



**IMPLEMENTASI EKSTRAKSI FITUR *BAG OF WORDS*  
MENGUNAKAN SELEKSI FITUR *CHI-SQUARE* DENGAN  
KLASIFIKASI *CATBOOST* PADA ANALISIS SENTIMEN MEDIA  
SOSIAL : COVID – 19**

**Skripsi**

**Untuk Memenuhi Persyaratan  
Dalam Menyelesaikan Strata-1 Ilmu Komputer**

**Oleh**

**Muhammad Zamzam**

**NIM 1711016210022**

**PROGRAM STUDI S-1 ILMU KOMPUTER  
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM  
UNIVERSITAS LAMBUNG MANGKURAT  
BANJARBARU**

**JUNI 2023**



**IMPLEMENTASI EKSTRAKSI FITUR *BAG OF WORDS* MENGGUNAKAN  
SELEKSI FITUR *CHI-SQUARE* DENGAN KLASIFIKASI *CATBOOST*  
PADA ANALISIS SENTIMEN MEDIA SOSIAL : COVID - 19**

**Skripsi**

**Untuk Memenuhi Persyaratan  
Dalam Menyelesaikan Strata-1 Ilmu Komputer**

**Oleh**

**Muhammad Zamzam**

**NIM 1711016210022**

**PROGRAM STUDI S-1 ILMU KOMPUTER  
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM  
UNIVERSITAS LAMBUNG MANGKURAT  
BANJARBARU**

**JUNI 2024**

**SKRIPSI**

**IMPLEMENTASI EKSTRAKSI FITUR *BAG OF WORDS* MENGGUNAKAN  
SELEKSI FITUR *CHI-SQUARE* DENGAN KLASIFIKASI *CATBOOST*  
PADA ANALISIS SENTIMEN MEDIA SOSIAL : COVID - 19**

Oleh :

**Muhammad Zamzam**

**1711016210022**

Telah dipertahankan di depan Dosen Penguji pada tanggal 28 Mei 2024

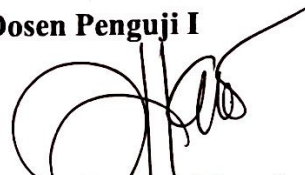
Susunan Dosen Penguji :

**Pembimbing I**



Muliadi, S.Kom., M.Cs.  
NIP. 197804222010121002

**Dosen Penguji I**



Dodon Turianto Nugrahadi, S.Kom., M.Eng.  
NIP. 198001122009121002

**Pembimbing II**



Friska Abadi, S.Kom M.Kom  
NIP. 198809132023211010

**Dosen Penguji II**



Triando Hamonangan Saragih, S.Kom M.Kom  
NIP. 199308242019031012



Banjarbaru, 28 Mei 2024

Ketua Program Studi Ilmu Komputer

Wan Budiman, S.T., M.Kom  
NIP. 197703252008121001

## PERNYATAAN

Dengan ini saya menyatakan bahwa dalam skripsi ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar kesarjanaan di suatu Perguruan Tinggi, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam Daftar Pustaka.

Banjarbaru, 13 Juli 2023



Muhammad Zamzam  
NIM. 1711016110007

## ABSTRAK

### **IMPLEMENTASI EKSTRAKSI FITUR *BAG OF WORDS* MENGGUNAKAN SELEKSI FITUR *CHI-SQUARE* DENGAN KLASIFIKASI *CATBOOST* PADA ANALISIS SENTIMEN MEDIA SOSIAL : *COVID - 19***

(Oleh: Muhammad Zamzam; Pembimbing: Muliadi, S.Kom., M.Cs. dan Friska Abadi, S.Kom M.Kom; 2023; 60 halaman )

Sejak awal tahun 2020, virus *SARS-CoV-2* yang menyebabkan *COVID-19* telah menjadi perhatian global dengan dampak signifikan terhadap masyarakat dan berbagai sektor kehidupan. Media sosial, khususnya Twitter, menjadi platform penting untuk mengumpulkan opini publik terkait *COVID-19* dan vaksinasi. Kebutuhan untuk memahami sikap dan pandangan masyarakat terhadap vaksinasi *COVID-19* di Twitter menuntut adanya analisis sentimen yang efektif. Berbagai metode analisis sentimen telah dilakukan sebelumnya, termasuk menggunakan *Bag of Words (BoW)* dan berbagai metode klasifikasi seperti *Naïve Bayes*, *Decision Tree*, *Random Forest*, dan *Catboost*. Penelitian ini berfokus pada implementasi ekstraksi fitur *Bag of Words* dengan seleksi fitur *Chi-Square* dan metode klasifikasi *Catboost* untuk menghasilkan model yang dapat mengklasifikasikan komentar *Twitter* terkait *COVID-19* secara akurat. Hasil penelitian menunjukkan bahwa penggunaan seleksi fitur *Chi-Square* meningkatkan akurasi klasifikasi. Rata-rata akurasi tertinggi dengan seleksi fitur *Chi-Square* adalah 75,2%, dibandingkan dengan 71,2% tanpa seleksi fitur, menunjukkan peningkatan akurasi sebesar 4%. Kesimpulannya, menggabungkan seleksi fitur *Chi-Square* dengan ekstraksi fitur *Bag of Words* dapat meningkatkan akurasi klasifikasi *Catboost* pada analisis sentimen *COVID-19*.

Kata kunci: *Sentiment analysis*, *Twitter*, *COVID-19*, *Bag of Words (BoW)*, *Chi-Square*, *Catboost*

## **ABSTRACT**

(By: Muhammad Zamzam; Supervisor: Muliadi, S.Kom., M.Cs. and Friska Abadi, S.Kom M.Kom; 2023; 60 pages)

*Since early 2020, the SARS-CoV-2 virus, which causes COVID-19, has garnered global attention with significant impacts on society and various sectors of life. Social media, particularly Twitter, has become an important platform for gathering public opinion regarding COVID-19 and vaccination. The need to understand public attitudes and views on COVID-19 vaccination on Twitter necessitates effective sentiment analysis. Various sentiment analysis methods have been conducted previously, including using Bag of Words (BoW) and classification methods such as Naïve Bayes, Decision Tree, Random Forest, and Catboost. This study focuses on implementing Bag of Words feature extraction with Chi-Square feature selection and the Catboost classification method to produce a model that can accurately classify Twitter comments related to COVID-19. The results show that using Chi-Square feature selection increases classification accuracy. The highest average accuracy with Chi-Square feature selection is 75.2%, compared to 71.2% without feature selection, showing an accuracy increase of 4%. In conclusion, combining Chi-Square feature selection with Bag of Words feature extraction can improve the accuracy of Catboost classification in COVID-19 sentiment analysis.*

**Keywords:** *Sentiment analysis, Twitter, COVID-19, Bag of Words (BoW), Chi-Square, Catboost*

## PRAKATA

Assalamualaikum Warahmatullahi Wabarakatuh.

Puji syukur penulis panjatkan ke hadirat Allah SWT karena atas berkat rahmat dan karunia-Nya penulis dapat menyelesaikan skripsi yang berjudul “*Implementasi Ekstraksi Fitur Bag of Words Menggunakan Seleksi Fitur Chi-Square Dengan Klasifikasi Catboost Pada Analisis Sentimen Media Sosial : Covid - 19*” untuk memenuhi syarat dalam menyelesaikan pendidikan program S1 Ilmu Komputer, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Lambung Mangkurat. Tak lupa pula penulis panjatkan shalawat dan salam ke hadirat Rasulullah Muhammad SAW beserta para sahabat, keluarga, dan pengikut beliau hingga *yaumul qiyamah*.

Pada lembar ini penulis ingin menyampaikan ucapan terimakasih kepada pihak-pihak yang sangat mendukung penulis dalam pembuatan dan penyusunan skripsi ini, adapun yang dimaksud adalah sebagai berikut :

1. Keluarga, terutama kedua orang tua tercinta untuk mama dan abah yang dari awal proses perkuliahan selalu memberikan bantuan, semangat, doa dan dukungan hingga sampai pada tahap proses penyelesaian skripsi ini.
2. Bapak Muliadi, S.Kom., M.Cs selaku dosen pembimbing utama yang turut serta membantu dan meluangkan waktu demi kelancaran penyelesaian skripsi ini.
3. Bapak Friska Abadi, S.Kom M.Kom selaku dosen pembimbing pendamping yang turut serta membantu dan meluangkan waktu demi kelancaran dalam penyelesaian skripsi ini.
4. Bapak Irwan Budiman, S.T., M.Kom selaku Koordinator Program Studi Ilmu Komputer FMIPA ULM, atas bantuan dan izin beliau skripsi ini dapat diselesaikan.
5. Seluruh Dosen dan Staff Program Studi Ilmu Komputer FMIPA ULM atas ilmu dan bantuan yang diberikan selama ini yang sangat bermanfaat.
6. Teman-teman yang telah bersedia membantu memecahkan kebingungan-kebingungan dan memberikan saran, dukungan, semangat, serta motivasi, selama proses penyelesaian skripsi.

7. Teman-teman keluarga Ilmu Komputer angkatan 2017 yang memberikan dukungan dan bantuan semasa kuliah dan selama proses penyelesaian skripsi.
8. Semua pihak yang tidak dapat disebutkan satu persatu yang telah turut membantu dalam penyelesaian skripsi ini.

Akhir kata penulis menyadari sepenuhnya bahwa penulisan ini jauh dari sempurna, namun penulis mengharapkan bantuan berupa saran dan kritik yang membangun dari semua pihak demi kesempurnaan dan mutu penulisan skripsi ini.

Semoga tulisan ini dapat bermanfaat bagi ilmu pengetahuan dan pembaca khususnya serta mendapat keridhaan Allah SWT.

Banjarbaru, 13 Mei 2024



Muhammad Zamzam

## DAFTAR ISI

HALAMAN JUDUL.....	i
HALAMAN PENGESAHAN.....	ii
HALAMAN PERNYATAAN .....	iii
ABSTRAK .....	iv
<i>ABSTRACT</i> .....	v
PRAKATA.....	vi
DAFTAR ISI.....	viii
DAFTAR TABEL.....	x
DAFTAR GAMBAR .....	xi
DAFTAR LAMPIRAN.....	xii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	2
1.3 Batasan Masalah .....	3
1.4 Tujuan.....	3
1.5 Manfaat Penelitian.....	3
BAB II TINJAUAN PUSTAKA.....	4
2.1 Kajian Terdahulu .....	4
2.2 Keaslian Penelitian .....	7
2.3 Landasan Teori .....	11
2.3.1 <i>COVID-19</i> .....	11
2.3.2 Analisis Sentimen.....	12
2.3.3 <i>Data Pre-Processing</i> .....	12

2.3.4	<i>Bag of Words</i> .....	13
2.3.5	<i>Chi-Square</i> .....	13
2.3.6	<i>Decision Tree</i> .....	14
2.3.7	<i>Boosting</i> .....	15
2.3.8	<i>Categorical Boosting</i> .....	15
2.3.9	<i>K-Fold Cross Validation</i> .....	16
2.3.10	<i>Confusion matrix</i> .....	17
BAB III METODE PENELITIAN.....		19
3.1	Bahan dan Alat Penelitian .....	19
3.2	Variabel Penelitian.....	19
3.3	Prosedur Penelitian .....	20
BAB IV HASIL DAN PEMBAHASAN .....		24
4.1	Hasil.....	24
4.1.1	Pengumpulan Data .....	24
4.1.2	<i>Preprocessing Data</i> .....	26
4.1.3	Ekstraksi Fitur <i>Bag of Words</i> .....	33
4.1.4	Seleksi Fitur <i>Chi-Square</i> .....	35
4.1.5	Klasifikasi Catboost .....	37
4.2	Pembahasan .....	40
BAB V PENUTUP.....		45
5.1	Kesimpulan .....	45
5.2	Saran .....	45
DAFTAR PUSTAKA .....		46
LAMPIRAN.....		59

## DAFTAR TABEL

Tabel 1. Tabel keaslian penelitian.....	7
Table 2. Perancangan Penelitian .....	11
Tabel 3. Contoh Kalimat Setiap Dokumen .....	13
Tabel 4. <i>Bag of Words</i> .....	13
Tabel 5. Tabel Confusion Matrix .....	17
Tabel 6. Tweet Positif .....	21
Tabel 7. Tweet Negatif.....	22
Tabel 8. Data Twitter .....	24
Tabel 9. Proses Data <i>Cleaning</i> .....	27
Tabel 10. Hasil Data <i>Cleaning</i> .....	28
Tabel 11. Hasil Data <i>Case Folding</i> .....	29
Tabel 12. Hasil Data <i>Stemming</i> .....	31
Tabel 13. Hasil Data <i>Stopwords</i> .....	33
Tabel 14. Hasil Ekstraksi Fitur <i>Bag of Words</i> .....	35
Tabel 15. Hasil Seleksi Fitur <i>Chi-Square</i> .....	36
Tabel 16. Hasil <i>Confusion Matrix</i> (Tanpa Seleksi Fitur).....	38
Tabel 17. Hasil Akurasi Fold 1 – 10 (Data Tanpa Seleksi Fitur).....	39
Tabel 18. Hasil <i>Confusion Matrix</i> (Dengan Seleksi Fitur) .....	39
Tabel 19. Hasil Akurasi Fold 1 – 10 (Data Dengan Seleksi Fitur) .....	40
Tabel 20. Perbandingan Hasil <i>Confusion Matrix</i> .....	42
Tabel 21. Perbandingan Hasil Akurasi Nilai Fold 1-10.....	43

## DAFTAR GAMBAR

Gambar 1. Struktur pohon pada decision tree.....	15
Gambar 2. Representasi 10 <i>folds cross validation</i> .....	17
Gambar 3. Alur Penelitian.....	20
Gambar 4. Perbandingan Hasil Akurasi Nilai Fold 1-10 .....	43

## **DAFTAR LAMPIRAN**

Lampiran 1 Kamus *Stopwords*

Lampiran 2 Kode Program