

# Topic Monitoring in the Pharmaceutical Industry

Master Team Project

presented by

Hailian Hou (123456789)  
Chia-Chien Hung (123456789)  
Lu Lifei (123456789)  
Olga Pogorelaya (123456789)  
Ngoc Nam Trung Nguyen (123456789)  
Md. Raziul Hasan Al Tariq (123456789)  
Alexander Weiß (123456789)

submitted to the

Data and Web Science Group  
Prof. Dr. Heiko Paulheim  
University of Mannheim

August 2017

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation and Goals . . . . .	2
1.2	Paper Structure . . . . .	3
<b>2</b>	<b>Fundamentals</b>	<b>4</b>
2.1	Social Media Data . . . . .	4
2.2	Sentiment Analysis . . . . .	5
2.3	Topic Detection . . . . .	6
2.4	Trend Detection . . . . .	6
<b>3</b>	<b>SentiTomo</b>	<b>7</b>
3.1	Overview . . . . .	7
3.1.1	Technical Explanation . . . . .	7
3.2	Data Collection . . . . .	8
3.2.1	Facebook API . . . . .	9
3.2.2	Twitter API . . . . .	11
3.2.3	Technology Stack . . . . .	17
<b>A</b>	<b>Program Code / Resources</b>	<b>21</b>
<b>B</b>	<b>Further Experimental Results</b>	<b>24</b>

# **List of Algorithms**

# List of Figures

# List of Tables

3.1	Anchor keywords for collecting data . . . . .	8
3.2	Common numbers of the Twitter dataset . . . . .	16
3.3	Twitter tweets per keyword . . . . .	17

# Listings

3.1	Twitter streaming endpoint . . . . .	13
3.2	Example Twitter API response . . . . .	13
3.3	Tweet object definition for the database . . . . .	15
A.1	TwitterCrawler.track . . . . .	21
A.2	FacebookCrawler.searchAndSaveFBPages . . . . .	22

# Acronyms

API     Application Programming Interface.

HTTP   Hypertext Transfer Protocol.

REST   Representational State Transfer.

SA     Sentiment Analysis.

TD     Topic Detection.

TRD   Trend Detection.

# Chapter 1

## Introduction

Today in 2017 social media is ubiquitous and is taken for granted in our day-to-day life. Nearly every single shop in the city to large companies own social media accounts to communicate with customers or try to engage new ones. The number of active users on the biggest social media platforms is tremendously high. On Facebook we have around 1.871 million active users as of January 2017. Whatsapp and Facebook messenger share the number of 1.000 million active users. For Twitter in comparison to Facebook, the numbers may seem small but with 317 million active users it belongs to the biggest social media platform of our time. (Chaffey 2017-04-2) More than ever people share their thoughts about current global events like conflicts between different countries, political events, terrorism, rising diseases as well as the opinions to companies and products.

In the course of the Master Team Project - Topic monitoring in the pharmaceutical industry of University of Mannheim for Master Data Science and Master Business informatics students, this work especially focuses on the monitoring of social media for the pharmaceutical context. It is processed in cooperation with AbbVie Inc. located in Ludwigshafen, Germany, a pharmaceutical company focused on both biopharmaceuticals and small molecule drugs. The goal of this project is to find efficient and effective ways to retrieve data from different social media platforms to analyze sentiments, public opinions, emotions towards different events, find suitable topics and detect rising trends in social media. To accomplish this task, this work utilizes different data mining algorithms to detect the requested features in the messages of users. The declared aim is to find a good combination of different of those algorithms to provide reliable results and findings so that especially AbbVie can use and incorporate them into their attempt to monitor social media.



This report describes how to gather data from different social media platforms and what restrictions will apply to them. It focuses on Facebook and Twitter, because this selection covers the most used and the most reactive social media platform today. Furthermore it deals with different machine learning algorithms, explains how they could be applied to different use cases and suggests situation when and how to use them.

## 1.1 Motivation and Goals

With the rapid growth of social media platforms like Twitter and Facebook in the recent years those has become interesting for companies to find out what customer think and talk about their company and products. A lot of people share their thoughts and experience about all manner of things. With Twitter we have one the highest reactive social media platforms of our time. Almost instantly people write messages (tweet) to react accordingly to new events in the world. Due to this fact companies try to effectively monitor Twitter and other platforms to get insights into the current customer situation regarding their product and company reputation. They do not need to organize big surveys among customers to find out these key values. It is sufficient to crawl social media platforms to get user insights and use this data in combination with data mining techniques to extract user sentiments and detect rising trends. This can all be done without the users concerned noticing or being able to influence such intrusion in a high manner. This work in particular focuses on the pharmaceutical industry to find out how customers react to different drugs and diseases. We want provide recommendations and workflow to utilize social media platforms to collect data and apply data mining algorithms to extract valuable results. The first goal is to detect the sentiment of users regarding drugs, diseases and different pharmacy companies. This is especially related to our partner AbbVie Inc. . With these result it shall be possible to get an overview about the reputation about AbbVie, its products and its competitors have on social media. Furthermore with the detection of sentiments in disease related texts we have the opportunity to get to know the The second goal of this project is the detection of topics and rising trends inside the social media data we gather from the different platforms. This is especially useful to find out what the users talk about and in what they are interested in in specific time ranges. These results can help to plan customer engagement strategies more accurate and related to current trends.

With this project we want to provide good insights into the structure of social media data and show suitable ways to extract sentiments, topic and trends out of it. Furthermore we want to provide you with some recommendations and hints you

may adopt to successfully mine social media data.

## 1.2 Paper Structure

In chapter 2 we discuss the fundamentals, which explains the basic understanding about social media platforms, sentiment analysis, topic detection and trend detection.

In chapter, we introduce and explain thoroughly step by step the final output of our project, namely SentiTomo. The product overview gives a general idea about SentiTomo and how it works and what it can achieve. Later we explain various modules it incorporates, like the data collection module, data preprocessing module, sentiment analysis module, topic detection module, trend detection module and finally the dashboard module.

In Chapter 3, we present our recommendations and conclude the report.

In the appendix, we have enclosed various tables and charts which we have used in this report, along with our bibliography.

## Chapter 2

# Fundamentals

At first we want to provide an explanation to some fundamental terms which will be used in this report. This should give an overview what the different terms mean and in which context they are used.

### 2.1 Social Media Data

The term *social media* refers to all web applications which are built on the technological foundations of the Web 2.0. In comparison to the beginning of the web, Web 2.0 is characterized by a huge improvement of communication channels between users and increased communications between those. Social media platforms heavily rely on this progressive development of user interaction. They offer a place where users can collaborate and interact with each other, share and create content with ease. Unlike traditional static websites, the content of social media platforms is nearly entirely made by the users, which lead to a very reactive and accurate picture of the community who is using the network. These platforms can be divided in different types. The following shows a non-exhaustive list of those. [11][12]

- Forums (Reddit)
- Microblogging services (Twitter)
- Social networks (Facebook, Google+)
- Wikis

- Live casting (YouTube, Twitch)
- Video sharing (YouTube, Twitch)

Social media is an integral part of everyones life and a world without it is nearly no longer imagineable. For companies it offers a broad variety of possibilities to engage and get to know their customers better. Social media analytics help to support business decisions, mostly based on the customer sentiment, which makes it a great datasource to analyze data. [11]

In this work we will mainly referring to the term social media in the context of two platforms, Twitter and Facebook, as those are building the base for our data gathering process.

## 2.2 Sentiment Analysis

Sentiment Analysis (SA) belongs to the broad field of data mining techniques. With the help of Natural Language Processing and text analytics, subjective information are extracted and quantified from texts and messages primarily from social media sources. SA methods are classification algorithms which are based on labelled data. This data is curated by human annotators who define pre-train classes to different types examples of data. These annotations later then are used to train the different algorithms which quantify the sentiment of texts based on these pre-annotations. The results of those techniques can quantify the sentiment of people's opinions and feelings towards a certain subject of the text. Mostly these are products, companies, other people or current events. Sometimes sentiment analysis is also known as opinion mining or emotion AI.

With the help of Sentiment Analysis companies can detect subjective key values for their organization and implement measures which take advantage of these findings. Some examples of these values are brand reputation along with the popularity of it, the company reputation overall and new product perception. This is essentially useful when new products are introduced and the acceptance of those needs to be measured as fast as possible. [14]

In this report sentiment analysis is used to determine the opinions of people regarding the pharmaceutical industry. It tries to apply suitable data mining algorithms to quantify the opinions regarding drugs, diseases and companies which are active in this sector.

### 2.3 Topic Detection

Topic Detection (TD) copes with the problem of identifying topics to unseen or unidentified events. It can be applied to a static collection of data ("retrospective detection") or even on real-time ("on-line detection") feed of data from different sources. To detect topics, miscellaneous data mining algorithms can be used. They are all based on using unlabeled data to train them. Unlabeled data means, that it is non curated by human annotators which detect topics beforehand. The algorithms just take raw data and try to cluster them into different groups (topics) based on common attributes among different data objects. [13]

Topic Detection within this work is primarily used on real-time news. It is first trained on pre-crawled data from social media platforms but later applied to the continuously stream of new data.

### 2.4 Trend Detection

Trend Detection (TRD) techniques try to find out what is suddenly increasing in popularity. This is especially applicable to highly reactive social media platforms. For some types organizations it is tremendously important to know what are trends among the people, especially for news organizations, governments and retailers. At most of the time a trend is coupled with a specific topic. If it possible to detect trending topics, organization can recognize what people think is noteworthy and can react to those events appropriately. [8]

## Chapter 3

# SentiTomo

In this section we want to discuss the different steps we took to build up our database and what algorithms we applied to them in order to fulfil the task of this work.

### 3.1 Overview

#### 3.1.1 Technical Explanation

As we decided to develop a small application which utilizes our findings to show a way how to integrate them into a real-world-like environment we need to find a suitable programming platform to use for this approach. For our data mining algorithms we used different programming languages. In total we had three which are R, Python and Java. But to let each of them cooperate with each other we had to find a base programming language which wraps around them. After a small evaluation we chose Javascript as our language to go. With the rise of Node.js, a framework for writing Javascript code on the server side, it is possible to invoke each of our machine learning languages inside Javascript. We used it to develop a Node.js based server application which can be run on any operating system which supports Node.js. Additionally to the plain server application it serves a front end view which is used to view the results of the machine learning algorithms. The whole application is backed by a MySQL database which saves the results sentiment analysis and topic detection algorithms. Additionally it is used to save raw tweets and their authors for a better evaluation of the results. We chose a SQL

based database because it was the easiest one for us to set up and in this use case it does not have any serious disadvantages compared to a NoSQL database.

Once the server is started it constantly uses the Twitter API to retrieve Tweets and their authors which match our keywords in table [keyword table]. These are then saved in the database. Parallel to this procedure two worker threads are running which detect sentiments and topics of tweets. The results of those are also getting saved into the database. They are separated into different tables inside the database provide a clear division between raw data and the result of data mining algorithms.

A detailed description of the implementation and structure of SentiTomo can be found in our separate sentitomo.md file. This file also includes detailed instruction of how to set up the application.

## 3.2 Data Collection

The first step we take is to collect data from related social media platforms, in our case Facebook and Twitter. We utilize the Application Programming Interface (API)s of those to retrieve posts, comments (Facebook) and tweets (Twitter) from users. This is done simultaneously for both Facebook and Twitter. We constructed a list of keywords related to pharmaceuticals industry, companies, diseases and products as "anchors" for the crawlers, to collect suitable and appropriate data. These keywords were used for both platforms.

Products	Companies	Diseases
adalimumab	abbvie	ankylosing spondylitis
humira	amgen	arthritis
enbrel	johnson&johnson	hepatitis
ibrutinib		psoriasis
		rheumatoid arthritis
		trilipix

Table 3.1: Anchor keywords for collecting data

### 3.2.1 Facebook API

#### Overview

Like every major web based platform, Facebook offers an Application Programming Interface (API) to programmatically interact with its data. This allows developers to build applications which can utilize the social connection and profile information to make them more involving. Furthermore they have access to public data on Facebook as well as the opportunity to publish posts and messages to the news feed. Officially the API is called Graph API, because it follows the concept of a Social Graph [3]. It consists of three key objects:

1. Nodes
  - Describe all things that are shown on the Facebook web page: User Pages, Posts, Comments and Images
2. Edges
  - Describe the relation between these things. For example the comments to an Image posted on a Page
3. Fields
  - Contain additional information to nodes. For example the name of pages or the relationship status of users

The Graph API is based on the simple Hypertext Transfer Protocol (HTTP) and can therefore work with every programming language which implements a HTTP library. Due to this fact it is very easy to implement the API into any application. It supports common request methods like, GET for retrieving information, POST for uploading data and DELETE for deleting data. With this report we exclusively use GET requests to get data from the Graph, as we do not want modify data on Facebook.

#### Restrictions

Facebook offers a highly efficient and very good documented API for getting data from the social network. So it would be actual a very good data source for social media monitoring, but unfortunately there are some restrictions which makes it hard to work. In the following the most serious ones are listed.

1. Data scraping terms of service



- Facebook owns a dedicated terms of service document for collecting data from Facebook through automated means. Before any automatic data collection takes place, a written permission has to be obtained by Facebook. Without this permission it is not allowed to scrape any data. Also other restrictions like not selling, or transferring the collected data apply at any given time. So it is absolutely necessary to read through these terms of services and contact Facebook before choosing it as a reliable datasource. [2]

## 2. Streaming API

- Facebook API offers no streaming endpoint, so no real time feed about new content on the platform can be obtained programmatically. It is up to the developer to create processes to constantly crawl for updated data on Facebook with static endpoints.

## 3. Rate limits

- The Graph API has some rate limits which will apply if too many requests are sent to it. During this work the rate limit only appeared once while a lot of pages for one keyword were crawled. So we assume that the rate limitation will not apply many times while using the API. [4]

## 4. No access to public posts

- On April 30th, 2015 made a breaking change to their search API endpoint. By then it was possible to search for public posts on Facebook through `/search?type=post&q=foobar`. After this day the endpoint was deprecated and exclusively available for approved companies. From then on, it was not possible to get any public user post. With this change in the API, it lost the most promising endpoint for companies for monitoring the network. Without the opportunity to crawl public posts from users the only available reliable source for getting posts is the page endpoint. With this it is possible to search the pages directory regarding specific keywords and retrieve posts from them. In fact, a company could try to apply for a license to get access to the post search endpoint, but it seems that this is only available to a limited set of media publishers. [5]

### **Workflow**

The restrictions explained in section 3.2.1 lead us to nearly drop Facebook from our workflow, but we tried to find a suitable way to get the most out of the remaining available endpoints. To overcome the deprecation of the post search endpoint, we tried to get as much information out of the pages search endpoint as we could. This serves the base endpoint for our work. It offers the opportunity to specify keywords which the pages need to match to be returned by the API. We used the keywords from Table 3.1 to crawl the pages directory. Afterwards it was necessary to restrict the retrieved pages to only those which were in a pharmaceutical context. Then the posts and comments from the pages can be crawled and saved for later work.

One major drawback of this approach is that the amount of pages and posts on these pages are very low. Only a few conversation about adverse drug reaction, opinions about drugs and companies take place on them. Furthermore the conversation inside those pages seems to be very limited and biased related to the keyword of the page. So this approach may not lead to a good overview about opinions on different keywords but for us it was the most suitable way to get data from Facebook.

### **The dataset**

### **Conclusion**

At first Facebook seemed to be a good and reliable data source for this project. But quickly it became clear that this is not the case. With the restriction of the public post search it is impossible to get all user posts about a keyword. Its only possible to get posts of pages and the corresponding comments to them. But it often occurs that posts on pages are very biased regarding the keyword of the page. Also a lot of pages in the pharmaceutical context seem to not get updated very often and some seem to be very abandoned. So the breaking point here definitely is the deprecation of the API endpoint for searching posts. This leads to the fact, that Facebook can not be used as a suitable data source for topic monitoring social media. Therefore we dropped the Facebook dataset out of the calculation and algorithms.

### **3.2.2 Twitter API**

## Overview

Like Facebook, Twitter also offers APIs to programmatically access data. They are divided into two different parts. One part are static Representational State Transfer (REST) conform ones which use simpleHTTP requests to create new tweets, read user profile, search for tweets and more.

The second part are dynamic APIs, better known as the Streaming API. They give applications access to the global real-time stream of Twitter data, without any overhead or by pulling data from the static API endpoints. In total Twitter offers three different streaming endpoints:

1. Public stream: For public Twitter data
2. User stream: For data associated with only one user
3. Site stream: Multi user version of user streams

This streaming endpoints differ significantly from the static endpoints. They are keeping a persistent HTTP connection open between client and server. The server then is streaming tweets as they occur and the client is responsible for processing and storing the result. This resolves into a realtime stream of twitter data which makes any application working with this model reactive and everytime up-to-date.[15]

## Restriction

Unlike Facebook, Twitter does not have that many restrictions we have to respect in our use case.

1. Date restriction in the search API (Static REST API)
  - When you want to use the static search endpoint, where you can search for tweets with specific keywords only the latest two weeks from the date of the request are available to search. So it is generally impossible to search for historical tweets through the API. There are some companies like GNIP [6], which are offering service to get access to historical data. To get access to this data you have to contact those companies to get a custom historical dataset, but this option was not suitable for this project as it would had involve payments and our goal is to avoid that.
2. Requests per minute (Static REST API)

- The REST endpoints, in contrast to the streaming endpoints are rate limited. So developers have to take care of the amount of requests that are sent to the API. A full overview about the different limits can be found at <https://dev.twitter.com/rest/public/rate-limits> [16]

### 3. Concurrent streaming APIs (Streaming API)

- Any application which is using any of the different streaming APIs can only open only one stream to Twitter. This have to be kept in the mind if a developer wants to use multiple streaming endpoint.

## Workflow

While getting to know the API better we decided to utilize the public streaming API to get the tweets for specific keywords mentioned in Table 3.1. This endpoint accepts a filter value, which defines what Tweets we will receive from the API. The keywords for that filter can either be separated by a comma, which will work as an **OR** concatenation or separated by an white space, which will work as an **AND** concatenation. An example request to the endpoint could look like the following:

```
1 https://stream.twitter.com/1.1/statuses/filter.json?&
   track=abbvie,humira,adalimumab...
```

Listing 3.1: Twitter streaming endpoint

This leads very convenient way to crawl data from Twitter only for specific topics. Another advantage of using the streaming API is that rate limits will not apply to them. Therefore is possible to receive a arbitrary number of tweets from it without getting rate limited, which is especially useful if peaks of tweets appear at certain events. The only restriction which applies is that it is not possible to open multiple data streams to the stream endpoint at a given time, but in our case this is not necessary as only the public one is used. Furthermore the response from this endpoint is very comprehensive so it is not necessary to get data from the static endpoints, which will benefit in not going to reach any rate limits. A truncated example response can be seen Listing 3.2.

```
1 { created_at: 'Mon Aug 21 13:17:14 +0000 2017',
2   id: 899621415120494600,
3   id_str: '899621415120494597',
```

```
4  text: 'RT @NRAS_UK: Rheumatoid Arthritis: What You
    Need to Know. Read article by @Treated_com here:
    https://t.co/Z036Mzepee https://t.co/bIh4kREAIi',
5  source: '<a href="http://twitter.com/#!/download/
    ipad" rel="nofollow">Twitter for iPad</a>',
6  truncated: false,
7  in_reply_to_status_id: null,
8  in_reply_to_status_id_str: null,
9  in_reply_to_user_id: null,
10 in_reply_to_user_id_str: null,
11 in_reply_to_screen_name: null,
12 user: [Object],
13 geo: null,
14 coordinates: null,
15 place: null,
16 contributors: null,
17 retweeted_status: [Object], // !null if this is a
    retweet, contains information about original
    tweet
18 is_quote_status: false,
19 retweet_count: 0,
20 favorite_count: 0,
21 entities:
22   { hashtags: [],
23     urls: [ [Object] ],
24     user_mentions: [ [Object], [Object] ],
25     symbols: [],
26     media: [ [Object] ] },
27 extended_entities: { media: [ [Object] ] },
28 favorited: false,
29 retweeted: false,
30 possibly_sensitive: false,
31 filter_level: 'low',
32 lang: 'en',
33 timestamp_ms: '1503321434179' }
```

Listing 3.2: Example Twitter API response

### The dataset

As of August 20th 2017 we stopped crawling data from the Twitter API. Our base dataset for Twitter then consists of 2 types of objects: Tweets and the authors of those. The most important one of these two is the tweet object with it's message field. The message is the fundamental part which is used to build models for sentiment analysis and topic detection. It contains the the tweet message along with URLs, Emojis and hashtags. The author object was only crawled for the sake of completeness and in case if it will be needed in some circumstances. The definition of the tweet object can be seen in Listing 3.3.

```
1  id: {
2      type: Sequelize.STRING,
3      primaryKey: true,
4  },
5  keywordType: {
6      type: Sequelize.STRING
7  },
8  keyword: {
9      type: Sequelize.STRING
10 },
11 created: { // when was the tweet created
12     type: Sequelize.DATE
13 },
14 createdWeek: {
15     type: Sequelize.INTEGER
16 },
17 toUser: {
18     type: Sequelize.STRING
19 },
20 language: {
21     type: Sequelize.STRING
22 },
23 source: {
24     type: Sequelize.STRING
25 },
26 message: {
27     type: Sequelize.STRING
28 },
29 hashtags: { //hashtags inside the message
```

```

30         type: Sequelize.STRING
31     },
32     latitude: { // latitude of the tweet author
33         type: Sequelize.STRING
34     },
35     longitude: { // longitude of the tweet author
36         type: Sequelize.STRING
37     },
38     retweetCount: {
39         type: Sequelize.INTEGER
40     },
41     favorited: {
42         type: Sequelize.BOOLEAN
43     },
44     favoriteCount: {
45         type: Sequelize.INTEGER
46     },
47     isRetweet: { // is the tweet a retweet
48         type: Sequelize.BOOLEAN
49     },
50     retweeted: { // was this tweet retweeted
51         type: Sequelize.INTEGER
52     }

```

Listing 3.3: Tweet object definition for the database

Type	Amount
Tweets	250732
Retweets	27631
Answer to tweets	597
User	48555

Table 3.2: Common numbers of the Twitter dataset

## Conclusion

After we had this massive disappointment with the Facebook dataset, Twitter is by far the better source for building up a dataset for data mining. It offers the oppor-

Keyword	Amount	Proportion in %
psoriasis	46070	25.5245
hepatitis c	35089	19.4406
rheumatoid arthritis	28940	16.0339
amgen	21404	11.8586
abbvie	14149	7.8391
bristol myers	13355	7.3992
johnson&johnson	5641	3.1253
arthritis	3513	1.9463
ankylosing spondylitis c	3473	1.9242
humira	3029	1.6782
ibrutinib	1851	1.0255
hepatitis	1662	0.9208
NULL	764	0.4233
enbrel	669	0.3707
adalimumab	667	0.3695
imbruvica	172	0.0953
trilipix	45	0.0249

Table 3.3: Twitter tweets per keyword

tunity to crawl tweets in real-time with no need to access the static endpoints. The API gives us access to every user tweet matching the keywords, which ensures that at most of the time only appropriate and suitable tweets are crawled. In addition to those advantages the streaming endpoint do not have an rate limitation, which makes it even easier to programmatically crawl tweets. This advantages make sure that Twitter is the primary data source for this project.

### 3.2.3 Technology Stack

The following description should emphasize how the tools we used to gather data from social media changed during the progress of this project.

In the first four months of the project we used corresponding R packages to access the data from Facebook and Twitter. For Facebook this was a package called *RFacebook* [1] and for Twitter the package was called *twitteR* [7]. The data was crawled manually by executing different R-Scripts which were saving their results into .csv files. These files then were uploaded to *Google Drive* to have common



database-like location for our data. But with the progress of this work the organization of the different files becomes unmanageable and we decided to save our data in a MySQL database. Simultaneously to that we switched our crawling stack and implemented it based on Node.js server which is directly connected to the database. With the database we had a structured place for our data which offered a better way to have access to our data than the *.csv* files. Furthermore the *.csv* files were encoding dependent which lead to a lot of issues during the initial crawl of data, the database solved this by not being affected of encodings. The implementation of the crawling processes was done with the help of different Node packages: *twitter* [9] for retrieving data from the streaming API of Twitter and *fbgraph* [10] for Facebook.

In contrast to the *twitteR* package *twitter* for Node.js is able to crawl the public stream of Twitter data. With that it was able to set up an automatic crawl mechanism which saved suitable tweets, regarding the keywords in Table 3.1, into our database. Because of that the manual data collection for Twitter was postponed and only served as a backup mechanism. A truncated (left out database saving) part of the code, which was responsible for the crawling and saving can be found in Listing A.1

As listed in section 3.2.1 Facebook does not offer a real-time stream to their data. Therefore the crawling of the Facebook data still had to be done manually by invoking methods from the *fbgraph* package on our server. The result of those manual crawls then was upserted into the database as well. A sample code of this process can be seen in Listing A.2

The update of our technology stack during the project led to a suitable and reliable way of crawling data from Twitter and Facebook. The centralised data consolidation inside the MySQL database additionally improved the workflow for our machine learning processes.

The final technology stack based on Node.js (JavaScript):

1. Node.js server running at least 8 hours a day
2. Twitter-Crawler for crawling data form the public stream API *twitter* package [9] (automatic)
3. Facebook-Crawler for crawling data from the GraphAPI *fbgraph* package [10] (manual)
4. MySQL database for saving tweets and authors (Twitter), pages,posts and comments (Facebook)

# Bibliography

- [1] Pablo Barbera et al. *Package "Rfacebook"*. May 25, 2017. URL: <https://cran.r-project.org/web/packages/Rfacebook/Rfacebook.pdf> (visited on 08/21/2017).
- [2] Facebook. *Facebook Site Scraping Terms of Service*. Ed. by Facebook. Apr. 10, 2010. URL: [https://facebook.com/apps/site\\_scraping\\_tos\\_terms.php?hc\\_location=ufi](https://facebook.com/apps/site_scraping_tos_terms.php?hc_location=ufi) (visited on 08/19/2017).
- [3] Facebook. *Graph API Overview*. Ed. by Facebook. Feb. 18, 2015. URL: <https://developers.facebook.com/docs/graph-api/overview/> (visited on 08/18/2017).
- [4] Facebook. *Graph API Rate Limit*. Ed. by Facebook. June 9, 2016. URL: [https://developers.facebook.com/docs/graph-api/advanced/rate-limiting?locale=de\\_DE](https://developers.facebook.com/docs/graph-api/advanced/rate-limiting?locale=de_DE) (visited on 08/19/2017).
- [5] Facebook. *Public Feed API*. Ed. by Facebook. Sept. 6, 2013. URL: [https://developers.facebook.com/docs/public\\_feed?locale=en\\_US](https://developers.facebook.com/docs/public_feed?locale=en_US) (visited on 08/19/2017).
- [6] Gnip. *Gnip - Unleash the Power of Social Media*. Ed. by Gnip. Oct. 20, 2015. URL: <https://gnip.com/> (visited on 08/19/2017).
- [7] Jeff Grenty. *Package "twitterR"*. July 29, 2015. URL: <https://cran.r-project.org/web/packages/twitterR/twitterR.pdf> (visited on 08/21/2017).
- [8] Jeff Kolb. *Trend detection in social data*. Ed. by Jeff Kolb. July 16, 2015. URL: [https://blog.twitter.com/official/en\\_us/a/2015/trend-detection-social-data.html](https://blog.twitter.com/official/en_us/a/2015/trend-detection-social-data.html) (visited on 08/22/2017).
- [9] Desmond Morris. *Node-Twitter*. Jan. 13, 2017. URL: <https://github.com/desmondmorris/node-twitter> (visited on 08/21/2017).

- [10] Cristiano Oliveira. *Node-fbgraph*. Jan. 19, 2017. URL: <https://github.com/criso/fbgraph> (visited on 08/21/2017).
- [11] Margaret Rouse. *Definition of Social Media*. Ed. by Margaret Rouse. Sept. 15, 2016. URL: <http://whatis.techtarget.com/definition/social-media> (visited on 08/21/2017).
- [12] Margaret Rouse. *Definition of Web 2.0*. Ed. by Margaret Rouse. Jan. 15, 2015. URL: <http://whatis.techtarget.com/definition/Web-20-or-Web-2> (visited on 08/21/2017).
- [13] Young-Woo Seo and Katia Sycara. *Text Clustering for Topic Detection*. Robotics Institute Carnegie Mellon University Pittsburgh, Pennsylvania 15213, Jan. 31, 2004. URL: <http://repository.cmu.edu/cgi/viewcontent.cgi?article=1705&context=robotics> (visited on 08/22/2017).
- [14] Techopedia. *What is Sentiment Analysis*. Ed. by Techopedia. Jan. 21, 2014. URL: <https://www.techopedia.com/definition/29695/sentiment-analysis> (visited on 08/22/2017).
- [15] Twitter. *Streaming API Overview*. Ed. by Twitter. May 15, 2012. URL: <https://dev.twitter.com/streaming/overview> (visited on 08/19/2017).
- [16] Twitter. *Twitter API Rate Limits Chart*. Ed. by Twitter. Nov. 22, 2016. URL: <https://dev.twitter.com/rest/public/rate-limits> (visited on 08/20/2017).

## Appendix A

### Program Code / Resources

```
1  /**
2  * @function track
3  * @param {String} filters The filters which are used
   to track tweets from the Twitter API
4  * @description Starts using the Twitter API with the
   specified filters. When a new tweet is crawled it
   upserts the author data
5  * and inserts the raw tweet data into the database.
6  * @see {@link module:Connectors~TwitterAuthor}
7  * @see {@link module:Connectors~Tweet}
8  * @memberof TwitterCrawler
9  * @return {void}
10 */
11 track(filters) {
12     logger.log('info', 'Start streaming Twitter tweets
   with filters: ' + filters);
13     this.client.stream(
14         'statuses/filter', {
15             track: filters
16         },
17         stream => {
18             stream.on('data', event => {
19                 if (event.lang == 'en') {
20                     //Save event (tweet) into the
   database
```

```

21         }
22     });
23     stream.on('error', function (error) {
24         logger.log('error', error);
25         //throw error;
26     });
27 }
28 );
29 }

```

Listing A.1: TwitterCrawler.track

```

1  /**
2  * @function searchAndSaveFBPages
3  * @param {String} keyword Keyword for searching
4  * @description Searches for Facebook pages which match
5  * the keyword. It also inserts the found pages and
6  * their posts and comments into the database.
7  * Only pages with matching categories (.env file
8  * config) are saved.
9  * @see File /server/.env
10 * @memberof FacebookCrawler
11 * @returns {void}
12 */
13 searchAndSaveFBPages(keyword) {
14     logger.info("Started Facebook search");
15
16     var pages = new Array();
17
18     var searchOptions = {
19         q: keyword,
20         type: "page",
21         fields: "name,id,feed{id,link,message,story,
22             likes,comments{id,message},created_time},
23         category"
24     }
25
26     graph.search(searchOptions, async (err, res) => {
27         //Iterate over results and save them into the db
28     })
29 }

```

```
23     } ) ;  
24 }
```

Listing A.2: FacebookCrawler.searchAndSaveFBPages

## **Appendix B**

# **Further Experimental Results**

In the following further experimental results are ...