# Document Classification: 20 Newsgroup

## INT3406 21 - Group 3

**Group 3**
Pham Truong Giang - 1802xxxx
Nguyen Manh Dung - 18020370
Nguyen Phuc Hai - 1802xxxx
Le Bang Giang - 1802xxxx

### Abstract

Your abstract should motivate the problem, describe your goals, and highlight your main findings. Given that your project is still in progress, it is okay if your findings are what you are still working on.

# 4 Classification

## 4.1 Linear Model

### 4.1.1 Naive Bayes

### 4.1.2 Logistic Regression (LR)

### 4.1.3 Ridge Classification

### 4.1.4 Perceptron

### 4.1.5 Passive-Aggressive

## 4.2 Non-parametric

### 4.2.1 K-nearest neighbor (KNN)

### 4.2.2 Support Vector Machine (SVM)

### 4.2.3 Linear Support Vector Machine (LinearSVC)

## 4.3 Tree-based Classifiers

### 4.3.1 K-nearest neighbor (KNN)

### 4.3.2 Support Vector Machine (SVM)

### 4.3.3 Linear Support Vector Machine (LinearSVC)

## 4.4 Graphical Classification

### 4.4.1 Conditional Random Fields (CRFs)

### 4.4.2 ...

## 4.5 Neural Network

# 5 Summary

Abstract

# References