

# Text Technologies for Data Science Assignment 3

S1140740

November 2, 2014

## Introduction

The purpose of this report is to provide an overview of decisions made in designing a plagiarism detection tool and cover different types of plagiarisms detected.

## Plagiarism types

*Type 1* Plagiarisms of this type are exact duplicates, e.g. documents t104 and t4172 which are classified as type 1 are exactly the same, word for word.

*Type 2* Plagiarisms of this type are near duplicates, e.g. documents t1088 and t5015 which are classified as type 2 are almost the same, except document t1088 is missing the word *room*.

## Plagiarism Detection

The Simhash algorithm together with the md5 hashing function were used to detect duplicates. For each document an md5 hash is computed and the document is then tokenized on white-space and non-alphanumeric characters with stop words removed.<sup>1</sup>

Furthermore, a 128-bit fingerprint is computed using the Simhash algorithm with md5 as a hashing function and term frequencies as weights. This fingerprint is then split into  $L$  chunks of size  $K$  and the  $i$ -th chunk is used as a key of a bucket containing documents having the same chunk in the  $i$ -th hash-table.

<sup>1</sup> List of English stop words obtained from <https://github.com/Alir3z4/stop-words> by Alireza Savand. No changes made.

### *Type 1 detection*

To detect plagiarisms of type 1, an md5 hash of a document is compared to every document in the same buckets. If these hashes are identical then we flag the documents as exact duplicates (type 1). Note, that two stories that are identical are guaranteed to end up in the same buckets as their vector space representation and thus their Simhash fingerprints are the same.

100% precision and recall compared to the file type1.truth are achieved using this technique.

*Type 2 detection*

To detect plagiarisms of type 2, a cosine similarity measure is computed for a document and all the other documents in the same buckets. If this measure is above a certain threshold, in this case **0.8 WHY?** we flag the documents as near duplicates.

Experiments were ran on different values of  $L$ . With  $L$  equal to 8 ( $K = 16$ ), 100% precision and recall compared to the file `type2.truth` are achieved.

*Type 3 detection*

To detect plagiarisms of type 3, all areas with high number densities are extracted from the data file. Then the same algorithms as for plagiarisms of types 1 and 2 are used to detect duplicates of type 3. This way, 13 duplicates were detected in the dataset `data.finn`.

*Notes*