so that it will do well across all train-test splits. (Number of documents of each class will be balanced for all train-test splits. The ratio of train to test documents will be 2-to-1.)

**Hint:** Recall the average precision formula:

$$AP@k = \frac{1}{number\_of\_relevant\_documents\_at\_k} \sum_i^k Precision@i \cdot rel(i)$$

In this assignment, we evaluate based on full AP, where  $k$ = total number of testing documents. Additionally, since this is a multi-class classification ,  $number\_of\_relevant\_documents\_at\_k$  always equal to  $k$ .