



Utrecht University

Data Mining: Assignment 2

**TEXT CLASSIFICATION FOR THE DETECTION OF
OPINION SPAM**

Ana Borovac
6584446

Argyro (Iro) Sfoungari
6528015

October, 2018

Contents

1	Introduction	3
2	Data	3
2.1	Data preprocessing	3
3	Methods	3
3.1	Naive Bayes	4
3.2	Logistic regression	4
3.3	Classification tree	4
3.4	Random forests	4
4	Results	4
5	Analysis	4
6	Conclusion	4

Abstract

1 Introduction

2 Data

Our data consists of 400 negative deceptive and 400 negative truthful hotel reviews that have been collected by Myble Ott and others

2.1 Data preprocessing

Before modeling we did some data preprocessing in order to get better models for prediction if a review is fake or real. First step in data preprocessing was to remove punctuation marks, after that we made every letter lower case, we also removed stopwords, numbers and excess whitespace. Next, we created a test and a training set. Because we wanted to use cross validation, we divided our data into 5 folds (each of size 160 samples - 80 fake reviews and 80 truthful reviews). So, 4 of the folds represented a training set, the remaining 5th fold was a test set. At that moment every unique word from the training reviews was a feature. A number of features was large therefore we removed the words as features which occur less than 5 % of training documents. Next, we created a new training set which consists all the features from the previous one and bigrams. We again removed bigrams that occur in less than 5 % of training documents. At the end we had two training sets ready to be used.

3 Methods

We used different classifiers to model our task:

- naive Bayes (generative linear classifier) – subsection 3.1,
- logistic regression (discriminative linear classifier) – subsection 3.2,
- classification tree (flexible classifier) – subsection 3.3,
- random forests (flexible classifier) – subsection 3.4.

3.1 Naive Bayes

For naive Bayes we used the code from the lectures ()

3.2 Logistic regression

3.3 Classification tree

3.4 Random forests

4 Results

5 Analysis

6 Conclusion