

Exercise 1 Classification

e12045110 Maria de Ronde
e12045110 Quentin Andre
e11921655 Fabian Holzberger

April 22, 2021

Introduction

Applied Algorithms

Performance Metrics

Amazon Reviews Dataset

Dataset Description

Pre-Processing

Parameter-Tuning

Performance-Analysis

Congressional Voting Dataset

Dataset Description

Pre-Processing

Parameter-Tuning

Performance-Analysis

Email Spam Dataset([link to dataset](#))

Dataset Description

The task of the email spam dataset is to predict if an email is spam or not. The link above contains three datasets from which we choose two, namely the `lingSpam.csv` and `enormSpamSubset.csv` for our project, since they have no missing values and the same layout. The dataset contains 12604 emails where 43.11 Every email has a binary target-label assigned, such that a 0 marks non-spam and a 1 marks spam emails. In figure 2 the distribution of the characters per email is shown. We see that most emails have a lenght in the range of 100 to 10.000 characters.

Index	Body	Label
100	Subject: inexpensive online medication here pummel wah springtail cutler bodyguard we ship quality medications overnight to your door !...	1
6006	Subject: organizational changes we are pleased to announce the following organizational changes : enron global assets and services in order to increase senior management focus on our international businesses...	0

Figure 1: Structure of the Email-Spam Dataset

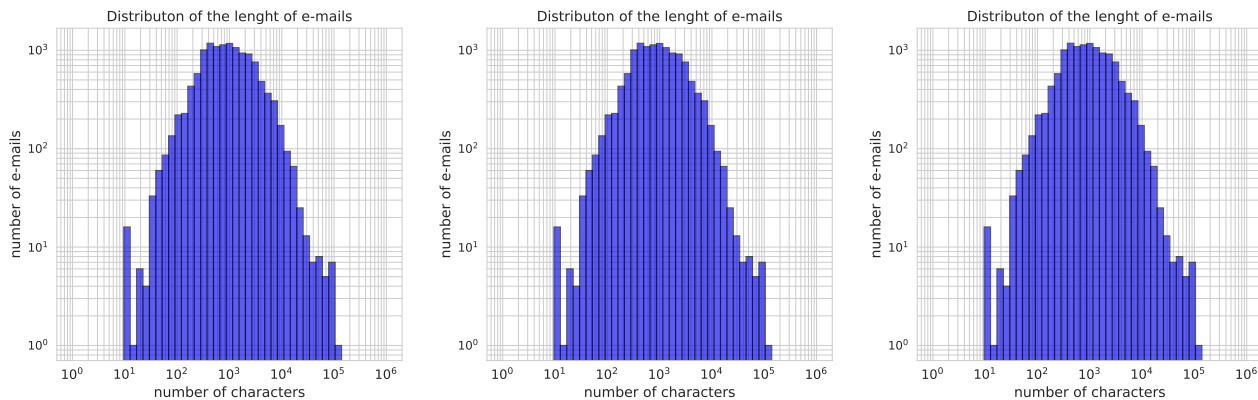


Figure 2: left: Distribution of e-mail lenghts

We apply the Bag of Words feature extractor to each email. The algorithm converts every email to a vector $v \in \mathbb{N}^+^N$ of intergers. First we create a list of all words and count their occurences in all emails. Then we take the N most common words and count the occurences of the most common words in every email ti get v . Before applying the Bag of word extractor we pre-process emails by the following steps:

Pre-Processing

Parameter-Tuning

Performance-Analysis

Bridges Dataset([link to dataset](#))

Dataset Description

Pre-Processing

Parameter-Tuning

Performance-Analysis

Conclusion