

Natural Language Processing

Project 1.4: Spam Filtering

Organization

- Please read the project description carefully. All the details are presented in this document.
- Submit your solution no later than July 4th, 23:55. Submit your solution using the corresponding ISIS section, where you have downloaded this document.
- You can make unlimited re-submissions until the deadline, but the latest received submission will be graded only.
- Submit your solution as a zip file with the following name format: 'NLP_project1_4_[Your Names separated by underscores].zip' containing your codes and report files.
- For those who decided to work on the project in a group, one submission of the solution would be enough.
- For each task, you see a percent that show the task's score in the whole project.
- Please use Python 3 for coding.
 - You are free to use any module/library for different tasks.
- Please put comments on your code, wherever you think it would help to improve the readability and understandability of your code.
 - You can submit your codes as a .py file or Jupyter notebook
- You should submit a short report, describing your approach to solve different tasks and also provide the obtained results (e.g., the evaluation result of your models).
 - For the report please use ACM proceedings template from here (<https://www.acm.org/publications/proceedings-template>). The office word and latex templates are provided in the page.
 - The report should be between 3 to 5 pages.
- There is a bonus task in the end of the document, as the name implies it is not a mandatory task, but you can take your time and solve it to get more scores.

Plagiarism Statement

Your project, includes the report and the code, will be checked against other submissions in the class. We trust you all to submit your own work only. Copying someone else's code and report and submit it with minor changes, will be treated as plagiarism.

If you have any further questions, please write to: salar.mohtaj@tu-berlin.de

Introduction:

Spam is any kind of unwanted, unsolicited digital communication that gets sent out in bulk. Often spam is sent via email, but it can also be distributed via text messages, phone calls, or social media. Natural Language Processing has been very promising in detecting spam messages in recent years. Nowadays all bid email service providers use NLP to automatically detect spam emails.

Data:

In this project you receive a zip file contains a spam detection dataset. The data includes 5730 instances of spam and normal (ham) email texts. Please apply the following tasks on the dataset.

Tasks:

Task 1: Extract insights from data (15%)

For this task you should **extract some insights** (i.e., some statistics and graphs) from the provided data. Here are some ideas:

- It could be a graph compared the length of spam and normal messages.
- Number of unique words in each category and also in the whole dataset
- ...

Please don't limit yourself to these two examples and try to summarize the data with numbers and graphs. Please highlight some of the most important findings in your report.

Task 2: Pre-processing (20%)

In this task, first, apply **all the necessary** pre-processing steps that you think they would help to better prepare your data for the next steps. You don't have to apply all the pre-processing tasks which are covered in the course. Regarding the report, you should briefly mention it in your report that why you decided to apply the chosen pre-processing steps (and why not the others).

Task 3: Text classification (50%)

In this task you should do the following sub-tasks:

- Split data into train and test sets. Use 20% of data as the test set.
- Train a **naïve Bayes model** on the training part and test it, using the test set.
 - o Compare the impact of different vectorization models (e.g., count vectorizer, TF-IDF and ...) on the final performance of your naïve Bayes model.
- Train a **feed forward neural network** model and report its performance (F1 score) on detecting spam messages on test data.
 - o Again, compare the impact of different vectorization approaches on the final performance of your model.

Please report all the achieved results with either models in your report document. Moreover, describe the hyper-parameters of your neural network model in the report.

Task 4: PMI based word similarity (15%)

Second, convert words to vectors using **PMI** (window size = 3). For 20 random words in the data, provide most similar words based on the computed word-word pointwise mutual information matrix. Please highlight the 20 randomly chosen terms and their most similar words in your report.

Bonus Task: Textual similarity (+20%)

As the bonus task, you should **fine-tune two pre-trained models** on the provided spam detection dataset. You can choose between the wide range of pre-trained models (e.g., BERT, RoBERTa, and so on). It is very important to be careful about the hyper-parameters and fine-tune the models as accurate as possible without too much overwriting of the current weights. Please describe your selected models and the training process in details in the report.