# VirginiaTech
## Invent the Future

CS 5604 Information Storage and Retrieval
Spring 2016

Interim Report 2

# Text Classification

Team Members:

Hossameldin Shahin <hshahin@vt.edu>
Matthew Bock <mattb93@vt.edu>
Michael Cantrell <mcantrell@vt.edu>

Project Advisor:

Prof. Edward A. Fox

3/31/2016
Virginia Tech, Blacksburg

## Abstract

In the grand scheme of a large Information Retrieval project, the work of our team will be that of performing text classification on both the tweet collections and their associated webpages. In order to accomplish this task, we will seek to complete three primary goals. We will begin by performing research to determine the best way to extract information that will be used to represent a given document. Following that, we will work to determine the best method to select features and then construct feature vectors. Our final goal is to use the information gathered previously to build an effective way to classify each document in the tweet and webpage collections. We intend for these classifiers to be built with consideration of the ontology developed for the IDEAL project.

The team assigned to perform this classification work last year researched various methods and tools that could be useful in accomplishing the goals we have set forth. Last year's team developed a system that was able to accomplish similar goals to those we have set forth with a promising degree of success. Our goal for this year is to improve upon their successes using new technologies such as Apache Spark. It is our hope that this will allow us to build a more optimized system capable of working with larger collections of tweets and webpages.

# Todo list

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

The classification team's goal is to take collections of tweets and webpages and classify them based on their relevance to given classes or topics. Our team fits into the grand scheme of the project by working between the Collection Management team and the Solr team. The Collection Management team will be responsible for taking the raw tweet and webpage data and filtering out any obvious spam, vulgarity, or otherwise unreadable and unwanted content. The Collection Management team will write the cleaned data into a table in our HBase instance. Once this is done, we will attempt to classify each document's relevance to a specific category. We will explore several methods for accomplishing this, starting with the methodology laid out by last year's Classification team [2]. Once we have classified the data, we will be able to pass it along to the Solr team by writing the newly classified data back into the primary HBase table as a column family. The Solr team will then be able to use the column family in the indexing of all the tweet and webpage data. Solr will then provide the indexes for the data to the Front End team, allowing anyone to make use to the systems we are creating.

This paper represents the third interim report from the Classification team working as a part of the larger scale project for CS5604 Information Storage and Retrieval. The current state of the work is a partial draft of the final report; therefore expect further additions and details in some chapters as the course progresses.

We begin by documenting our understanding of some of the pertinent literature, namely the course textbook and the Classification team report from last year; this can be found in Chapter 2, our literature review. Chapter 3 is the primary section of the document and includes our discussion of the project requirements, an outline for our design for the classification portion of the larger project, and finally a breakdown of our current progress and future plans for the text classification project.

These chapters are then followed by a User Manual in Chapter 4, where we will discuss details in which a user of our methods and programs would have interest. We then include a separate Developer Manual in Chapter 5, which will document and include details of the code base so that it might be extended and leveraged by other developers.

Following this we include our conclusions about the state of the project thus far and present our thoughts for future work in Chapter 8.

# Chapter 2

# Literature Review

## 2.1 Textbook

The textbook [4] introduces the classification problem. That is, given a set of classes, the goal is to determine what subset of those classes a document may belong to. To that end, the textbook describes a number of methodologies for selecting features. These features are then used by one of the classification methods discussed. The primary methods discussed are Naïve Bayes, Vector Space Classification, and the Support Vector Machine.

## 2.2 Papers

> could better expand on the state of the art

The primary paper that has been used as a reference is the final report of the Classification team from last year [2]. We have read through this report to understand the progress that the Classification team made last year as well as using the paper to gain an understanding of the task and interactions of the systems that we are using to perform our work. At a high level, the paper described a methodology employed by the team in which they were able to apply the Naïve Bayes method to classify sets of data. The team used Apache Mahout machine learning technology to generate Naïve Bayes classifiers to predict for new data. The biggest difference from last year's work to this year's work is that we will be primarily using Apache Spark, and so while we will be able to reference their work, we are using entirely different technologies and will need to modify our approach accordingly.

Since we were advised to look at Frequent Pattern Mining (FPM) we have looked at a paper by Han et al. [3]. This paper presents a novel frequent pattern tree (FP-tree) structure and FP-tree based mining method called FP-growth. This method allows for mining the complete set of frequent patterns by pattern fragment growth. In this paper they also demonstrate that their new method is an order of magnitude faster than the typical apriori algorithm. The paper goes on to compare the performance against other popular data mining algorithms and discusses the use of the algorithm on some large industrial databases.

We make use of the FP-growth algorithm and data structure in our classification work. Our use for this will

> Give more justification for using FPM besides the advice and one paper

be described later in the report.

# Chapter 3

# Requirements, Design, And Implementation

## 3.1   Requirements

Based upon group discussions in a classroom setting, it was decided what information that the Collection Management team would provide to the other teams that need to work with cleaned tweet data. A full specification of these decisions can be found by viewing the Collection Management team's report; however we will briefly discuss the points that were relevant to us.

Given the cleaned tweet data from Collection Management, our team will then be able to perform the methodologies we describe later to classify whether a document is relevant or non-relevant to a given class. A detailed layout of the project with our interactions highlighted is provided by Figure 3.1.

We will then place our classification information in the HBase datastore to mark the relevance of each document. The Solr team will then index this class information that we provide, allowing for more robust query results and more useful facets to be created.

## 3.2   Design

Our current design is based primarily off of recommendations from the GRAs assisting the class. We have also taken substantial input from last year's Classification team [2] and the Professor.

We have designed our current solution around pulling training data from and testing on a small collection of around 100,000 tweets. This was originally going to be performed on the small collection that was assigned to our team, `#germanwings`. However, due to some changes and discussion among the other teams, we have decided to continue with designing and testing our solution using a cleaned dataset provided by the Collection Management team. However this dataset was not made available until after our current solution was mostly implemented using our original small dataset. Therefore for the rest of this document we will be using the small dataset `z_602`, `#germanwings`. All future work will be done using cleaned data from the Collection Management team.

**Figure 3.1:** Layout of the project teams, as provided by Professor Fox and his GRAs.

Provide some rationale on the choice of using FPM

## 3.3 Implementation

Our implementation can be easily laid out in a step by step workflow, with only one step requiring true manual input and evaluation from a person. Figure 3.2 illustrates the general flow of data that we have implemented thus far. Below we will discuss each step in detail.

Our methodology primarily revolves around building a training set. This training set will be used for training a machine learning model, a.k.a. a classifier.

### 3.3.1 Environment Set–Up

Our group decided to avoid working directly on the cluster, granting us full administrative rights to a machine, which allowed us to set up and test different configurations and algorithms beyond what might currently be supported by the Hadoop cluster being used for class.

A prime example, and the primary motivating factor for our choice, was our decision to use Frequent Pattern Mining (FPM) in the construction of our training sets. The cluster provided for the class was running Spark version 1.2.0, but FPM is not supported by Spark until version 1.5.0. Since FPM was a methodology

**Collection Management Team**

Collection management team starts with raw data and transforms it into a more useful form

Cleaned data

**Frequent Pattern Mining**
Determine keywords based on frequent patterns across data set

Keywords

**Build training data**
Use frequent patterns to select positive samples, randomly select negative samples from the rest

Training set

**Build classifier**
Train a classifier on Apache spark with the training data

Classifier

Classified data is put back into HBase for use by other teams

**HBase**

Classified data

**Classify data set**
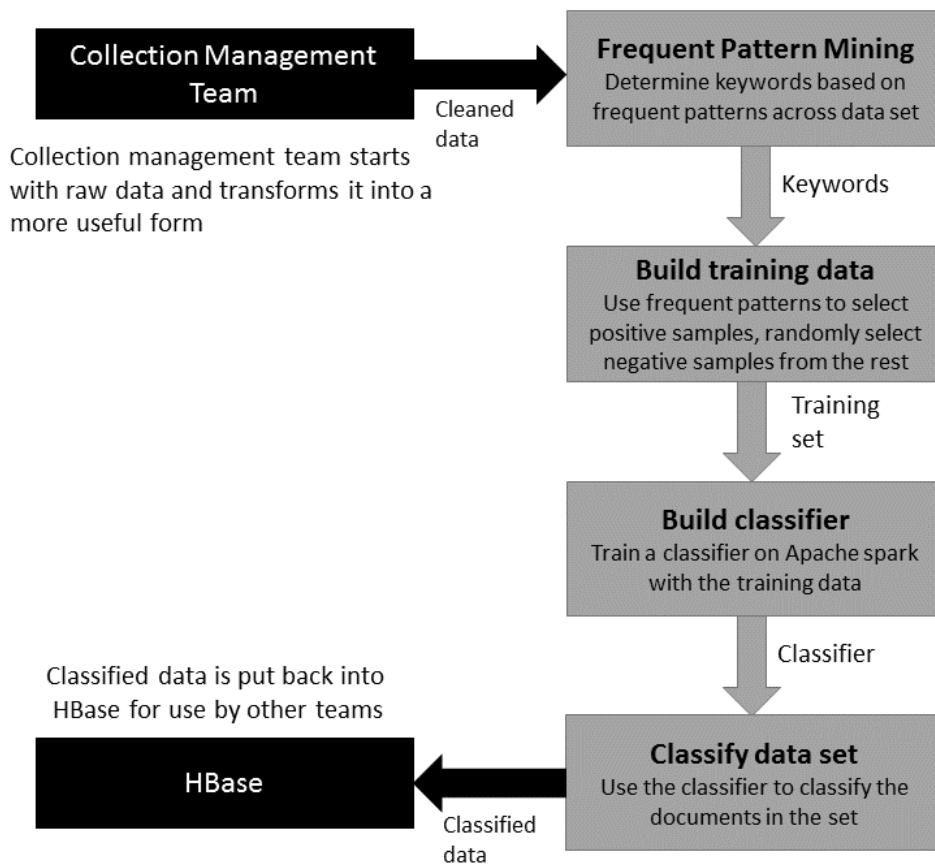Use the classifier to classify the documents in the set

**Figure 3.2:** High level overview of the flow of data through the classification system.

suggested by the GRAs, and they assured us that they could upgrade the cluster to a more current version of Spark, we decided to proceed with this method.

Due to time constraints, we created a Debian Virtual Machine with sufficient specs to begin writing and testing our methodology. This machine is hosted on a VMware vSphere cluster that belongs to the Virginia Tech IT Security Office/Lab. Since one of the group members works in this lab, we were able to provision this VM. It has been assigned 32 GB of RAM, 16 processor cores, and 40 GB of storage. This has proven to be more than adequate for the testing that we have performed on our small data set.

The Hadoop cluster has now been upgraded to Spark 1.5.0 allowing us to transition our code and methods over and work directly with the HBase instance. As this was a recent change, we are still in the process of this transition and cannot report on what effect it will have, ie whether there are any issues encountered that break our implementation.

### 3.3.2 Building the Training Data

In order to begin the classification process we have to prepare a set of training data to use for the machine learning classification algorithm. In order to do this, we assume that we are working with data that has been

cleaned of profanity and non-printable characters.

For our methodology we then need to take the content of each tweet and webpage as a string and tokenize it. This means that we will remove any duplicate words and have the result be a set of the vocabulary in the tweet (or webpage). At this stage we also remove any stop words that are in the vocabulary to make sure our algorithms in the next phase are not skewed by taking stop words into account. During this phase we also strip the # from any hashtags that are present in the tweet.

In our initial design and testing we did not yet have access to cleaned tweets and webpages from the Collection Management team, so we worked primarily to build the methodology that would be used once those resources became available. Therefore, some of the steps mentioned previously, such as the removal of the stop words and the removal of the # from hashtags are unnecessary when working with the data provided by Collection Management in HBase.

The next step in our algorithm involves the use of Frequent Pattern Mining (FPM) to determine the most frequently used patterns of words in the text of a series of tweets or webpages. FPM looks at the set of existing vocabulary for each text and determines which tokens appear together within a tweet's vocabulary most often.

This is the stage where manual intervention in necessary. We look at a file containing a sorted list of all the frequent pattern;, from this we choose a set of features to train on. In our attempts at training a classifier for our small data collection, we chose to select a frequent pattern of four words. This was essentially an arbitrary choice, but we believe that it strikes a good balance between being too specific and missing some relevant tweets and being too broad and pulling in a number of irrelevant tweets.

At this stage we need to create positive and negative training samples. To do this, we pull a random subset of tweets that contain our frequent pattern and classify those as positive samples.

> Have you checked whether the results of using FPM to determine labeling of tweets as positive or negative is correct? That is, did you manually check (some) of the label assignments? Will you use the same approach for webpages too?

For our negative sample we pull a random subset of tweets that do not contain our frequent pattern.

### 3.3.3 Training the Classifier

We then use the sets of positive and negative samples to train a classifier. We are currently using a logistic regression classifier in our implementation, primarily due to its ease of implementation in Spark.

FPM allows us to develop a training set by filtering the tweets and webpages down to some subset of those that do and do not contain the selected frequent patterns. We take these subsets as the positive training data and the negative training data.

At this point we then feed the selected documents into the classifier; this step uses all the words in each document, not just the words used for FPM. As we progress we intend to try using several different classifiers to determine which one provides us with the best results.

### 3.3.4   Predicting the Class

After training the classifier, we apply the classifier to all the tweets across our small data collection. This results in a binary scoring for each tweet as relevant or non-relevant to the collection based upon the training data we selected.

### 3.3.5   Evaluating the Classifier

We can evaluate the accuracy of our model by judging how well it classifies some of the data that could have been in our positive or negative samples. This is the most intuitive evaluation however it requires a great deal of manual effort to perform. We need to pull out a sampling of classified data and look through it manually; marking whether or not the classification was correct.

Another evaluation that we can perform looks at how well each small collection does at being relevant to the topic at hand. If a large number of the documents in the collection are classified as non-relevant, than the collection as a whole has not done a good job capturing the event in question.

Finally, we need to perform an evaluation of how well our method of using FPM to determine the training documents works. This is much like the first evaluation in that it can be very manual. We will need to dump the generated training sets out to files and then manually decide whether a given set has all (or primarily) relevant or non-relevant documents, as appropriate.

### 3.3.6   Writing Results to HBase

# Chapter 4

# User Manual

Provide screenshots

## 4.1 Installation Requirements

In this documentation, we make the assumption that you have already installed at least Spark version 1.5.0 and have it configured properly. If you have not performed this step, you can find installation and setup instructions at `http://spark.apache.org/downloads.html`.

The next step will be to clone the git repository of our team. This can be found hosted on github at: `https://github.com/hosamshahin/Spring2016_IR_Project.git`

Finally, you will need to extract the `small_data.zip` file that can be found in the `Spring2016_IR_Project/data` directory

## 4.2 Preparing the Server

In order for some of the Machine Learning libraries in Spark to work, you need to ensure that you have `libgfortran3` installed.

On a Debian based system, you can run the following command:

```
$ sudo apt-get install libgfortran3
```

Next install Anaconda which is a completely free Python distribution. It includes more than 400 of the most popular Python packages for science, math, engineering, and data analysis. For a complete instructions on how to download and install Anaconda go to:

`https://www.continuum.io/downloads`

Anaconda is shipped with IPython notebook, A powerful interactive shell for Python. Now to set up Spark and IPython to run together and provide a nice web interface, there are a number of things that need to be done.

We will need to begin by creating an IPython profile for PySpark. This can be done using the command:

```
$ ipython profile create pyspark
```

If you are running other IPython profiles, you may want to change the default port. This can be done by modifying ~/.ipython/profile_pyspark/ipython_notebook_config.py.

You can then change the default port within that file.

```
c = get_config()
# Simply find this line and change the port value
c.NotebookApp.port = <your port number>
```

Next the following lines need to be placed in .bashrc.

```
export PYTHONPATH=$SPARK_HOME/python:$PYTHONPATH
export PYTHONPATH=$SPARK_HOME/python/lib/py4j-0.8.2.1-src.zip:$PYTHONPATH
```

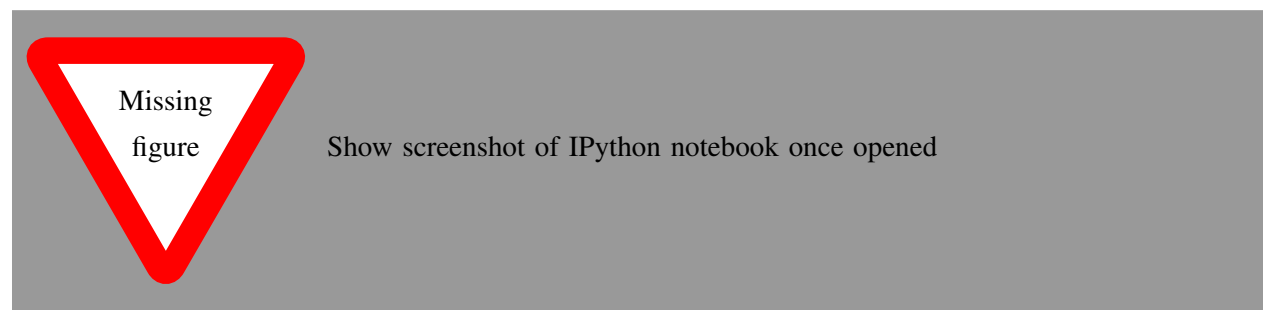After this the .bashrc file should be reloaded by relogging in, or sourcing the .bashrc file.

To run the IPython server you should run the following command from within the cloned git repository, replacing <IP> with the IP address of your server.

```
ipython notebook --profile=pyspark --no-browser --ip <IP>
```

You can then navigate to <IP>:<PORT> in your browser to get access to the IPython notebook. By default the port will be 8888.

## 4.3   Using the IPython Notebook

Once the IPython notebook has been opened in the browser, open the `Tweets_Classification-V2.ipynb` IPython notebook.



Missing
figure              Show screenshot of IPython notebook once opened

In order to use the script, you will need to modify the `base_dir` variable in the sixth cell down. Set this to point at the `data` directory.

### 4.3.1   Using the Configuration File

To go further in using the script, you will need to make modifications to the `collections_config.json` file under `data`. At this point the only section of the file that needs to be worried about is the first collection object. That is being used for the binary classifications we are currently using. The second object will be used for multiclass classification in the future.

To classify your data, you need to specify the table ID of your collection in the `collections_config.json` configuration file. You should model your entry off of one of the sample entries in the file, or modify a sample directly. It is also suggested that you update the name field to the name of your collection. At the moment, don't worry about the FP values.

### 4.3.2   Selecting a Frequent Pattern

After setting the configuration file, begin at the first cell and press Shift + Enter to execute each cell in order. Continue this until you reach "Manually choose frequent patterns and write them in the configuration file."

You will now need to open the frequent pattern output file located at:
`data/FPGrowth/«collectionId».txt`

You should now inspect the file for the tokens frequently found together. Look for high frequency patterns. Choose a pattern that seems highly relevant to your collection. We suggest a pattern of two to four words to strike a balance between over and under specification.

Take the pattern and copy it as the value of "FP" in the configuration file.

### 4.3.3   Building a Training Set

After the frequent pattern has been selected and input into the configuration file, continue to step through the script using Shift + Enter in each cell.

This will take you through the process of building a training set by identifying a positive sample set and a negative sample set and placing those into a DataFrame.

### 4.3.4   Training a Classifier

Further progression through the script will actually use the training data constructed in the last step to train a logistic regression classifier.

We also evaluate the logistic regression model on the training data before moving forward.

> How do you do the evaluation?

### 4.3.5   Running the Classifier

Finally apply the classification model to our test data and provide a prediction of relevance (1.0) or non–relevance (0.0) for each tweet in the test data.

# Chapter 5

# Developer Manual

> more concrete instructions such as on how the code is being developed, module overview, and how the data was processed, etc.

> collaborate more on the design and implementation steps

The developer manual will grow and evolve with the project. It documents the tools we are currently using, not tools we are considering for future steps.

## 5.1 Algorithms

### 5.1.1 Frequent Pattern Mining

Our methodology for constructing training data is based on frequent pattern mining algorithm. Spark.mllib version 1.6 provides a parallel implementation of FP-growth, a popular algorithm to mining frequent itemsets.

### 5.1.2 Logistic Regression

Logistic regression is a popular method to predict a binary response. It is a special case of Generalized Linear models that predicts the probability of the outcome. For more background and more details about the implementation, refer to the documentation of the logistic regression in spark.mllib.

## 5.2   Environment Setup

### 5.2.1   Dependencies

**Java**

You will need the latest version of Java, which at the time of this writing is Java Version 8 Update 73. Download and installation instructions for your environment can be found here:

https://java.com/en/download/

**Python**

You will need Python 2.7.9 (it may work with newer versions of the 2.7 line, but we have not tested it). Download and installation instructions for your environment can be found here:

https://www.python.org/downloads/release/python-279/

### 5.2.2   Apache Spark

This system is built on Apache Spark 1.6.0, which can be downloaded here:

http://spark.apache.org/downloads.html

# Chapter 6

# Plan

> Fix a lack of description on the management of the team work

> Provide breakdown of individuals contributions

Our work on this project will happen in two phases. In the first phase, we will attempt to create an initial prototype working with a small set of sample data. This phase will be accomplished mainly in our own external environment so that we can experiment more freely with the technologies to gain a better understanding of how they work and interact together. Currently we are in the finishing stages of this phase and preparing to move on to phase two.

Phase two of the project will involve taking the work done in phase one and expanding it out to the other small sets of data provided to us, and eventually to larger data sets. We will explore other potential techniques as we deem them necessary. This phase will also include a definition and implementation of an evaluation approach. Phase one is small enough that we feel a subjective evaluation is sufficient, but as we expand it will become essential to have a more concrete, quantitative measure of effectiveness.

Please see Table 6.1 for a weekly breakdown of work.

> Be sure to cite the tutorials and list them in references

**Table 6.1:** Weekly breakdown of work to be done.

| Weeks | End Date | Tasks |
|-------|----------|-------|
| Week 1 | 22 Jan. | Understanding the classification task |
| Week 2 | 29 Jan. | • Understanding the classification task<br><br>• Read about Hadoop streaming using Python |
| Week 3 | 5 Feb. | Start online tutorials about Hadoop and Apache Spark. |
| Week 4 | 12 Feb. | Continue online tutorials about Hadoop and Apache Spark. |
| Week 5 | 19 Feb. | **Phase 1 will include only tweets small data set:**<br><br>• Understanding the classification task<br><br>• Read about Hadoop streaming using python |
| Week 6 | 26 Feb. | • Prepare training data using Solr<br><br>• Build classifier using Apache Spark |
| Week 7 | 4 March | • Build classifier using Apache Spark<br><br>• Get output data format from Solr team |
| Week 8 | 11 March | Optimize classifier performance |

| Weeks | End Date | Tasks |
|---|---|---|
| Week 9 | 18 March | **Phase 2 will include tweets and webpages:**<br>**Once Spark on the cluster is updated to version 1.6 we will do the following:**<br><br>• Run our methodology to classify the tweets on the cluster.<br><br>• Apply the Frequent Pattern methodology on the cleaned webpages provided by Collection Management team.<br><br>• Develop HBase interface through which the classifier prediction output will be saved on HBase. |
| Week 10 | 25 March | Design an evaluation approach to test and evaluate our methodology. |
| Week 11 | 1 April | Train multiple classifiers and pick the best model per collection. |
| Week 12 | 8 April | More research on hyper-parameter optimization. Check the feasibility of Integrating hyper-parameters optimization library [1] output with Spark. |
| Week 13-15 | TBD | Search for known approaches to select the most representative data samples for each collection. Then check whether training a classifier using these data samples will enhance the performance or not. |

# Chapter 7

# Conclusion

The current state of the project is one of great promise. We have developed a functional first prototype of the classification system on our local environment. Our preliminary testing on our sample data indicates that the system is performing well from a subjective standpoint. From here, our main roadblocks are some updates that we need to be done to the cluster. The main one being that in order to use Frequent Pattern Mining, we need to be using Spark version 1.5.0 or newer. Once this is in place, we will be able to move the classification system over to the cluster and start branching out to the other small sets of data. This will allow us to optimize our classifier for general use, and ensure that any changes made affect the performance of the classification in general rather than its performance on our one sample set.

Give details that show the system is performing well

# Chapter 8

# Future Work

Chapter stub. We will expand this section as our project progresses.

# Bibliography

[1] Dan Bergstra James; Yamins and David D. Cox. "Hyperopt: A Python library for optimizing the hyper-parameters of machine learning algorithms". In: *Proceedings of the 12th Python in Science Conference*. 2013, pp. 13–20.

[2] Xuewen Cui, Rongrong Tao, and Ruide Zhang. *Classification Team Project for IDEAL in CS5604, Spring 2015*. http://vtechworks.lib.vt.edu/bitstream/handle/10919/52253/ReportClassify.pdf. 2015.

[3] Jiawei Han, Jian Pei, and Yiwen Yin. "Mining frequent patterns without candidate generation". In: *ACM Sigmod Record*. Vol. 29. 2. ACM. 2000, pp. 1–12.

[4] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *An Introduction to Information Retrieval*. Vol. 1. 1. Cambridge University Press Cambridge, 2008.