

Coursework 1: Sentiment analysis from tweets

Due: Monday 10th November, 2025 (4.00pm GMT)

Instructions

This assignment is marked out of 100 and counts for 40% towards your final grade for the module.

Note: This is an individual coursework. You must attempt the questions yourself and provide strong evidence for understanding your method in both your notebooks and your report.

Objective: In this coursework, you will implement a Support Vector Machine classifier (or SVM) that classifies tweets according to their sentiment, i.e. positive or negative. You will derive features from the text of the tweets to build and train the classifier on part of the dataset. You will then test the accuracy of your classifier on an unseen portion of the dataset. Much of the background for this part is in Lecture 2 and Lab 2 on Text Classification, though you should use all of your knowledge across the module. The dataset is provided in the file ‘sentiment-dataset.tsv’.

How to submit: You will submit:

- Well documented code as part of **one or more IPython files** (.ipynb), with a recommended separate IPython file for Question 4. These IPython files will be built on the provided template file “NLP_Assignment_1.ipynb” as your starting point (Python 3.7+). Ensure your code runs from top to bottom without errors before submission. If you do use more than one IPython file, it must be clear which file corresponds to which questions – use informative names for the files.
- **A 2-page report** where you describe how you achieved each question briefly, with the relevant observations asked for below. The report is limited to 2 pages, and any additional pages will not be considered during marking. It is also not acceptable to put a note in the report saying that explanations continue as comments in the code in the IPython files. It is part of the exercise to provide a concise and clear 2-page report that summarises your contribution, results and analysis.

The template file contains some functions to load in the dataset, as well as for simple data input and pre-processing, but there are some missing parts that you are going to fill in as per the questions below.

Question 1 (10 marks)

Simple feature extraction. The first step is to implement the `to_feature_vector()` function. Given a preprocessed text (that is, a list of tokens), it will return a Python dictionary that has as its keys the tokens/words, and for the values a weight for each of those tokens in the preprocessed texts. The weight could simply be 1 for each word, or the number of occurrences of a token in the preprocessed text (i.e. a bag-of-words representation), or it could give more weight to specific words. While building up this feature vector, you may want to incrementally build up the global `global_feature_dict`, which should be a list or dictionary that keeps track of all the tokens/feature names which appear in the whole dataset. While a global feature dictionary is not strictly required for this coursework, it will help you understand which features (and how many!) you are using to train your classifier and can help understand possible performance issues you encounter on the way.

Hint: you can start by using binary feature values; 1 if the feature is present, by default the sklearn learn vectorization function will give it 0 if it's not.

Question 2 (20 marks)

Cross-validation on training data. Using the `load_data()` function already present in the template file, you are now ready to process the texts for sentiment analysis. In order to train a good classifier, finish the implementation of the `cross_validate()` function to do a 10-fold cross validation on

the training data (leave the test data split of 20% alone for now). Make use of the given functions `train_classifier()` and `predict_labels()` to do the cross-validation. Make sure that your program stores the precision, recall, f1 score, and accuracy of your classifier in a variable `cv_results`, which should contain average scores for all folds and be returned by this function.

Hint: the package `sklearn.metrics` contains many utilities for evaluation metrics - you could try `precision`, `recall`, `f1score`, `support` to start with.

Question 3 (20 marks)

Error analysis. Look at the performance of the classes using a confusion matrix through the method provided `confusion_matrix_heatmap()` to see what the balance of false positives and false negatives is for the positive and negative labels. Carry out an error analysis on a simple train-test split of the training data (e.g. the first fold from your cross-validation function). For this you should print out (or better, print to file) all the false positives and false negatives for the positive label to try to understand why the classifier is not getting these correct and write in your report some observations and examples of where it is getting confused.

Hint: see Lab 2.

Question 4 (50 marks)

Question 4 will consist in improving the existing sentiment analysis model, and marks will be awarded according to different objectives, with the 50 marks broken down as indicated next.

Optimising pre-processing and feature extraction (30 marks). Now that you have performance scores of your classifier and have done some initial error analysis, think of ways to improve this performance score. Some ideas as to how to do this:

- Improve the preprocessing. Which tokens might you want to throw out or preserve?
- What about punctuation? Do not forget normalisation, lemmatising, stop word removal - what aspects of this might be useful?
- Think about the features: what could you use other than unigram tokens? It may be useful to look beyond single words to combinations of words or characters. Also the feature weighting scheme: what could you do other than using binary values?
- You could add extra stylistic features like the number of words per sentence.
- You could consider playing with the parameters of the SVM (cost parameter? per-class weighting?)
- You could do some feature selection, limiting the numbers of features through different controls on e.g. the vocabulary.
- You could use external resources like the opinion lexicon available at <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#lexicon>.

Report what methods you tried and what the effect was on the classifier performance in your report and evidence the exploration in your notebook. Marks will be given depending on how many of the above suggestions you have tested.

If you implement any of the suggestions above, but for any reason it doesn't lead to improved performance so you choose not to keep it as part of your final model, then: (i) don't remove this implementation from your code, just comment it out so that we can see you implemented it, and (ii) briefly discuss it in the report.

Improve model performance (20 marks). Marks will be given depending on the performance of your final model within a range between a baseline model and an upper-bound score, measured with the F1 score. We will use 0.85 as the baseline F1 score achieved by a simple pre-processing / feature extraction model, and the upper-bound score will be determined by the best-performing model among all coursework submissions. As a reference for the upper-bound score, we anticipate that the best-performing model will be slightly over the 0.9 F1 score mark.