# CS534 Implementation Assignment 2: Naive Bayes

Amit Bawaskar, Michael Lam
EECS, Oregon State University

May 1, 2013

**Abstract**

In this assignment, we implemented the Naive Bayes classifier with the Bernoulli model and Multinomial model, and compared their performance.

# 1 Introduction

We implemented the Naive Bayes classifier to solve a document classification problem on the 20-newsgroup data set. Two models were implemented and compared for performance: Bernoulli model and Multinomial model.

# 2 Naive Bayes

Paragraph about Naive Bayes.

## 2.1 Model

Naive Bayes assumes that the features are independent.

## 2.2 Inference

Inference is performed by using Bayes rule with the learned likelihood and prior probabilities, and using decision theory to select the class that maximizes the posterior probability.

## 2.3  Learning

Learning the likelihood probabilities for each feature and class is done by maximum a posteriori estimation. Equivalently this amounts to applying Laplace smoothing to the maximum likelihood estimator. The prior probabilities for each class is the maximum likelihood estimator for it.

## 2.4  Implementation Details

In this project we operated with the log of probabilities in order to avoid underflow issues. That is, for every multiplication and division operation, we instead used addition and subtraction of the log of the operands. We also stored the log of probabilities. For decision theory, we simply selected the class that maximizes the log of the posterior probability since the log function is a monotonically increasing function. Thus we operated with the log of probabilities for every circumstance.

We also applied Laplace smoothing to the likelihood probability of each feature and class in order to assign a default prior to words that have not been encountered.

# 3  Bernoulli Model

Paragraph about Bernoulli model.

The overall test accuracy for the Bernoulli model is **0.772152**. Figure 1 shows the confusion matrix.

# 4  Multinomial Model

Paragraph about Multinomial model.

Report overall testing accuracy.

Report the confusion matrix.

# 5  Heuristics for Reducing Vocabulary Size

We explored several heuristic strategies for reducing the vocabulary size in hopes of improving classification accuracy.
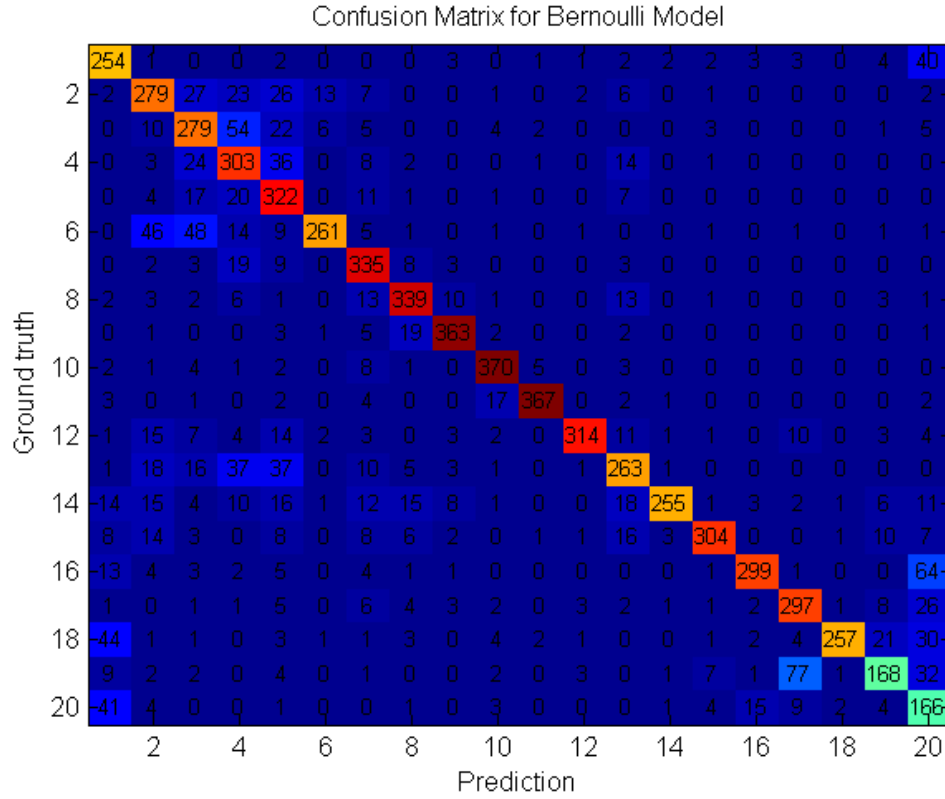
Figure 1: Confusion matrix for the Bernoulli model.

We tried eliminating words in the dictionary that have length less than or equal to a certain threshold value. We explored thresholds from 0 to 5, meaning keeping all words to keeping only words with length greater than 5. The results in Figure 2 indicate that the classification accuracy increased negligibly when eliminating words up to 2 letters long, but decreased noticeably as the threshold increases after 2 letters.
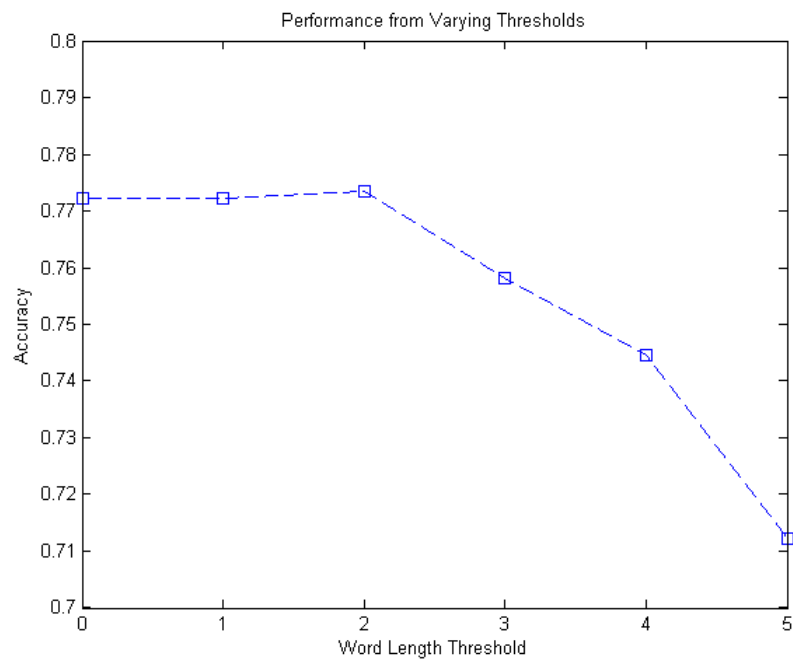
Figure 2: Classification accuracy using the Bernoulli model and varying word length thresholds. Words of length less than or equal to the threshold value are omitted in the vocabulary for learning and inference.