

**PERBANDINGAN KOMBINASI MODEL ALGORITMA *NAIVE*  
*BAYES* DENGAN TEKNIK PEMBOBOTAN KATA  
DALAM ANALISIS SENTIMEN**

**SKRIPSI**



disusun oleh

**Muttafi'ah**

**17.11.1236**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2021**

**PERBANDINGAN KOMBINASI MODEL ALGORITMA *NAÏVE*  
*BAYES* DENGAN TEKNIK PEMBOBOTAN KATA  
DALAM ANALISIS SENTIMEN**

**SKRIPSI**

untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Muttafi'ah**

**17.11.1236**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
2021**

**PERSETUJUAN**

**SKRIPSI**

**PERBANDINGAN KOMBINASI MODEL ALGORITMA *NAÏVE*  
BAYES DENGAN TEKNIK PEMBOBOTAN KATA  
DALAM ANALISIS SENTIMEN**

yang dipersiapkan dan disusun oleh

**Muttafi'ah  
17.11.1236**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 21 November 2020

**Dosen Pembimbing,**

**Rizqi Sukma Kharisma, M.kom  
NIK. 190302215**

**PENGESAHAN**

**SKRIPSI**

**PERBANDINGAN KOMBINASI MODEL ALGORITMA *NAÏVE*  
BAYES DENGAN TEKNIK PEMBOBOTAN KATA  
DALAM ANALISIS SENTIMEN**

yang dipersiapkan dan disusun oleh

**Muttafi'ah**

**17.11.1236**

telah dipertahankan di depan Dewan Penguji  
pada tanggal 17 Juni 2021

**Susunan Dewan Penguji**

**Nama Penguji**

**Wiwi Widayani, M.Kom**  
**NIK. 190302272**

**Muhammad Tofa Nurcholis, M.Kom**  
**NIK. 190302281**

**Rizqi Sukma Kharisma, M.Kom**  
**NIK. 190302215**

**Tanda Tangan**

Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 18 Juni 2021

**DEKAN FAKULTAS ILMU KOMPUTER**

**Hanif Al Fatta, M.Kom**  
**NIK. 190302096**

## PERNYATAAN

Saya yang bertandatangan dibawah ini menyatakan bahwa, skripsi ini merupakan karya saya sendiri (ASLI), dan isi dalam skripsi ini tidak terdapat karya yang pernah diajukan oleh orang lain untuk memperoleh gelar akademis di suatu institusi pendidikan tinggi manapun, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis dan/atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Segala sesuatu yang terkait dengan naskah dan karya yang telah dibuat adalah menjadi tanggungjawab saya pribadi.

Yogyakarta, 26 Juni 2021



Muttafi'ah

NIM. 17.11.1236

## MOTTO

“Hargai waktumu, karna 5 menit dalam hidupmu dapat mengubah duniamu”



## KATA PENGANTAR

Puji dan syukur penulis persembahkan atas kehadiran Allah SWT yang telah melimpahkan rahmat, taufik, dan hidayah sehingga penulis dapat menyelesaikan skripsi ini. Tak lupa pula shalawat serta salam penulis haturkan kepada junjungan Nabi Muhammad SAW, yang telah menjadi teladan bagi seluruh umatnya.

Skripsi ini disusun untuk menyelesaikan studi jenjang Strata Satu (S1) pada program studi Informatika fakultas Ilmu Komputer Universitas Amikom Yogyakarta. Selain itu juga sebagai bukti bahwa mahasiswa telah memenuhi salah satu syarat untuk memperoleh gelar Sarjana Komputer.

Dalam penyusunan skripsi ini, penulis menyadari bahwa masih banyak kekurangan dan jauh dari kata sempurna. Selain itu dengan terselesaikannya skripsi ini, penulis ingin menyampaikan rasa terima kasih kepada:

1. Allah SWT atas segala berkah dan karunia serta ridho-Nya sehingga penulis dapat menyelesaikan skripsi ini dengan baik.
2. Nabi Muhammad SAW yang selalu menjadi panutan dan suri tauladan.
3. Kedua orang tua dan keluarga penulis yang selalu memberikan dukungan dalam bentuk apapun.
4. Bapak Prof. Dr. M. Suyanto, M.M., selaku Rektor Universitas Amikom Yogyakarta.
5. Bapak Rizqi Sukma Kharisma, M.kom, selaku dosen pembimbing yang telah memberikan bimbingan, saran dan arahan, serta masukan kepada penulis.

6. Seluruh staf pengajar dan karyawan Universitas Amikom Yogyakarta yang telah memberikan ilmu yang bermanfaat.
7. Dia yang selalu menemani, memberi semangat dan bantuan kepada saya.
8. Teman-teman 17 IF 05, dan gimang yang menemani dan memberikan semangat.
9. Serta semua pihak yang telah membantu penulis baik dukungan moral maupun materi, pikiran dan tenaga dalam menyelesaikan skripsi ini.

Semoga Allah SWT memberikan balasan lebih kepada semua yang telah ikut membantu penulis hingga terselesainya skripsi ini. Semoga skripsi ini dapat bermanfaat bagi penulis maupun pembaca.

Yogyakarta, 14 Juni 2021



Muttafi'ah



## PERSEMBAHAN

Dengan penuh rasa syukur, segala puji Allah SWT atas rahmat dan karunia yang telah diberikan. Skripsi ini saya persembahkan kepada semua pihak yang terlibat secara langsung maupun tidak langsung dalam proses pembuatan skripsi.

1. Kedua orang tua dan kakak – kakak saya yang selalu mendoakan, memberi motivasi dan dukungan kepada saya.
2. Bapak Rizqi Sukma Kharisma, M.kom selaku dosen pembimbing yang telah membimbing saya dari awal hingga akhir pembuatan skripsi ini.
3. Dosen-dosen Universitas Amikom Yogyakarta yang telah memberikan banyak ilmu pengetahuan kepada saya selama perkuliahan.
4. Teman seperjuangan saya, Benedicta Kristi yang rela meluangkan waktu menemani saya dan membantu memberikan solusi, semangat, serta motivasi selama pembuatan skripsi.
5. Teman – teman Gimang, Destri Herliana Irianti, Akbar Hari Mukti, Muhammad Gufron Hawaly, Herlandro Tribiakto, Jordan Kurnia, dan Rhaka Noviansyah D. yang telah menjadi sahabat dan keluarga baru, dari awal perkuliahan hingga saat ini.
6. Teman-teman kelas 17 IF 05 yang selalu menemani perkuliahan.
7. Martha Presina yang selalu menemani, memberikan semangat, dorongan, motivasi dan menjadi tempat berkeluh kesah bagi saya.

## DAFTAR ISI

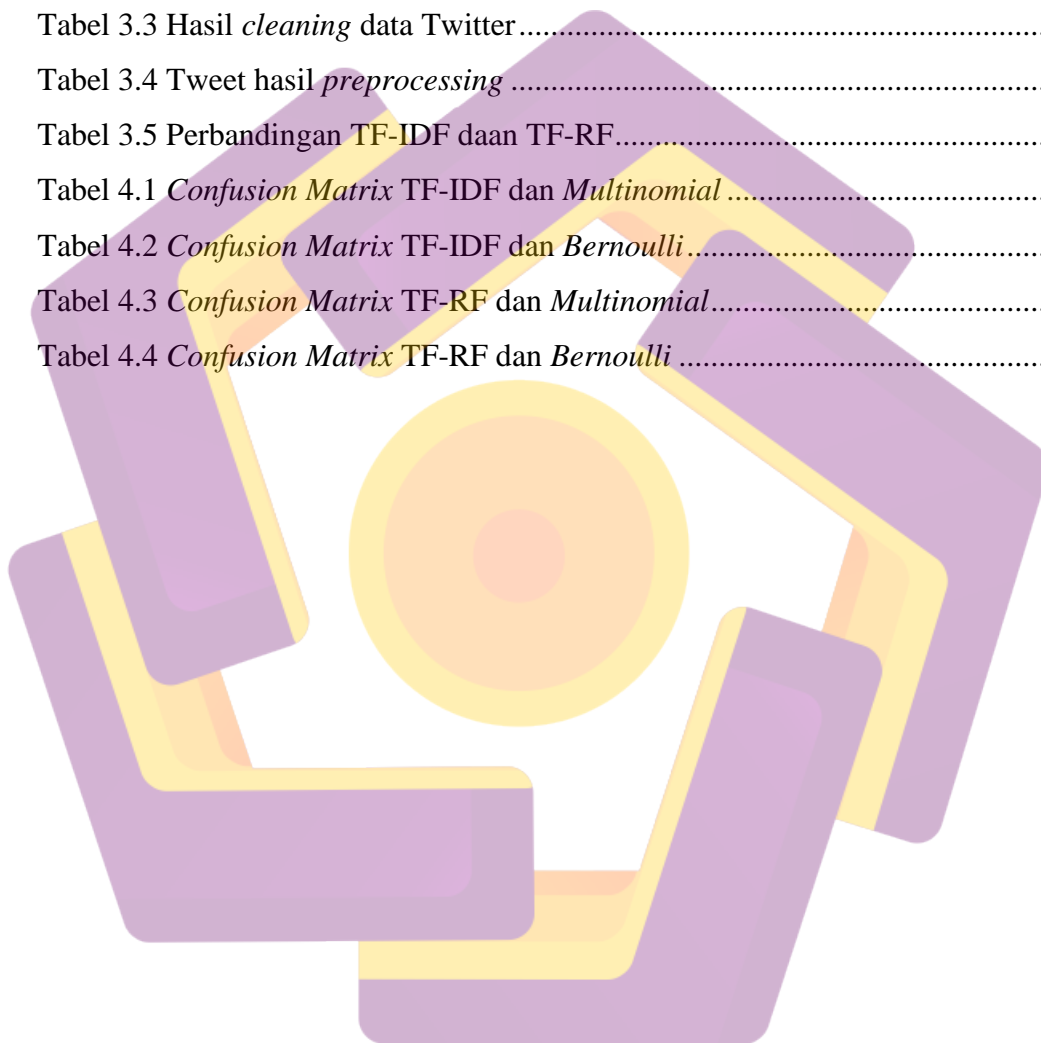
HALAMAN JUDUL .....	i
PERSETUJUAN .....	ii
PENGESAHAN .....	iii
PERNYATAAN .....	iv
MOTTO .....	v
PERSEMBAHAN .....	viii
KATA PENGANTAR .....	vi
DAFTAR ISI .....	viii
DAFTAR TABEL .....	xii
DAFTAR GAMBAR .....	xiii
INTISARI .....	xv
<i>ABSTRACT</i> .....	xvi
BAB I PENDAHULUAN .....	1
1.1 Latar Belakang .....	1
1.2 Rumusan Masalah .....	3
1.3 Batasan Masalah .....	4
1.4 Tujuan Penelitian .....	4
1.5 Manfaat Penelitian .....	5
1.6 Metode Penelitian .....	5
1.6.1 Pengumpulan Data .....	6
1.6.2 Metode Analisis .....	7
1.6.3 Metode Pengujian .....	8

1.7	Sistematika Penulisan.....	9
<b>BAB II LANDASAN TEORI.....</b>		<b>11</b>
2.1	Kajian Pustaka.....	11
2.2	Dasar Teori.....	15
2.2.1	Analisis Sentimen .....	15
2.2.2	Data Mining .....	17
2.2.3	Machine Learning .....	19
2.2.4	Text Mining .....	20
2.2.5	Preprocessing .....	21
2.2.6	Pembobotan kata ( <i>Term Weighting</i> ) .....	22
2.2.7	<i>Naive Bayes Classifier</i> .....	27
2.2.8	<i>Confusion Matrix</i> .....	31
2.2.9	<i>K-Fold Cross Validation</i> .....	35
<b>BAB III METODE PENELITIAN .....</b>		<b>36</b>
3.1	Gambaran Umum Penelitian .....	36
3.2	Alat dan Bahan .....	37
3.2.1	Alat.....	37
3.2.2	Bahan .....	38
3.3	Pengumpulan Data .....	38
3.3.1	Crawling Data .....	38
3.4	Pelabelan Data.....	41
3.5	Preprocessing Data.....	44
3.5.1	Cleaning Data.....	44
3.5.2	Tokenization .....	46
3.5.3	Stopword Removal.....	46

3.5.4	Stemming .....	47
3.6	Pembobotan Kata ( <i>Term Weighting</i> ).....	48
3.6.1	<i>Term Frequency–Inverse Document Frequency</i> (TF-IDF).....	49
3.6.2	<i>Term Frequency–Relevance Frequency</i> (TF–RF) .....	50
3.7	Klasifikasi <i>Naive Bayes Classifier</i> .....	52
3.7.1	<i>Tuning Parameter</i> .....	52
3.7.2	Mengaplikasikan Hasil Tuning Parameter pada Model Klasifikasi	54
3.7.3	Klasifikasi Algoritma Menggunakan Model Tuning .....	56
3.7.4	Save Model .....	58
3.8	Perancangan Sistem.....	59
3.8.1	Flowchart .....	60
3.8.2	Perancangan Antarmuka Pengguna .....	60
BAB IV HASIL DAN PEMBAHASAN .....		65
4.1	Pengujian.....	65
4.1.1	<i>Confusion Matrix</i> .....	65
4.1.2	<i>K-Fold Cross Validation</i> .....	70
4.2	Antarmuka .....	73
4.2.1	Halaman Beranda.....	73
4.2.2	Halaman Prediksi Sentimen.....	74
4.2.3	Halaman Komparasi Model .....	76
BAB V PENUTUP .....		80
5.1	Kesimpulan.....	80
5.2	Saran.....	81
DAFTAR PUSTAKA .....		82
LAMPIRAN.....		85

## DAFTAR TABEL

Tabel 3.1 Hasil <i>crawling</i> data <i>Tweet</i> .....	40
Tabel 3.2 Hasil <i>labeling</i> data Twitter .....	42
Tabel 3.3 Hasil <i>cleaning</i> data Twitter.....	45
Tabel 3.4 <i>Tweet</i> hasil <i>preprocessing</i> .....	47
Tabel 3.5 Perbandingan TF-IDF dan TF-RF.....	50
Tabel 4.1 <i>Confusion Matrix</i> TF-IDF dan <i>Multinomial</i> .....	66
Tabel 4.2 <i>Confusion Matrix</i> TF-IDF dan <i>Bernoulli</i> .....	67
Tabel 4.3 <i>Confusion Matrix</i> TF-RF dan <i>Multinomial</i> .....	68
Tabel 4.4 <i>Confusion Matrix</i> TF-RF dan <i>Bernoulli</i> .....	69



## DAFTAR GAMBAR

Gambar 1.1 Alur metode penelitian.....	6
Gambar 3.1 <i>Script crawling</i> data Twitter .....	39
Gambar 3.2 <i>Script labeling</i> data Twitter.....	42
Gambar 3.3 <i>Script cleaning</i> data Twitter.....	45
Gambar 3.4 <i>Script tokenizing, stopword removal, dan stemming</i> data Twitter ....	47
Gambar 3.5 <i>Script</i> pembagian data latih dan data uji .....	49
Gambar 3.6 <i>Script</i> pembobotan TF-IDF.....	49
Gambar 3.7 <i>Script</i> pembobotan kata TF-RF.....	50
Gambar 3.8 Tuning parameter .....	52
Gambar 3.9 Hasil tuning skema 1 .....	53
Gambar 3.10 Hasil tuning skema 2.....	53
Gambar 3.11 Hasil tuning skema 3.....	54
Gambar 3.12 Hasil tuning skema 4.....	54
Gambar 3.13 Pengaplikasian hasil tuning parameter.....	55
Gambar 3.14 Klasifikasi TF-IDF dan <i>Multinomial</i> .....	56
Gambar 3.15 Klasifikasi TF-IDF dan <i>Bernoulli</i> .....	57
Gambar 3.16 Klasifikasi TF-RF dan <i>Multinomial</i> .....	57
Gambar 3.17 Klasifikasi TF-RF dan <i>Bernoulli</i> .....	58
Gambar 3.18 <i>Script</i> penyimpanan model.....	59
Gambar 3.19 <i>Flowchart</i> sistem.....	60
Gambar 3.20 Rancangan halaman antarmuka beranda.....	61
Gambar 3.21 Halaman antarmuka hasil prediksi .....	62
Gambar 3.22 Halaman antarmuka prediksi sentimen .....	62
Gambar 3.23 Halaman antarmuka komparasi model.....	63
Gambar 3.24 Halaman antarmuka hasil komparasi model .....	64
Gambar 3.25 Halaman antarmuka tentang.....	64
Gambar 4.1 Hasil performa TF-IDF dan <i>Multinomial</i> .....	66
Gambar 4.2 <i>Classification report</i> TF-IDF dan <i>Multinomial</i> .....	66
Gambar 4.3 Hasil performa TF-IDF dan <i>Bernoulli</i> .....	67

Gambar 4.4 <i>Classification report</i> TF-IDF dan <i>Bernoulli</i> .....	68
Gambar 4.5 Hasil performa TF-RF dan <i>Multinomial</i> .....	69
Gambar 4.6 <i>Classification report</i> TF-RF dan <i>Multinomial</i> .....	69
Gambar 4.7 Hasil performa TF-RF dan <i>Bernoulli</i> .....	70
Gambar 4.8 <i>Classification report</i> TF-RF dan <i>Bernoulli</i> .....	70
Gambar 4.9 <i>K-Fold Cross Validation</i> TF-IDF dan <i>Multinomial</i> .....	70
Gambar 4.10 <i>K-Fold Cross Validation</i> TF-IDF dan <i>Bernoulli</i> .....	71
Gambar 4.11 <i>K-Fold Cross Validation</i> TF-RF dan <i>Multinomial</i> .....	72
Gambar 4.12 <i>K-Fold Cross Validation</i> TF-RF dan <i>Bernoulli</i> .....	72
Gambar 4.13 Antarmuka halaman beranda .....	74
Gambar 4.14 Antarmuka prediksi sentimen .....	75
Gambar 4.15 Antarmuka halaman hasil Prediksi .....	76
Gambar 4.16 Antarmuka halaman komparasi model .....	77
Gambar 4.17 Antarmuka halaman hasil upload data .....	77
Gambar 4.18 Antarmuka hasil preprocessing .....	78
Gambar 4.19 Antarmuka hasil komparasi model .....	79

## INTISARI

*Text mining* merupakan konsep dasar analisis sentimen dan disiplin ilmu yang menggabungkan ilmu bahasa dan ilmu komputer dengan teknik pembelajaran mesin (*machine learning*). *Text mining* digunakan untuk mengubah teks menjadi lebih terstruktur. Sedangkan pembelajaran mesin berfokus untuk mencari dan mengembangkan algoritma untuk membangun sebuah sistem yang dapat mensimulasikan atau meniru sebuah pola dari kumpulan data. Dalam penelitian ini menggunakan teknik supervised learning yang merupakan teknik dasar pembelajaran mesin dengan mengkomparasikan model algoritma *Naive Bayes Classifier* yaitu *Multinomial Naive Bayes* dan *Bernoulli Naive Bayes* dengan objek sentimen dari Twitter. Pada penelitian ini juga menggunakan teknik pembobotan kata yaitu TF-IDF dan TF-RF pada masing – masing model. Penelitian ini dilakukan untuk mengetahui kombinasi terbaik dari setiap model dengan pembobotan kata.

Tahap awal dalam penelitian ini adalah *crawling data* menggunakan API Twitter, kemudian data tersebut *labeling*. Setelah data diberi label, data tersebut akan masuk ke tahap penting dalam penelitian, yaitu *preprocessing* dan pembobotan kata. Data yang telah dilabeli dibersihkan dan diubah menjadi data yang terstruktur sehingga data siap untuk dianalisis. Data hasil *preprocessing* diberi bobot dengan teknik TF-IDF dan TF-RF, kemudian diklasifikasi satu-persatu menggunakan 2 model NBC sehingga dalam penelitian ini terdapat 4 skema model, yaitu Multinomial dan TF-IDF, Bernoulli dan TF-IDF, Multinomial dan TF-RF serta Bernoulli dan TF-RF. Tahap terakhir dari penelitian ini adalah pengujian menggunakan *Confusion Matrix* kemudian divalidasi dengan *K-Fold Cross Validation*, pengujian dilakukan untuk melihat performa terbaik dari 4 skema tersebut.

Dari 4 skema yang dilakukan, skema TF-IDF dan TF-RF dengan *Bernoulli Naive Bayes* dari hasil pengujian *Confusion Matrix* menghasilkan akurasi terbaik yaitu 61%, dan rata - rata nilai akurasi dari validasi *5-fold* sebesar 60%. Dan yang memiliki nilai akurasi terendah terletak pada model *Multinomial Naive Bayes* dan TF-IDF yaitu 58% dari *Confusion Matrix*, dengan rata-rata nilai akurasi 59% dari validasi *5-fold*.

**Kata-kunci:** Multinomial Naive Bayes, Bernoulli Naive Bayes, TF-IDF, TF-RF, Analisis sentimen



## ***ABSTRACT***

*Text mining is a basic concept of sentiment analysis and a discipline that combines linguistics and computer science with machine learning techniques. Text mining is used to change the text to be more structured. While machine learning focuses on finding and developing algorithms to build a system that can simulate or imitate a pattern from a dataset. In this study, supervised learning is used which is a basic machine learning technique with comparing the Naive Bayes Classifier algorithm model, namely Multinomial Naive Bayes and Bernoulli Naive Bayes with sentiment objects from Twitter. This study also uses Term Weighting techniques, namely TF-IDF and TF-RF in each model. This study was conducted to determine the best combination of each model with Term Weighting.*

*The first step in this research is crawling the data using the Twitter API, then the data is labeled. After the data is labeled, the data will enter an important step in the research, namely preprocessing and term weighting. The data that has been labeled is cleaned and converted into structured data so that the data is ready for analysis. The preprocessing data are weighted using the TF-IDF and TF-RF techniques, then classified one by one using 2 NBC models, so in this study there are 4 model schemes, namely Multinomial and TF-IDF, Bernoulli and TF-IDF, Multinomial and TF-RF as well as Bernoulli and TF-RF. The last stage of this research is testing using Confusion Matrix, and then validated with K-Fold Cross Validation, testing is carried out to see the best performance of the 4 schemes.*

*Of the 4 schemes, TF-IDF and TF-RF with Bernoulli Naive Bayes schemes from the results of the Confusion Matrix test produce the best accuracy 61%, and the average accuracy value of the 5-fold validation is 60%. And the one with the lowest accuracy value lies in the Multinomial Naive Bayes model and TF-IDF which is 58% from Confusion Matrix, with an average value of 59% from the 5-fold validation.*

**Keywords:** *Multinomial Naive Bayes, Bernoulli Naive Bayes, TF-IDF, TF-RF, Sentiment analysis*