



FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

Political Polarity in US Twitter

S e m i n a r B i g D a t a

at the
Friedrich Schiller University of Jena
Faculty of Mathematics and Computer Science
Graduate Degree Computer Science

submitted to
Prof. Dr. Bucker,
Dr. rer. nat. Bosse and
Herrn Schoder

submitted by
Kenny Gozali,
Chris Gerlach and
Walter Ehrenberger

Jena, January 27, 2023

Abstract

In der vorliegenden Arbeit behandeln wir eine Sentimentalitätsanalyse von US amerikanischen Politikern aus dem *House of Representatives*. Dazu haben wir Daten von Twitter der letzten 12 Jahren zu den genannten Repräsentanten *gescrap*t und mithilfe des Big Data Frameworks Spark verarbeitet. Ziel der Sentimentalitätsanalyse war es, Unterschiede der beiden Parteien (Republikaner und Demokraten) zu bestimmten politischen und auch allgemeinen Themen herauszufiltern. Jedoch haben sich in den gegebenen Daten weniger Diskrepanzen zwischen den beiden Parteien erkennen lassen, als zu Beginn erwartet, wie im Laufe dieser Arbeit deutlich wird.

Contents

1	Introduction	1
2	Getting data (good title missing)	3
2.1	Background	3
2.1.1	GetOldTweets3-Pakage	3
2.1.2	NLTK-Natural language Toolkit	4
2.1.3	TextBlob	4
2.1.4	Sparks	4
2.2	Scraping and Sanitization	5
2.2.1	Scraping	5
2.2.2	Sanitization	5
2.3	MapReduce (good title missing)	6
3	Analysing data (good title missing)	7
3.1	data1 (good title missing)	7
3.2	data2 (good title missing)	7
4	Schluss	8
4.1	Resume	8

List of Figures

1.1	Entwicklung der Polarität politisch engagierter Amerikaner.	1
2.1	Ine Beschreibung f	4
2.2	Ine Beschreibung f	5
2.3	Ine Beschreibung f	5

1 Introduction

Als mächtigste Weltmacht beeinflussen die Vereinigten Staaten nahezu jeden Teil des Globus. Mitunter deshalb und aufgrund der enormen Präsenz in den Medien sowie des Einflusses auf diese fallen Diskrepanzen in der Bevölkerung schneller auf als in anderen Ländern. Aufgrund dieser Stellung wirkt sich die dortige Sentimentalität somit auch auf das Leben in anderen Ländern aus. Der Kapitolsanschlag sowie die Black Lives Matter Proteste der letzten Jahre sind ein Zeichen für die zunehmende Polarität und Unzufriedenheit in der Bevölkerung, wie sich auch in folgender Grafik erkennen lässt [al.14].

Democrats and Republicans More Ideologically Divided than in the Past

Distribution of Democrats and Republicans on a 10-item scale of political values

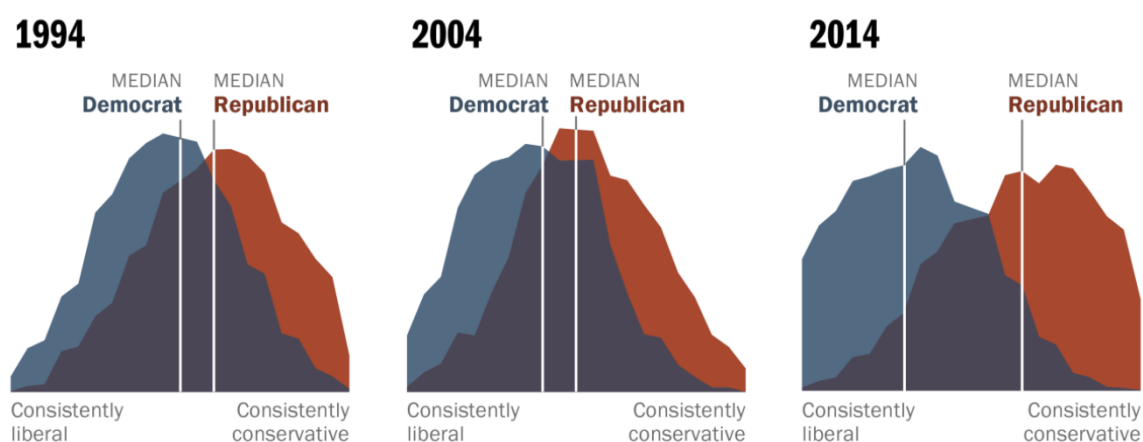


Figure 1.1: Entwicklung der Polarität politisch engagierter Amerikaner. Basierend auf 10 politischen Metriken werden Demokraten (blau) und Republikaner (rot) hier verglichen. Wie zu erkennen bewegen sich die zu Beginn teils noch überlappenden Ideologien in den letzten 10 Jahren auseinander [al.14].

Aufgabe dieser Arbeit ist es nicht, sich mit den komplexen und vielschichtigen Hintergründen für diese Entwicklung auseinanderzusetzen. Vielmehr wird hiermit ver-

sucht, eben diese Polarität in den oberen Reihen der amerikanischen Politik genauer zu analysieren.

Unsere Zielsetzung bestand darin, mit den Tweets der letzten 12 Jahre von 420 Politikern des Repräsentantenhauses eine Sentimentalitätsanalyse durchzuführen. Dabei handelt es sich um ein Mittel der natürlichen Sprachverarbeitung, bei dem die Ansicht beziehungsweise Gefühlslage eines Textes quantifiziert wird.

2 Getting data (good title missing)

2.1 Background

- Allgemeine Vorstellung der Packages ihr Vor und Nachteile

2.1.1 GetOldTweets3-Pakage

- Was ist der Vorteil der Scraping Bibliothek und Nachteile

GetOldTweets3 ist ein kostenlose Python 3 Packages mit welchen Twitterdaten ohne API-Schlüssel abgerufen werden können. Mit GetOldTweets3 können Sie Tweets mit einer Vielzahl von Suchparametern wie Start-/Enddatum, Benutzername(n), Textabfrage und Referenzortbereich durchsuchen. Außerdem können Sie angeben, welche Tweet-Attribute Sie einbeziehen möchten. Einige Attribute sind: Nutzername, Tweettext, Datum, Retweets und Hashtags.[?] Die offizielle API von Twitter hat eine lästige Zeiteinschränkung, weshalb man keine Tweets älter als eine Woche abrufen kann. Es gibt einige Tools die Zugang zu älteren Tweets anbieten, diese sind jedoch meistens kostenpflichtig. Das Forscher- team hat nach einem andere Tool gesucht die diese Aufgaben übernehmen, wodurch die Wahl auf das Package GetOldTweets3 gefallen ist.[?] Die Analyse des Codes von GetOldTweets3 und die Funktionsweise des Searchthrough Browsers von Twitter zeigt wie das Packages auch an alte Tweets kommt. Grundsätzlich, wenn sie auf Twitter seiten eingeben oder User suchen, startet ein Scroll-Loader. D.h. wenn sie dann weiter nach unten scrollen, bekommen sie immer mehr Tweets zu den Suchbegriffen. Diese ganzen Tweets bekommen sie durch Abfragen an einen JSON-Provider. GetOldTweets3 imitiert den Searchthrough Browser von Twitter um den Scroll-Loader zu starten und zieht sich dann anhand, der Abfragen an einen JSON-Provider die JSON-Datei und gibt diese decodiert zurück um dann alle Twitts anhand der oben gegebenen Parameter her-

auszufiltern. Dies kann man in Quelle [?], dem Github-Repositorium gut nachvollziehen. Somit ist es möglich sowohl aktuelle als auch sehr alte Tweets zu scrapen.[?]

```
1  #!/bin/bash
2  # cat user_list.csv
3
4  while IFS="," read -r rec_column1 rec_column2 rec_column3 rec_column4 rec_column5
5  do
6      echo "Writing to data/$rec_column3"
7      python GetOldTweets3/cli.py --username $rec_column3 > data/$rec_column3
8
9  done < <(tail -n +2 user_list.csv)
```

Figure 2.1: Code Beispiel für das Scrapen der Tweets

Ist eine Python Bibliothek mit der Twitter Daten durch den Scroll-Loader des Searchthrough Browsers von Twitter als JSON-Datei abgerufen werden können.

So kann durch eine paar Zeilen Code, wie man in der Abbildung sieht, eine bash-Datei erstellt werden, durch welche die Daten gesucht und abgespeichert werden. Das Scraping an sich kann durch die Größe der JSON-Datei einige Zeit in Anspruch nehmen. Wir haben gerade mal 2 Millionen Tweets insgesamt bei einer Laufzeit von ca. 35 Stunden.

2.1.2 NLTK-Natural language Toolkit

- Vllt Klären warum wir nicht NLTK verwendet haben, sondern Textblob - Vorteile von Textblob gegenüber NLTK - Nachteile von Nltk

2.1.3 TextBlob

- Textblob: Vllt. herausfinden wie die Berechnung stattgefunden hat - Vorteile von Textblob gegenüber NLTK - Warum haben wir das Tool genutzt.

2.1.4 Sparks

- Aus der Vorlesung Vorteile von Spark finden und einbauen

2.2 Scraping and Sanitization

2.2.1 Scraping

Democrats and Republicans More Ideologically Divided than in the Past
Distribution of Democrats and Republicans on a 10-item scale of political values

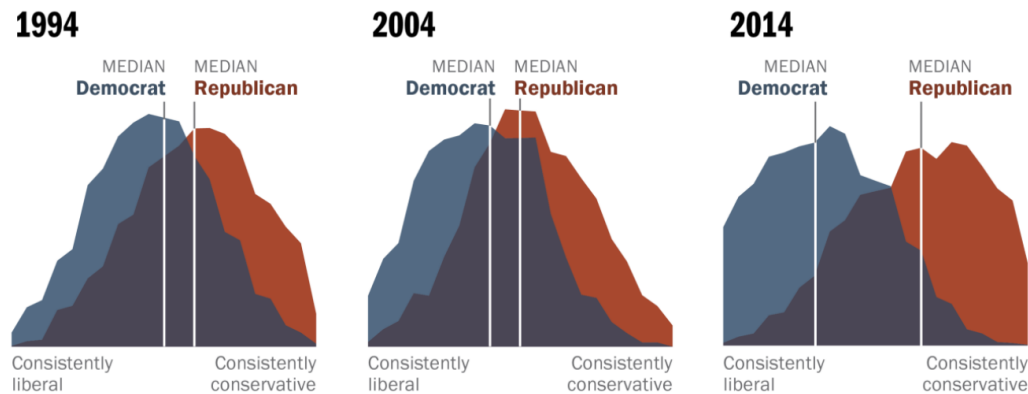


Figure 2.2: Ine Beschreibung f

2.2.2 Sanitization

Democrats and Republicans More Ideologically Divided than in the Past
Distribution of Democrats and Republicans on a 10-item scale of political values

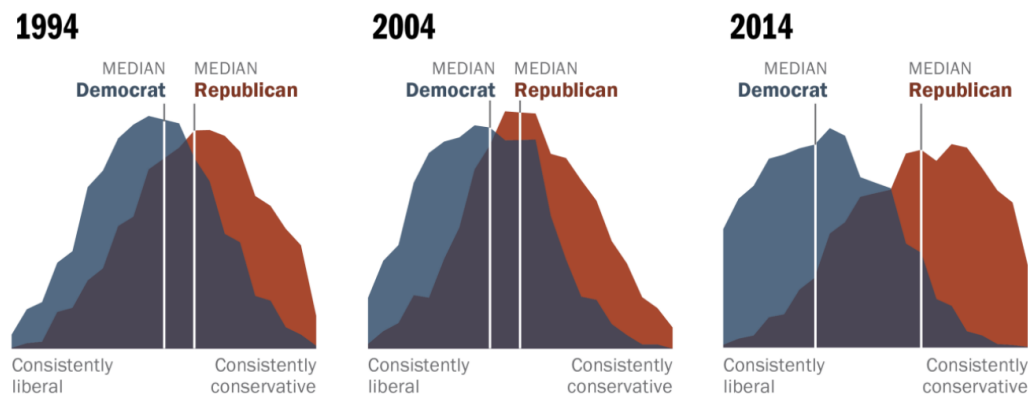


Figure 2.3: Ine Beschreibung f

2.3 MapReduce (good title missing)

Hallo Palmoooooooo und Walta

3 Analysing data (good title missing)

3.1 data1 (good title missing)

Hey Kennyyyyyy

3.2 data2 (good title missing)

Hey Kennyyyyyy

4 Schluss

4.1 Resume

Wie sie sehen sehen sie nichts.

Bibliography

[al.14] AL., Michael D.: Political Polarization in the American Public. (2014)

Internet sources

- [Moo06] MOOR, J.: *The Dartmouth College Artificial Intelligence Conference: The Next Fifty years*. 2006. – AI Magazine, Vol 27, No., 4, Pp. 87-9