

EL-GY-9133 Machine Learning for Cyber-Security
Lab 1: E-mail Spam Filtering
Release Date: 10/11/2023; Due Date: Midnight, 10/28/2023

Overview

In this lab, you will design an e-mail spam filter using a Naïve Bayes and SVM based classification on the ling-spam dataset. You will explore the impact of feature selection and compare the performance of different variants of an NB classifier and also implement your own SVM based classifier. (Note: You may use the scikitt learn classifiers to only compare the accuracy of their model to yours).

Dataset

The ling-spam corpus contains e-mails from the Linguist mailing list categorized as either legitimate or spam emails. The corpus is divided into four sub-folders that contain the same emails that are pre-processed with/without lemmatization and with/without stop-word removal. The e-mails in each sub-folder partitioned into 10 "folds." In this lab, we will use the first 9 folds from the ling-spam corpus as training data, and the 10th fold as test data.

What You Have to Do

You will implement your e-mail spam filters in Python. You are free to use any Python libraries that are relevant to the problem.

- Download the ling-spam dataset from http://www.aueb.gr/users/ion/data/lingspam_public.tar.gz

Please use the "lingspam_public01" corpus with both lemmatization and stop-word enabled (under the lemm_stop folder).

- Your first goal is to perform feature selection using the information gain (IG) metric. From the training data, select the top-N features for $N = \{10, 100, 1000\}$. Note that feature selection based on the IG metric only accounts for the occurrence of (and not frequency with which terms appear) in the dataset.
- Next, implement the following classifiers:
 - Bernoulli NB classifier with binary features;
 - Multinomial NB with binary features; and
 - Multinomial NB with term frequency (TF) features.
- For each of the three classifiers above and for $N = \{10, 100, 1000\}$ report the spam precision and spam recall. Spam precision is defined as the fraction of true spam e-mails among all e-mails predicted as spam, and spam recall is defined as fraction of true spam e-mails predicted as spam. Also report the latency of each model. Design a Support Vector Machine (SVM) based spam filter. This problem is open ended: for instance, you can choose to use either BF or TF and any feature selection method. Note that you should NOT use the test dataset in picking the hyper-parameters of your spam filter; instead use cross-validation on the training dataset.

What to Submit

Your Python code in the form of a Google Colab notebook. Please also include a PDF of your colab notebook. Your Colab notebook should print:

- a list of the top-10 words identified from Part (1) above, and
- a list of spam precision and spam recall values for each of the three classifiers for $N = \{10, 100, 1000\}$. That is, your list should have 9 rows, one for each classifier and N combination.
- For the SVM based spam filter, please describe your methodology, i.e., what kind of features you used, how many features you used and how you selected them, the parameters of the SVM and finally the precision and recall on the test dataset.
- For Your Neural network model, please report on the model design, input parameters, number of layers, number of neurons in each layer and finally the accuracy and latency of the model.
- For the adversarial attack, report the False Negative rate of the baseline NB classifier before and after the attacker's modifications to test emails. Also, report the average "cost" of the attacker's modifications, averaged over all spam emails in the test set. Finally, report the False Negative and False Positive rate of the updated NB classifier.