



# IMapBook Collaborative Discussion Classification

Uroš Polanc

## Abstract

Automatic text classification has been considered as a vital method to manage and process a vast amount of documents in digital form, which is widespread and continuously increasing. That brought a wave of new tools for student education, such as the IMapBook platform. It gives students a way to communicate and discuss the books they are reading. In our assignment, we will try to classify those discussions based on the book it's relevant to, and the type of message.

## Keywords

IMapBook, Text Classification, TF-IDF, Word Cloud

Advisors : Slavko Žitnik

## Introduction

With the rapid development of the internet and big data, digital information is increasing at a high rate. With this technologies and applications such as IMapBook [1], a web-based application that allows for discussions on reading materials as well as interactive games, also started emerging.

In our assignment, we will look at some IMapBook collaborative discussions on different books. For that, we need to look firstly at some vital information, that we extracted from the discussions and the book texts. Firstly we will create two different message classifications, based on :

- book ID (3-class),
- message type (16-class).

First, we applied some fixes to the manual classifications in the dataset, since not all fields were correctly filled. The errors such as excess whitespace, lower and upper case, etc... were removed. Classifications where there were multiple classes, were simplified into one class, to simplify the classification process. The final fixed distribution is not intraclass uniform as can be seen in figures 1 and 2.

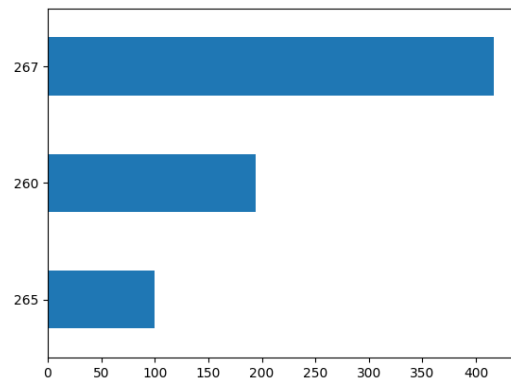


Figure 1. The distribution of book IDs in the dataset.

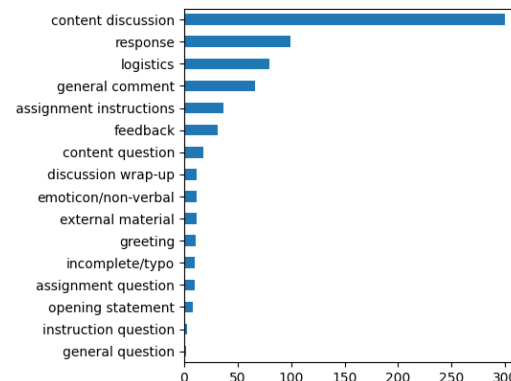


Figure 2. The distribution of message types in the dataset.

