

## TUGAS AKHIR

# DETEKSI DAN ANALISIS KARAKTERISTIK AKUN BUZZER MEDIA SOSIAL TWITTER PADA SISTEM TERSEBAR SPARK



Axel Joseph Yang

NPM: 6181901025

PROGRAM STUDI INFORMATIKA  
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS  
UNIVERSITAS KATOLIK PARAHYANGAN  
2024

**FINAL PROJECT**

**DETECTION AND ANALYSIS OF TWITTER SOCIAL MEDIA  
BUZZER ACCOUNT CHARACTERISTICS BASED ON SPARK  
DISTRIBUTED SYSTEM**



**Axel Joseph Yang**

**NPM: 6181901025**

**DEPARTMENT OF INFORMATICS  
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES  
PARAHYANGAN CATHOLIC UNIVERSITY  
2024**

## LEMBAR PENGESAHAN

### DETEKSI DAN ANALISIS KARAKTERISTIK AKUN BUZZER MEDIA SOSIAL TWITTER PADA SISTEM TERSEBAR SPARK

Axel Joseph Yang

NPM: 6181901025

Bandung, 11 Januari 2024

Menyetujui,

Pembimbing

Digitally signed  
by Veronica Sri  
Moertini

Prof. Dr. Veronica Sri Moertini

Ketua Tim Penguji

Digitally signed  
by Maria V.  
Claudia M.

Maria Veronica, M.T.

Anggota Tim Penguji

Digitally signed  
by Gede Karya

Gede Karya, M.T.

Mengetahui,

Ketua Program Studi

Digitally signed  
by Lionov

Lionov, Ph.D.

## PERNYATAAN

Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa tugas akhir dengan judul:

### **DETEKSI DAN ANALISIS KARAKTERISTIK AKUN BUZZER MEDIA SOSIAL TWITTER PADA SISTEM TERSEBAR SPARK**

adalah benar-benar karya saya sendiri, dan saya tidak melakukan penjiplakan atau pengutipan dengan cara-cara yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan.

Atas pernyataan ini, saya siap menanggung segala risiko dan sanksi yang dijatuhkan kepada saya, apabila di kemudian hari ditemukan adanya pelanggaran terhadap etika keilmuan dalam karya saya, atau jika ada tuntutan formal atau non-formal dari pihak lain berkaitan dengan keaslian karya saya ini.

Dinyatakan di Bandung,  
Tanggal 11 Januari 2024



Axel Joseph Yang  
NPM: 6181901025

## ABSTRAK

Terdapat bermacam-macam akun dalam media sosial, salah satunya adalah akun *buzzer*. Akun *buzzer* adalah akun pengguna media sosial yang secara sengaja menyebarkan informasi, opini, atau pesan tertentu dengan tujuan memengaruhi opini publik. Pada penelitian ini digunakan media sosial Twitter dengan domain Pemilihan umum Presiden Indonesia tahun 2024, karena Twitter merupakan salah satu media sosial dimana masyarakat bebas berpendapat dan kegiatan ini banyak dinanti oleh masyarakat Indonesia. Oleh karena itu ingin dideteksi akun *buzzer* untuk membedakan mana akun *buzzer* yang cenderung menggunakan kata-kata bernada positif (bersifat mendukung) dan yang bernada negatif (bersifat menyerang atau tidak sopan).

Data media sosial Twitter dikumpulkan dengan menggunakan teknologi *big data* bernama Apache Kafka. Apache Kafka adalah sebuah sistem pengumpul data dalam skala besar yang juga menjadi sarana dalam penyederhanaan *data stream*, nantinya hasil dari pengumpulan data akan disimpan pada Hadoop Distributed File System (HDFS). Lalu digunakan Apache Spark sebagai alat untuk mendeteksi dan menganalisis akun *buzzer* menggunakan data yang sudah tersimpan pada HDFS. Apache Spark sendiri merupakan sebuah *tools* yang dapat digunakan untuk memproses data secara paralel.

Sebelum dilakukan deteksi, datanya akan disiapkan terlebih dahulu sesuai kebutuhan. Deteksi pada penelitian ini terdiri dari 5 tahap dengan menggunakan fitur-fitur berdasarkan dari hasil studi literatur dan eksplorasi secara nyata, contoh salah satu fiturnya adalah *FollowersCount* atau jumlah pengikut. Diambil akun dengan jumlah pengikut di atas 500 pengikut, fitur ini didapatkan dari hasil studi literatur yang dilakukan. Selanjutnya hasil deteksi dianalisis untuk memberi contoh akun-akun yang tergolong sebagai *buzzer*, dari 59 akun yang didapatkan dari hasil deteksi, terdapat 51 akun yang merupakan akun *buzzer*. Oleh karena itu, ketika dihitung tingkat keberhasilan deteksinya, didapatkan angka sebesar 86.4%, dapat dikatakan bahwa deteksi yang dilakukan sudah berhasil.

Hasil dari deteksi ini kemudian akan dilanjutkan dengan klasifikasi untuk memprediksi sentimen *buzzer* untuk memahami sentimen yang dihasilkan oleh akun *buzzer* tersebut. Pada prediksi sentimen ini dilakukan 2 eksperimen dengan perbedaan pada jumlah datanya. Terdapat 2 metode yang digunakan untuk prediksi sentimen, yaitu metode *lexicon based* dan metode *machine learning* (*decision tree* dan *random forest*). Model prediksi akan dievaluasi menggunakan *confusion matrix* untuk mengetahui seberapa bagus hasil yang dihasilkan oleh model. Hasil prediksi dianalisis untuk memberi contoh akun hasil prediksi yang didapatkan beserta dengan contoh *tweet* yang dilakukannya.

Kesimpulannya, berdasarkan eksperimen yang sudah dilakukan, metode tersebut dapat mendeteksi akun *buzzer* menggunakan fitur-fitur yang didapatkan dari hasil studi literatur dan eksplorasi secara nyata. Hasil deteksinya dapat dimanfaatkan untuk melakukan analisis sentimen agar mendapatkan model yang dapat memprediksi dengan baik. Sentimen yang sudah didapatkan dapat digunakan untuk membedakan mana akun *buzzer* yang bernada positif dan negatif. Untuk menampilkan hasil dari penelitian ini, dibangun sebuah perangkat lunak berbasis *website* yang dapat menerima masukkan sebuah *file csv* untuk dideteksi dan diprediksi lalu menampilkan hasilnya.

**Kata-kata kunci:** *Buzzer*, Twitter, *Data Stream*, Apache Kafka, Apache Spark, Klasifikasi, Sentimen

## ABSTRACT

There are various types of accounts on social media, one of which is the buzzer account. A buzzer account is a social media user account that intentionally spreads information, opinions, or specific messages with the aim of influencing public opinion. In this research, Twitter social media platform is used with the domain of the 2024 Indonesian Presidential Election, because Twitter is one of the social media platforms where people can freely express their opinions and this activity is eagerly awaited by the Indonesian people. Therefore, buzzer accounts will be detected to distinguish between buzzer accounts that tend to use positively toned words (supportive) and those with negatively toned words (aggressive or impolite).

Twitter social media data is collected using a big data technology called Apache Kafka. Apache Kafka is a large-scale data collection system that also serves as a means to simplify data streams. The collected data will be stored on the Hadoop Distributed File System (HDFS). Then Apache Spark is used as a tool to detect and analyze buzzer accounts using the data stored on HDFS. Apache Spark itself is a unified computing engine that can be used to process data in parallel.

Before the detection process, the data will be prepared according to the needs. The detection in this research consists of 5 stages using features based on literature studies and real-world exploration, such as one of the features being the FollowersCount or the number of followers. Accounts with more than 500 followers are selected, this feature is obtained from the literature study conducted. Next, the detection results are analyzed to provide examples of accounts classified as buzzer. Out of 59 accounts obtained from the detection, 51 of them are identified as buzzer accounts. Therefore, when calculating the detection success rate, it is found to be 86.4%. It can be said that the detection performed has been successful.

The results of this detection will then be followed by classification to predict the sentiment of buzzer accounts to understand the sentiment generated by these accounts. In this sentiment prediction, two experiments will be conducted with a difference in the amount of data. Two methods will be employed for sentiment prediction, that is lexicon based and machine learning (decision tree and random forest algorithms). The prediction models will be evaluated using a confusion matrix to assess the performance of the models. The prediction results will be analyzed to provide examples of predicted accounts along with sample tweets they have made.

In conclusion, based on the conducted experiments, the methods can successfully detect buzzer accounts using features obtained from data collection and exploration. The detection results can be utilized for sentiment analysis to build a model capable of making accurate predictions. The obtained sentiment can be used to distinguish between buzzer accounts with positive and negative tone. To showcase the results of this research, a web-based software has been developed that can accept input of a csv file to be detected and predicted, and then display the results.

**Keywords:** Buzzer, Twitter, Data Stream, Apache Kafka, Apache Spark, Classification, Sentiment

*Dipersembahkan untuk kedua orang tua. . .*





## KATA PENGANTAR

Skripsi ini merupakan hasil dari perjalanan panjang yang penuh dengan tantangan dan pengalaman yang tak terlupakan. Maka dari itu dengan rasa syukur Penulis ingin menyampaikan terima kasih ke hadirat Tuhan yang Maha Esa, atas berkat-Nya dan karunia-Nya yang melimpah sehingga Penulis dapat menyelesaikan skripsi ini. Pada kesempatan kali ini, Penulis juga ingin mengungkapkan rasa terima kasih yang mendalam kepada semua pihak yang telah membantu dan mendukung selama menjalani kuliah dan selama proses penulisan skripsi ini, yaitu sebagai berikut:

1. Kepada kedua orang tua dan adik Penulis yang selalu mendukung, menyemangati dan memberikan doa.
2. Kepada Ibu Prof. Dr. Veronica Sri Moertini, Ir., M.T. selaku dosen pembimbing, yang telah banyak memberikan ilmu, saran, dan bantuan selama proses pengerjaan skripsi ini.
3. Kepada Ibu Maria Veronica Claudia, ST, MT dan Bapak Gede Karya, S.T., M.T., CISA, IPM, selaku dosen penguji yang telah memberikan masukan dan saran untuk skripsi ini.
4. Kepada teman-teman yang menemani Penulis selama masa kuliah, Jeremy, Febri, Rio, Gavin, Ferell, dan Indra.
5. Kepada teman-teman yang berjuang bersama selama skripsi, Filipus, Michael, Daryl, Brian, Kinan, dan Alma.
6. Kepada seluruh pihak yang mendukung Penulis yang tidak dapat disebutkan satu-satu.

Sebagai penutup, Penulis menyadari bahwa skripsi ini masih belum sempurna, oleh karena itu Penulis memohon maaf jika terdapat kesalahan penulisan atau metode yang digunakan pada skripsi ini. Penulis juga ingin berterima kasih kepada pembaca yang telah membaca skripsi ini, Penulis berharap skripsi ini dapat bermanfaat dan membantu untuk kemajuan ilmu pengetahuan, dan dapat dijadikan panduan untuk penelitian selanjutnya.

Bandung, Januari 2024

Penulis



# DAFTAR ISI

<b>KATA PENGANTAR</b>	<b>xv</b>
<b>DAFTAR ISI</b>	<b>xvii</b>
<b>DAFTAR GAMBAR</b>	<b>xix</b>
<b>DAFTAR TABEL</b>	<b>xxi</b>
<b>DAFTAR KODE PROGRAM</b>	<b>xxiv</b>
<b>1 PENDAHULUAN</b>	<b>1</b>
1.1 Latar Belakang . . . . .	1
1.2 Rumusan Masalah . . . . .	3
1.3 Tujuan . . . . .	3
1.4 Batasan Masalah . . . . .	3
1.5 Metodologi . . . . .	3
1.6 Sistematika Pembahasan . . . . .	4
<b>2 LANDASAN TEORI</b>	<b>5</b>
2.1 <i>Buzzer</i> . . . . .	5
2.2 Teknik Analisis Data Terstruktur . . . . .	7
2.2.1 <i>Distinct</i> . . . . .	7
2.2.2 <i>Filtering</i> . . . . .	8
2.2.3 <i>Sorting</i> . . . . .	8
2.2.4 <i>Aggregate</i> . . . . .	8
2.3 Statistika Dasar . . . . .	8
2.4 Dasar-Dasar <i>Text Mining</i> . . . . .	12
2.5 Hadoop . . . . .	13
2.6 Apache Spark . . . . .	14
2.6.1 Definisi dan Arsitektur Spark . . . . .	14
2.6.2 Konsep RDD . . . . .	15
2.7 Apache Kafka . . . . .	16
2.8 Analisis Sentimen . . . . .	18
2.9 Algoritma Klasifikasi . . . . .	19
2.9.1 <i>Decision Tree Classifier</i> . . . . .	19
2.9.2 <i>Random Forest Classifier</i> . . . . .	21
2.9.3 Metrik Evaluasi . . . . .	21
<b>3 STUDI EKSPLORASI</b>	<b>25</b>
3.1 Instalasi dan Konfigurasi Kafka di Lingkungan Hadoop/Spark . . . . .	25
3.2 Eksplorasi Kafka Cluster . . . . .	34
3.3 Eksplorasi <i>Buzzer</i> . . . . .	34
3.4 Komputasi Manual Metode pada Python . . . . .	50

<b>4</b>	<b>PENGUMPULAN, EKSPLORASI, DAN PENYIAPAN DATA</b>	<b>53</b>
4.1	Pengumpulan Data Twitter . . . . .	53
4.2	Penyiapan Data . . . . .	54
4.2.1	Penyiapan Data untuk Deteksi <i>Buzzer</i> . . . . .	54
4.2.2	Penyiapan Data untuk Deteksi Tahap Selanjutnya . . . . .	57
4.2.3	Penyiapan Data untuk Prediksi Sentimen <i>Buzzer</i> . . . . .	57
4.3	Eksplorasi Data . . . . .	60
<b>5</b>	<b>ANALISIS DATA DAN HASIL</b>	<b>63</b>
5.1	Deteksi dan Analisis Akun <i>Buzzer</i> . . . . .	63
5.1.1	Deteksi dan Analisis Akun <i>Buzzer</i> . . . . .	64
5.1.2	Analisis Hasil Deteksi Akun <i>Buzzer</i> . . . . .	67
5.2	Prediksi dan Analisis Sentimen <i>Buzzer</i> . . . . .	74
5.2.1	Pembuatan Model Prediksi Eksperimen 1 . . . . .	74
5.2.2	Pembuatan Model Prediksi Eksperimen 2 . . . . .	84
<b>6</b>	<b>PERANCANGAN PERANGKAT LUNAK DAN IMPLEMENTASINYA</b>	<b>95</b>
6.1	Fitur Perangkat Lunak . . . . .	95
6.2	Diagram <i>Use Case</i> . . . . .	95
6.3	Implementasi Perangkat Lunak . . . . .	97
<b>7</b>	<b>KESIMPULAN DAN SARAN</b>	<b>101</b>
7.1	Kesimpulan . . . . .	101
7.2	Saran . . . . .	101
	<b>DAFTAR REFERENSI</b>	<b>103</b>
	<b>A KODE PROGRAM</b>	<b>105</b>
	<b>B HASIL EKSPERIMEN</b>	<b>125</b>

## DAFTAR GAMBAR

1.1	Karakteristik <i>big data</i> [1]	2
2.1	Contoh kultwit pada Twitter [2]	6
2.2	<i>Scatterplot</i> [3]	10
2.3	<i>Symmetric histogram</i>	10
2.4	<i>Right Skewed histogram</i>	11
2.5	<i>Left Skewed histogram</i>	11
2.6	<i>Boxplot</i>	11
2.7	Komponen Apache Spark [4]	14
2.8	Arsitektur Kafka [5]	17
2.9	Ilustrasi <i>decision tree</i> <sup>1</sup>	19
2.10	Ilustrasi <i>Random Forest</i> <sup>2</sup>	21
2.11	Ilustrasi <i>Confusion Matrix</i> <sup>3</sup>	22
3.1	Download Kafka	26
3.2	Zookeeper properties	26
3.3	Kafka Manager	27
3.4	Kafka Manager create	27
3.5	Kafka Manager create topic	28
3.6	Topik pada Kafka Manager	28
3.7	Kafka Connect UI	29
3.8	Kafka Connect stream	29
3.9	HDFS home	30
3.10	HDFS user	30
3.11	HDFS vsm	31
3.12	HDFS data	31
3.13	HDFS twit-politik	32
3.14	HDFS result	32
3.15	HDFS year	33
3.16	HDFS month	33
3.17	HDFS day	34
3.18	Konfigurasi pada Kafka Connect	34
3.19	Contoh <i>tweet</i> yang dilakukan oleh Denny Siregar	35
3.20	Contoh <i>tweet</i> yang dilakukan oleh Denny Siregar	36
3.21	Contoh <i>tweet</i> yang dilakukan oleh Chusnul Chotimah	37
3.22	Contoh <i>tweet</i> yang dilakukan oleh Chusnul Chotimah	38
3.23	Contoh <i>tweet</i> yang dilakukan oleh Amel	39
3.24	Contoh <i>tweet</i> yang dilakukan oleh Amel	39
3.25	Contoh <i>tweet</i> yang dilakukan oleh Relawan Anies Baswedan	40
3.26	Contoh <i>tweet</i> yang dilakukan oleh Relawan Anies Baswedan	40
3.27	Contoh <i>tweet</i> yang dilakukan oleh Maudy Asmara	41
3.28	Contoh <i>tweet</i> yang dilakukan oleh Maudy Asmara	42
3.29	Contoh <i>tweet</i> yang dilakukan oleh Jarnas ABW Balikpapan	43

3.30	Contoh <i>tweet</i> yang dilakukan oleh Jarnas ABW Balikpapan . . . . .	44
3.31	Contoh <i>tweet</i> yang dilakukan oleh B. Prasetya . . . . .	45
3.32	Contoh <i>tweet</i> yang dilakukan oleh B. Prasetya . . . . .	46
3.33	Contoh <i>tweet</i> yang dilakukan oleh Dekade 08 . . . . .	47
3.34	Contoh <i>tweet</i> yang dilakukan oleh Dekade 08 . . . . .	48
3.35	Contoh <i>tweet</i> yang dilakukan oleh Yoyok . . . . .	49
3.36	Contoh <i>tweet</i> yang dilakukan oleh Yoyok . . . . .	49
4.1	Twitter API <i>key</i> . . . . .	53
4.2	Contoh data hasil <i>streaming</i> . . . . .	54
4.3	Memberi label secara manual pada <i>dataset</i> . . . . .	57
4.4	Contoh <i>CountVectorizer</i> . . . . .	58
4.5	Contoh nilai sentimen . . . . .	60
4.6	5 akun terbanyak berdasarkan <i>FollowersCount</i> . . . . .	61
4.7	5 akun terbanyak berdasarkan <i>StatusesCount</i> . . . . .	61
5.1	<i>Flowchart</i> deteksi akun <i>buzzer</i> . . . . .	63
5.2	<i>Boxplot</i> frekuensi kemunculan akun (atas) dan <i>Boxplot</i> secara lebih besar (bawah) . . . . .	66
5.3	Contoh perwakilan <i>tweet</i> . . . . .	67
5.4	Hasil label manual . . . . .	68
5.5	<i>Flowchart</i> prediksi dengan <i>lexicon based</i> . . . . .	75
5.6	Hasil akhir perhitungan nilai terletak pada kolom “Nilai NLP” . . . . .	76
5.7	Hasil sentimen terletak pada kolom “Sentimen NLP” . . . . .	77
5.8	<i>Flowchart</i> prediksi dengan <i>machine learning</i> . . . . .	78
5.9	Contoh 7 perwakilan <i>tweet</i> . . . . .	85
5.10	Hasil akhir perhitungan nilai terletak pada kolom “Nilai NLP” . . . . .	86
5.11	Hasil sentimen terletak pada kolom “Sentimen NLP” . . . . .	86
6.1	Diagram <i>use case</i> perangkat lunak . . . . .	96
6.2	Tampilan awal antarmuka perangkat lunak . . . . .	98
6.3	Tampilan antarmuka perangkat lunak saat menekan “ <i>Choose File</i> ” . . . . .	98
6.4	Tampilan antarmuka perangkat lunak setelah memilih <i>file</i> . . . . .	99
6.5	Tampilan antarmuka perangkat lunak hasil deteksi . . . . .	99
6.6	Tampilan antarmuka perangkat lunak hasil deteksi dan prediksi . . . . .	100

## DAFTAR TABEL

2.1	Data terstruktur . . . . .	7
2.2	Contoh daftar <i>stopwords</i> . . . . .	12
2.3	Contoh hasil <i>lowercasing</i> . . . . .	12
2.4	Contoh hasil <i>stemming</i> . . . . .	13
2.5	Contoh hasil <i>lemmatization</i> . . . . .	13
3.1	Contoh data yang digunakan untuk komputasi manual . . . . .	50
4.1	Kolom awal sebelum penyiapan data . . . . .	55
4.2	Kolom hasil sesudah penyiapan data . . . . .	55
4.3	Daftar akun berita dan akun radio yang dihapus . . . . .	56
4.4	<i>Hashtags</i> yang paling banyak digunakan . . . . .	62
5.1	Jumlah frekuensi kemunculan akun . . . . .	68
5.2	Hasil <i>buzzer</i> 1 . . . . .	69
5.3	Hasil <i>buzzer</i> 2 . . . . .	69
5.4	Hasil <i>buzzer</i> 3 . . . . .	69
5.5	Hasil <i>buzzer</i> 4 . . . . .	70
5.6	Hasil <i>buzzer</i> 5 . . . . .	70
5.7	Hasil <i>buzzer</i> 6 . . . . .	71
5.8	Hasil <i>buzzer</i> 7 . . . . .	71
5.9	Hasil <i>buzzer</i> 8 . . . . .	71
5.10	Hasil <i>buzzer</i> 9 . . . . .	72
5.11	Hasil <i>buzzer</i> 10 . . . . .	73
5.12	<i>Confusion matrix</i> dari metode <i>lexicon based</i> . . . . .	77
5.13	<i>Confusion matrix</i> dari <i>decision tree</i> . . . . .	79
5.14	<i>Confusion matrix</i> dari <i>decision tree</i> . . . . .	79
5.15	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	80
5.16	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	81
5.17	Hasil evaluasi model <i>decision tree</i> . . . . .	81
5.18	Hasil evaluasi model <i>random forest</i> . . . . .	81
5.19	Hasil evaluasi metode <i>lexicon based</i> . . . . .	82
5.20	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	82
5.21	Akun hasil prediksi . . . . .	82
5.22	<i>Tweet</i> akun bengkeldodo . . . . .	83
5.23	<i>Tweet</i> akun awaluddin7S . . . . .	83
5.24	<i>Tweet</i> akun elsadea15 . . . . .	84
5.25	<i>Confusion matrix</i> dari metode <i>lexicon based</i> . . . . .	87
5.26	<i>Confusion matrix</i> dari <i>decision tree</i> . . . . .	88
5.27	<i>Confusion matrix</i> dari <i>decision tree</i> . . . . .	89
5.28	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	89
5.29	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	90
5.30	Hasil evaluasi model <i>decision tree</i> . . . . .	90

5.31	Hasil evaluasi model <i>random forest</i> . . . . .	91
5.32	Hasil evaluasi metode <i>lexicon based</i> . . . . .	91
5.33	<i>Confusion matrix</i> dari <i>random forest</i> . . . . .	91
5.34	Contoh hasil prediksi . . . . .	91
5.35	Akun hasil prediksi . . . . .	92
5.36	Tweet akun aldhopriskap . . . . .	92
5.37	Tweet akun SintaPratiwi01 . . . . .	93
5.38	Tweet akun sedekahtissue . . . . .	93
6.1	Skenario kasus memasukkan <i>file csv</i> . . . . .	96
6.2	Skenario kasus mendeteksi akun <i>buzzer</i> . . . . .	97
6.3	Skenario kasus memprediksi sentimen akun <i>buzzer</i> . . . . .	97
B.1	Akun <i>buzzer</i> hasil deteksi . . . . .	125
B.2	Hasil prediksi sentimen eksperimen 1 . . . . .	125
B.3	Hasil prediksi sentimen eksperimen 2 . . . . .	126



## DAFTAR KODE PROGRAM

2.1	Contoh kode untuk membuat <i>DataFrames</i>	7
2.2	Contoh kode untuk melakukan <i>distinct</i>	7
2.3	Contoh kode untuk melakukan <i>filtering</i>	8
2.4	Contoh kode untuk melakukan <i>sorting</i>	8
2.5	Contoh kode untuk melakukan <i>aggregate</i>	8
2.6	Kode untuk membuat RDD	15
2.7	Contoh kode untuk melakukan <i>distinct</i>	15
2.8	Contoh kode untuk melakukan <i>filter</i>	15
2.9	Contoh kode untuk melakukan <i>map</i>	15
2.10	Contoh kode untuk melakukan <i>sort</i>	15
2.11	Contoh kode untuk melakukan <i>reduce</i>	16
2.12	Contoh kode untuk menghitung jumlah baris	16
2.13	Contoh kode untuk mengeluarkan nilai pertama	16
2.14	Contoh kode untuk mengeluarkan nilai maksimum	16
2.15	Contoh kode untuk mengeluarkan nilai minimum	16
3.1	Kode untuk mengaktifkan <i>zookeeper</i>	26
3.2	Kode untuk mengaktifkan <i>Kafka</i>	26
3.3	Kode untuk menghidupkan <i>web interface</i> Kafka Manager	26
3.4	Kode untuk menjalankan Kafka Connect	28
3.5	Kode untuk menghidupkan <i>web interface</i> Kafka Connect UI	28
3.6	Kode untuk menggunakan <i>consumer</i>	29
4.1	Kode untuk menggabungkan <i>file</i>	54
4.2	Kode untuk konversi <i>float</i>	56
4.3	Kode untuk memisahkan kalimat dalam teks	57
4.4	Kode untuk menghapus yang tidak penting	57
4.5	Kode untuk membuang <i>stopword</i>	57
4.6	Kode untuk melakukan <i>stemming</i>	57
5.1	Kode deteksi tahap 1	64
5.2	Kode deteksi tahap 2	64
5.3	Kode deteksi tahap 3	64
5.4	Kode deteksi tahap 4	66
5.5	Kode deteksi tahap 5	67
5.6	Kode untuk membaca korpus kata positif dan negatif	75
5.7	Kode untuk membandingkan korpus kata dengan teks	75
5.8	Kode untuk menghitung hasil akhir	76
5.9	Kode untuk memasukkan nilai ke <i>dataframe</i>	76
5.10	Kode untuk mengubah presentasi nilai menjadi sentimen	76
5.11	Kode untuk menggunakan <i>confusion matrix</i>	77
5.12	Kode untuk memilih fitur dan label	78

5.13	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	79
5.14	Kode untuk memilih fitur dan label . . . . .	79
5.15	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	79
5.16	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	80
5.17	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	80
5.18	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	80
5.19	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	81
5.20	Kode untuk membaca korpus kata positif dan negatif . . . . .	85
5.21	Kode untuk membandingkan korpus kata dengan teks . . . . .	85
5.22	Kode untuk menghitung hasil akhir . . . . .	85
5.23	Kode untuk memasukkan nilai ke <i>dataframe</i> . . . . .	86
5.24	Kode untuk mengubah presentasi nilai menjadi sentimen . . . . .	86
5.25	Kode untuk menggunakan <i>confusion matrix</i> . . . . .	87
5.26	Kode untuk memilih fitur dan label . . . . .	87
5.27	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	88
5.28	Kode untuk memilih fitur dan label . . . . .	88
5.29	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	88
5.30	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	89
5.31	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	89
5.32	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	89
5.33	Kode untuk menentukan jumlah <i>train</i> dan <i>test</i> data . . . . .	90
6.1	Kode untuk menampilkan hasil deteksi <i>buzzer</i> . . . . .	99
6.2	Kode untuk menampilkan hasil prediksi sentimen . . . . .	100
A.1	Kode untuk menggabungkan data per hari . . . . .	105
A.2	Kode untuk menggabungkan data per 7 hari . . . . .	105
A.3	Kode <i>import</i> yang digunakan . . . . .	105
A.4	Kode untuk eksplorasi dan penyiapan data . . . . .	106
A.5	Kode untuk deteksi akun <i>buzzer</i> . . . . .	107
A.6	Kode untuk menganalisis akun <i>buzzer</i> . . . . .	110
A.7	Kode untuk prediksi sentimen eksperimen 1 . . . . .	110
A.8	Kode untuk prediksi sentimen eksperimen 2 . . . . .	116
A.9	Kode perangkat lunak . . . . .	121
A.10	Kode perangkat lunak untuk <i>styling</i> . . . . .	122

# BAB 1

## PENDAHULUAN

Bab ini membahas mengenai latar belakang dilakukannya penelitian ini dengan mendeskripsikan gambaran besar permasalahan yang ada. Lalu dibahas rumusan masalah yang berisi mengenai inti permasalahan dari penelitian ini, dilanjutkan dengan tujuan yang berisi mengenai tujuan akhir yang ingin dicapai dalam penelitian ini sebagai solusi masalah. Kemudian terdapat batasan masalah yang membahas mengenai asumsi yang digunakan untuk membatasi ruang lingkup penelitian ini. Selanjutnya dibahas mengenai metodologi yang berisi mengenai tahapan pekerjaan dan eksperimen yang dilakukan, lalu diakhiri dengan pembahasan mengenai sistematika pembahasan yang berisi ringkasan apa saja yang dikerjakan pada setiap bab di buku ini.

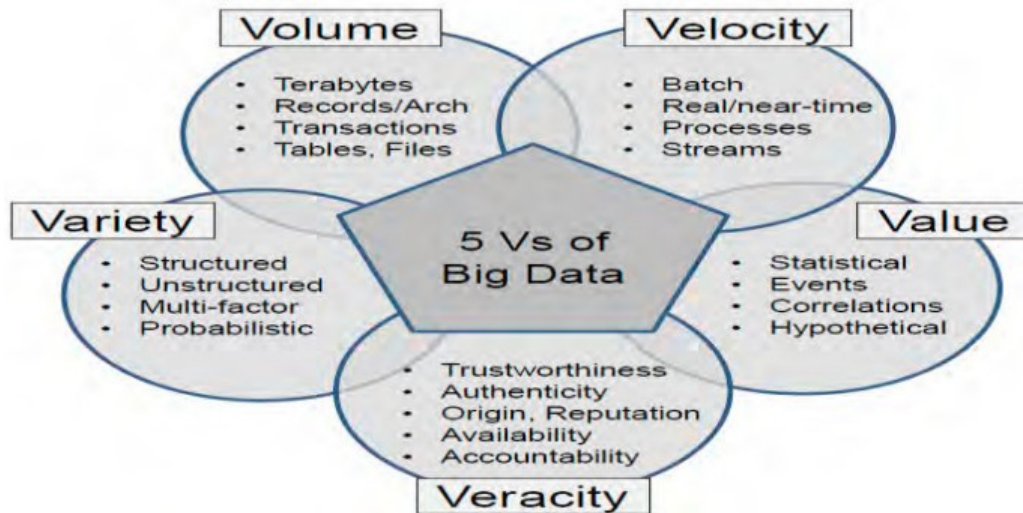
### 1.1 Latar Belakang

Media sosial adalah sebuah sarana digital yang memfasilitasi penggunaannya untuk saling bersosialisasi satu sama lain atau membagikan konten seperti foto dan video yang dilakukan secara *online* yang memungkinkan manusia untuk saling berinteraksi tanpa dibatasi oleh ruang dan waktu. Beberapa contoh dari media sosial adalah Instagram, Tiktok, dan Twitter. Instagram adalah sebuah aplikasi yang memungkinkan penggunaannya untuk mengunggah foto dan video ke dalam *post* dan *story* yang dapat diedit dengan berbagai filter, pengguna juga dapat menyukai foto dan video serta mengikuti pengguna lain untuk menambahkan konten pada *home* mereka. Tiktok adalah sebuah aplikasi yang memberi wadah bagi penggunaannya untuk dapat berekspresi melalui konten video. Twitter adalah sebuah media sosial yang memungkinkan penggunaannya untuk mengirimkan sebuah pesan singkat atau yang disebut dengan *tweet* kapanpun dan dimanapun mereka berada. *Tweet* sendiri bisa terdiri dari pesan teks maupun foto dan video, tetapi mayoritas isi dari *tweet* adalah pesan teks.

Pengguna Twitter terdiri dari berbagai macam tipe, seperti akun untuk berjualan, untuk berbagi tentang informasi suatu acara, maupun digunakan sebagai akun *buzzer*. Istilah *Buzzer* sendiri adalah orang atau sekumpulan orang yang pendapatnya memiliki pengaruh terhadap hal-hal tertentu. *Buzzer* berada di mana-mana, salah satunya pada media sosial Twitter. Akun *buzzer* pada Twitter adalah pengguna akun Twitter yang mayoritas mempunyai jumlah *followers* yang banyak dan memiliki pengaruh terhadap suatu hal atau acara. Akun *buzzer* biasanya mempunyai tugas untuk mempromosikan sesuatu, baik berupa produk, orang, dan lainnya secara berulang.

Banyaknya akun *buzzer* ini menyebabkan kebingungan kepada para pembacanya, seperti akun mana yang benar-benar memberi informasi yang informatif dan dapat dipercaya dan akun mana yang hanya melakukan *tweet* yang kurang informatif. Oleh karena itu dilakukan deteksi dan analisis karakteristik akun *buzzer* pada media sosial Twitter.

Data Twitter memenuhi karakteristik 5V dari *big data*. *Big Data* merupakan salah satu kemajuan di bidang teknologi yang berkembang dan dibutuhkan saat ini. *Big data* adalah sekumpulan data yang berukuran besar yang jumlahnya terus bertambah, yang terdiri dari berbagai macam jenis data dan terbentuk dengan terus menerus dengan kecepatan tertentu [6]. *Big data* memiliki karakteristik khusus 5V, yaitu *volume*, *velocity*, *variety*, *veracity*, dan *value* seperti pada Gambar 1.1.



Gambar 1.1: Karakteristik *big data* [1]

*Volume* mengacu pada sekumpulan data dalam jumlah dan ukuran yang sangat besar. *Velocity* adalah kecepatan penambahan data dalam kurun waktu tertentu. *Variety* mendefinisikan berbagai jenis data yang ada, mulai dari data terstruktur seperti tabel basis data, data semi-terstruktur seperti *Extensible Markup Language* (XML) dan *JavaScript Object Notation* (JSON), dan data tidak terstruktur seperti citra, suara, dan teks. *Veracity* mengarah kepada seberapa akurat dan dapat dipercaya suatu data. *Value* merupakan seberapa bernilainya atau bermaknanya suatu data.

Data Twitter merupakan data terstruktur yang memenuhi karakteristik 5V dari *big data* dari segi *volume* dan *velocity*, karena seperti yang disebutkan pada penelitian Gupta dan Hewett [7], diperkirakan setiap harinya ada 500 juta *tweets* yang ditulis di Twitter, yang dapat menyebabkan *data stream* yang tidak terbatas dan terus berkembang. Oleh karena itu, Twitter menghasilkan *volume* data yang sangat besar dengan *velocity* yang sangat tinggi.

*Big data* memiliki beberapa masalah seperti metode penyimpanan dan cara pemrosesannya, maka dari itu harus diproses menggunakan *tools* khusus untuk menggunakannya. Hal ini disebabkan karena data yang digunakan jumlahnya sangat banyak, sehingga kurang memadai jika diolah menggunakan *tools* biasa atau tradisional. *Big data* memerlukan *tools* atau teknologi yang memadai dengan kapasitas penyimpanan yang besar atau dengan nama lainnya adalah teknologi *big data*.

Teknologi *big data* adalah teknologi khusus untuk menangani masalah *big data*. Walaupun *big data* memiliki karakteristik 5V, tetapi masalah yang harus ditangani oleh teknologi *big data* hanya *volume*, *velocity*, dan *variety*. Pertama untuk masalah *volume*, digunakan teknik penyimpanan dan pemrosesan data terdistribusi. Kedua untuk masalah *velocity*, digunakan pemrosesan *stream* dan terdistribusi. Ketiga untuk masalah *variety*, digunakan teknik integrasi data dan penyimpanan data tidak terstruktur [8]. Beberapa teknologi *big data* yang memadai adalah Hadoop, Spark, Hive, Kafka dan masih banyak lagi, dalam tugas akhir ini digunakan Kafka untuk pengumpulan data media sosial Twitter dan digunakan Spark untuk menganalisis akun *buzzer*. Oleh karena itu akan dibahas lebih lanjut tentang Kafka dan Spark.

Kafka diciptakan untuk menangani aliran data secara *real-time* yang aman dan efisien. Seiring dengan berjalannya waktu kebutuhan untuk pemrosesan data secara *real-time* semakin meningkat, oleh karena itu Kafka menjadi semakin banyak digunakan dan populer. Kafka sendiri adalah sistem pengumpul data dan sistem penyederhanaan *data stream*, yang dapat menangani jumlah *volume* data yang sangat banyak secara *real-time*. Komponen dari Kafka terdiri dari *topics*, *partitions*, *cluster*, *broker*, *producer*, dan *consumer*.

Spark adalah sebuah mesin komputasi dan sekumpulan *library* yang dapat digunakan untuk memproses data secara paralel [9]. Spark memiliki keunggulan yang berupa *in-memory processing*, dimana komputasi dilakukan di dalam *memory*. Spark juga dapat digunakan pada berbagai bahasa

pemrograman seperti Python, Java, Scala, dan R, lalu Spark juga mendukung menggunakan kueri *Structured Query Language* (SQL). Spark memungkinkan penggunaannya untuk memproses data dalam jumlah yang besar, dengan berbagai macam *library* yang relatif mudah digunakan, seperti Spark SQL, Spark Streaming, Spark MLlib, dan GraphX.

## 1.2 Rumusan Masalah

Rumusan masalah dari latar belakang yang telah dipaparkan adalah sebagai berikut:

1. Bagaimana mengumpulkan data Twitter menggunakan Kafka?
2. Bagaimana eksplorasi dan penyiapan data Twitter yang sudah dikumpulkan?
3. Bagaimana mendeteksi akun *buzzer* pada media sosial Twitter pada sistem tersebar Spark dengan pendekatan analisis *dataset batch*?
4. Bagaimana analisis akun *buzzer* dari data Twitter yang sudah dikumpulkan?
5. Bagaimana menampilkan informasi *buzzer* dan hasil analisisnya dari perangkat lunak dengan menarik dan informatif?

## 1.3 Tujuan

Tujuan penelitian dari rumusan masalah yang telah dipaparkan adalah sebagai berikut:

1. Mengumpulkan data Twitter menggunakan Kafka.
2. Melakukan eksplorasi dan penyiapan data Twitter yang sudah dikumpulkan.
3. Mendeteksi akun *buzzer* pada media sosial Twitter pada sistem tersebar Spark dengan pendekatan analisis *batch dataset*.
4. Melakukan analisis terhadap akun *buzzer* dari data Twitter yang sudah dikumpulkan.
5. Menampilkan informasi *buzzer* dan hasil analisisnya dari perangkat lunak dengan menarik dan informatif.

## 1.4 Batasan Masalah

Batasan masalah untuk penelitian ini adalah:

1. *Tweet* yang dianalisis hanya *tweet* yang berupa teks saja.
2. Teks *tweet* yang dianalisis hanya berasal dari domain Pemilihan umum Presiden Indonesia tahun 2024.

## 1.5 Metodologi

Metodologi yang digunakan dalam penelitian ini adalah sebagai berikut:

1. Melakukan studi literatur mengenai akun *buzzer*.
2. Melakukan studi literatur mengenai teknik analisis data terstruktur, statistika dasar, dan dasar-dasar *text mining*.
3. Melakukan studi literatur mengenai Apache Spark dan Apache Kafka.
4. Melakukan studi literatur mengenai analisis sentimen dan algoritma klasifikasi.
5. Melakukan instalasi atau konfigurasi Kafka pada kluster Hadoop/Spark dan pengumpulan data media sosial Twitter menggunakan Kafka.
6. Melakukan eksplorasi terhadap fungsi-fungsi yang dapat digunakan untuk mengumpulkan data deteksi akun *buzzer*.
7. Melakukan eksplorasi *buzzer* secara nyata pada Twitter.
8. Melakukan eksplorasi dan penyiapan data media sosial Twitter.
9. Melakukan deteksi akun *buzzer* dan analisis terhadap *buzzer* yang didapatkan.
10. Melakukan prediksi sentimen akun *buzzer* dan analisis terhadap hasil prediksinya.

11. Mengevaluasi hasil analisis yang didapatkan.
12. Membangun perangkat lunak yang dapat mendeteksi *buzzer*, memprediksi sentimen *buzzer*, dan menampilkan hasilnya.

## 1.6 Sistematika Pembahasan

Sistematika penulisan tugas akhir ini adalah sebagai berikut:

1. Bab 1 Pendahuluan  
Pada Bab 1 membahas tentang pengantar tugas akhir yang berisi latar belakang, rumusan masalah, tujuan, batasan masalah, metodologi, dan sistematika pembahasan.
2. Bab 2 Landasan Teori  
Pada Bab 2 membahas tentang studi literatur yang mencakup penjelasan dari konsep-konsep yang digunakan, seperti teori *buzzer*, teori teknik analisis data terstruktur, teori statistika dasar, teori dasar *text mining*, teori Apache Spark, teori Apache Kafka, teori analisis sentimen dan teori mengenai algoritma klasifikasi.
3. Bab 3 Studi Eksplorasi  
Pada Bab 3 membahas tentang langkah-langkah dan eksperimen untuk melakukan konfigurasi pada teknologi yang digunakan, serta hasil eksplorasi *buzzer* secara nyata dan contoh komputasi manual metode pada Python.
4. Bab 4 Pengumpulan, Eksplorasi, dan Penyiapan Data  
Pada Bab 4 membahas tentang bagaimana pengumpulan data dilakukan, lalu eksplorasi terhadap data yang didapatkan, dan diakhiri dengan penyiapan data agar dapat digunakan dengan baik.
5. Bab 5 Analisis Data  
Pada Bab 5 membahas tentang tahap-tahap untuk mendeteksi akun *buzzer* beserta dengan hasil analisis dan evaluasinya, lalu dilanjutkan dengan prediksi sentimen *buzzer* dan hasil analisis serta evaluasinya.
6. Bab 6 Perancangan Perangkat Lunak dan Implementasinya  
Pada Bab 6 membahas tentang fitur yang tersedia pada perangkat lunak dan juga memperlihatkan cara kerja dari perangkat lunak yang dibangun beserta hasilnya.
7. Bab 7 Kesimpulan dan Saran  
Pada Bab 7 membahas tentang kesimpulan dari keseluruhan penelitian yang sudah dilakukan dan saran yang membangun yang bisa diterapkan dipengembangan selanjutnya.