

[Lab] Filtering spam messages using Naïve Bayes

Jae Yun JUN KIM*

Due: Before the end of today lab session

Evaluation: Show your Python code and results to the Professor and answer his questions on the code and on the results.

Remark:

- Only groups of two or three people accepted (preferably three). Forbidden groups of fewer or larger number of people.
 - Show your lab and explain it to the Professor before the end of today lab session.
 - No plagiarism. If plagiarism happens, both the “lender” and the “borrower” will have a zero.
 - Code yourself from scratch **following the theory explained in class**. No lab work will be considered if you solve the problem using any ML library.
 - Do thoroughly all the demanded tasks.
 - Study the theory for the questions.
 - There is NO make-up session.
-

1 Tasks

1. Divide the data in two groups: training and test examples.
2. Parse both the training and test examples to generate both the spam and ham data sets.
3. Generate a dictionary from the training data.
4. Extract features from both the training data and test data.
5. Implement the Naïve Bayes from scratch, and fit it to the training data.
6. Make predictions for the test data.
7. Measure the spam-filtering performance for each approach through the confusion matrix, accuracy, precision, recall, and F1 score.
8. Plot a graph with true positive rate on the vertical axis and with false positive rate on the horizontal axis.
9. Discuss your results.

*ECE Paris Graduate School of Engineering, 10 Rue Sextius Michel 75015 Paris, France; jae-yun.jun-kim@ece.fr