



**KLASIFIKASI DAN ANALISIS SENTIMEN PADA DATA
TWITTER MENGGUNAKAN ALGORTIMA NAÏVE BAYES
(STUDI KASUS: PEKAN OLAHRAGA NASIONAL XX 2021)**

SKRIPSI

**KRISNA JONATHAN SITORUS
1810511090**

**UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
FAKULTAS ILMU KOMPUTER
PROGRAM STUDI INFORMATIKA
2022**



**KLASIFIKASI DAN ANALISIS SENTIMEN PADA DATA
TWITTER MENGGUNAKAN ALGORTIMA NAÏVE BAYES
(STUDI KASUS: PEKAN OLAHRAGA NASIONAL XX 2021)**

SKRIPSI

**Diajukan Sebagai Salah Satu Syarat untuk Memperoleh Gelar
Sarjana Komputer**

KRISNA JONATHAN SITORUS

1810511090

**UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
FAKULTAS ILMU KOMPUTER
PROGRAM STUDI INFORMATIKA
2022**

PERNYATAAN ORISINALITAS

Skripsi ini adalah hasil karya sendiri, dan sumber yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar.

Nama : Krisna Jonathan Sitorus
NIM : 1810511090
Tanggal : 13 Juni 2022

Bilamana dikemudian hari ditemukan ketidaksamaan dengan pernyataan ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Jakarta, 13 Juni 2022



Krisna Jonathan Sitorus

PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS

Sebagai civitas akademik Universitas Pembangunan Nasional Veteran Jakarta, saya yang bertanda tangan dibawah ini:

Nama : Krisna Jonathan Sitorus

NIM : 1810511090

Fakultas : Ilmu Komputer

Program Studi : Informatika

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Pembangunan Nasional Veteran Jakarta Hak Bebas Royalti Non Eksklusif (*Non-exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

Klasifikasi dan Analisis Sentimen Pada Data Twitter Menggunakan Algoritma Naïve Bayes. (Studi Kasus: Pekan Olahraga Nasional XX 2021)

Beserta perangkata yang ada (jika diperlukan). Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional Veteran Jakarta berhak menyimpan, mengalih-meida/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan Tugas akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Jakarta

Pada tanggal : 22 Juli 2022

Yang menyatakan,



(Krisna Jonathan Sitorus)

LEMBAR PENGESAHAN

Dengan ini dinyatakan bahwa Skripsi berikut:

Nama : Krisna Jonathan Sitorus
NIM : 1810511090
Program Studi : S1 Informatika
Judul : Klasifikasi dan Analisis Sentimen Pada Data Twitter Menggunakan Algoritma Naïve Bayes. (Studi Kasus: Pekan Olahraga Nasional XX 2021)
Telah berhasil dipertahankan di hadapan Tim penguji dan diterima sebagai bagian persyaratan yang diperlukan untuk memperoleh gelar Sarjana Komputer pada Program Studi S1 Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional Veteran Jakarta.

Henki Bayu Setu, S.Kom., MTI.

Penguji 1

Mayanda Mega Santoni, S.Kom., M.Kom.

Penguji 2

Anita Muliawati, S.Kom., MTI.

Dosen Pembimbing 1



Dr. Ermawita, M.Kom.

Dekan

Sarika, S.Kom., M.Kom.

Dosen Pembimbing 2

Desta Sandya Prasvita, S.Kom., M.Kom.

Ketua Program Studi

Ditetapkan di : Jakarta
Tanggal Pengesahan : 20 Juli 2022



KLASIFIKASI DAN ANALISIS SENTIMEN PADA DATA TWITTER MENGGUNAKAN ALGORITMA NAÏVE BAYES. (STUDI KASUS: PEKAN OLAHRAGA NASIONAL XX 2021)

Krisna Jonathan Sitorus

ABSTRAK

Dunia teknologi sekarang ini sedang mengalami perkembangan yang cukup signifikan, masyarakat kini bisa dengan bebas mengakses dan memanfaatkan media sosial sebagai contoh yakni Twitter. Banyak sekali pengguna Twitter mencurahkan pendapatnya melalui tweet-tweet yang mereka kirimkan pada media sosial tersebut, khususnya mengenai Pekan Olahraga Nasional (PON) XX 2021 silam. Banyak sekali tweet yang bersifat mendukung, tetapi tidak jarang juga ada tweet yang bersifat keluhan mengenai penyelenggaraan PON XX tersebut. Dari masalah tersebut, dilaksanakanlah penelitian mengenai analisis sentimen pada data twitter yang berkaitan dengan PON XX dan mempergunakan metode Naïve Bayes. Jumlah data yang di crawling sebanyak 1000 data lalu data melalui proses *dataduplicate removal* sehingga memperoleh hasil data sebanyak 218 data tweet dan belum terlabelkan. Sebelum proses mengklasifikasikan data yang diperoleh, data harus dilaksanakan pemberian label pada datanya serta pembersihan data terlebih dulu sebelum masuk pada tahapan *text processing*, selanjutnya data diberi bobot pada tiap kata dengan *Term Frequency– Inverse Document Frequency* (TF-IDF) yang akan kedepannya kata tersebut akan dijadikan sebagai fitur. Lalu, karena data berlabel positif serta negatif mempunyai jumlah yang jauh berbeda, maka dimanfaatkan metode *Synthetic Minority Oversampling Technique* (SMOTE) guna melaksanakan penyeimbangan terhadap datanya. Tahapan selanjutnya dilaksanakan pembagian data yang besarnya yakni 80% 20% dan diklasifikasikan dengan metode Naive Bayes. Hasil yang diperoleh dari pelaksanaan penelitiannya tersebut ialah diperoleh bahwa data uji memperoleh accuracy yang besaran persentasenya yakni 99%, precision dengan besaran persentasenya 100%, recall dengan besaran persentasenya 98%.

Kata Kunci: Analisis Sentimen, Klasifikasi, PON XX, Naïve Bayes

CLASSIFICATION AND SENTIMENT ANALYSIS ON TWITTER DATA USING NAÏVE BAYES. (CASE STUDY: PEKAN OLAHRAGA NASIONAL XX 2021)

Krisna Jonathan Sitorus

ABSTRACT

The world of technology is currently developing very rapidly, people can now freely access and express themselves using social media as an example, namely Twitter. Many Twitter users express their opinions through the tweets they send on social media, especially regarding the XX 2021 National Sports Week (PON) ago. Several types of tweets that are appreciative of the event are often seen, but not infrequently there are also tweets that are complaining about the implementation of the XX PON. From this problem, a research was conducted on sentiment analysis on twitter data regarding PON XX using the Naïve Bayes method. The amount of data that is crawled is 1000 data and then the data goes through a duplicate data removal process so that it produces 218 tweets and has not been labeled. Before the process of classifying the data obtained, the data must be labeled and cleaned first before entering the text processing stage, then the data is given a weight for each word with a Term Frequency–Inverse Document Frequency (TF-IDF) which will be used in the future as a feature. Then, because the data with positive and negative labels have significantly different amounts, the Synthetic Minority Oversampling Technique (SMOTE) method was used to balance the data. The next stage is the distribution of data by 80% 20% and classified by the Naïve Bayes method. The results obtained from this study are obtained that the test data get 99% accuracy, 100% precision, 98% recall.

Keyword: Sentiment Analysis, Classification, PON XX, Naïve Bayes

KATA PENGANTAR

Segala puji dan syukur penulis panjatkan kehadiran Tuhan Yang Mahakuasa atas limpahan rahmat dan karunia-Nya sehingga penulis dapat menyelesaikan Skripsi (Tugas Akhir). Skripsi ini disusun sebagai syarat kelulusan Studi Informatika Fakultas Ilmu Komputer Universitas Pembangunan Nasional Veteran Jakarta.

Dalam penyelesaian karya tulis ilmiah ini tidak lepas dari bantuan banyak pihak yang telah memberikan masukan kepada penulis. Untuk itu penulis mengucapkan terima kasih kepada:

1. Dr. Ermatita, M.Kom., selaku dekan Fakultas Ilmu Komputer
2. Desta Sandya Prasvita, S.Kom, M.Kom., selaku Ketua Program Studi Sarjana Jurusan S1 Informatika.
3. Anita Muliawati, S.Kom., MTI. selaku dosen pembimbing 1 dari pihak jurusan.
4. Sarika, S.Kom., M.Kom. selaku dosen pembimbing 2 dari pihak jurusan.
5. Orang tua yang telah memberikan dukungan baik secara moril maupun materil.
6. Seluruh pihak yang terlibat dalam kelancaran pembuatan makalah karya ilmiah ini dan yang belum disebutkan di atas, penulis ucapkan terima kasih.

Penulis menyadari bahwa masih banyak kekurangan dari laporan ini, baik dari materi maupun teknik penyajiannya. Oleh karena itu, kritik dan saran yang membangun penulis harapkan.

Jakarta, 8 Januari 2021

Krisna Jonathan Sitorus

DAFTAR ISI

PERNYATAAN ORISINALITAS	ii
LEMBAR PERSETUJUAN	iii
LEMBAR PENGESAHAN	iv
ABSTRAK	v
ABSTRACT	vi
KATA PENGANTAR	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	xii
DAFTAR GAMBAR	xii
DAFTAR SIMBOL	xiii
BAB I PENDAHULUAN.....	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah	2
1.3. Ruang Lingkup.....	2
1.4. Tujuan	2
1.5. Manfaat	3
1.6. Luaran yang Diharapkan	3
1.7. Sistematika Penulisan.....	3
BAB II TINJAUAN PUSTAKA	5
2.1. Twitter	5
2.2. Analisis Sentimen.....	5
2.3. Fleiss Kappa	6
2.4. Text Analysis	7
2.4.1 Text Mining.....	7

2.4.2	Pra Proses Data	7
2.5.	Pembobotan Kata	9
2.5.1	TF (<i>Term Frequency</i>).....	9
2.5.2	IDF (<i>Inverse Document Frequency</i>)	9
2.6.	Metode Klasifikasi	9
2.6.1	Algoritma Naïve Bayes	9
2.7.	Python	10
2.8.	Penelitian Terdahulu	11
	BAB III METODOLOGI PENELITIAN	12
3.1.	Kerangka Pikir	12
3.1.1.	Identifikasi masalah	13
3.1.2.	Studi Literatur	13
3.1.3.	Akuisisi Data.....	13
3.1.4.	Pelabelan Data.....	13
3.1.5.	Praproses Data.....	14
3.1.6.	Pembobotan TF IDF.....	16
3.1.7.	Synthetic Minority Oversampling Technique (SMOTE).....	16
3.1.8.	Klasifikasi	17
3.1.9.	Hasil Klasifikasi.....	18
3.1.10.	Evaluasi.....	18
3.1.11.	Analisis Hasil Klasifikasi.....	19
3.2.	Perangkat Penelitian.....	19
	BAB IV HASIL DAN PEMBAHASAN	20
4.1	Pengumpulan Data	20
4.2	Pelabelan Data.....	20
4.3	Praproses	25
4.3.1	Case Folding	25
4.3.2	Data Cleaning.....	26
4.3.3	Normalization.....	28
4.3.4	Stemming	29

4.3.5	Stopwords Removal	30
4.3.6	Tokenizing.....	31
4.4	Pembobotan TF-IDF	32
4.5	SMOTE (Synthetic Minority Oversampling Technique).....	33
4.6	Klasifikasi Naïve Bayes	33
4.6.1	Data Latih.....	33
4.6.2	Data Uji	36
4.7	Evaluasi.....	41
4.8	Analisis Hasil Kegiatan.....	43
4.8.1	<i>Wordcloud</i> Sentimen Positif terhadap Data Tweet PON XX 2021	43
4.8.2	<i>Wordcloud</i> Sentimen Negatif terhadap Data Tweet PON XX 2021	44
BAB V	PENUTUP	45
5.1	Kesimpulan	45
5.2	Saran.....	45
	DAFTAR PUSTAKA	47
	RIWAYAT HIDUP	50
	LAMPIRAN.....	51

DAFTAR TABEL

Tabel 3. 1 Tabel Cleaning	15
Tabel 3. 2 Stopwords Removal.....	16
Tabel 3. 3 Confusion Matrix	18
Tabel 4. 1 Tweet dan Hasil Pelabelan.....	21
Tabel 4. 2 Hasil Penilaian Anotator	23
Tabel 4. 3 Tabel Hasil Perhitungan Nilai Pi	23
Tabel 4. 4 Case Folding	25
Tabel 4. 5 Data Cleaning	27
Tabel 4. 6 Normalization	28
Tabel 4. 7 Stemming	29
Tabel 4. 8 Stopwords Removal.....	30
Tabel 4. 9 Tokenization	31
Tabel 4. 10 Contoh Data TF-IDF.....	32
Tabel 4. 11 Perhitungan TF-IDF.....	32
Tabel 4. 12 Jumlah Data Sebelum dan Sesudah SMOTE.....	33
Tabel 4. 13 Contoh Data Latih.....	34
Tabel 4. 14 Nilai Probabilitas Pada Data Latih.....	36
Tabel 4. 15 Contoh Data Uji	37
Tabel 4. 16 Data Uji yang Setelah Praproses.....	37
Tabel 4. 17 Probabilitas Sampel Data Uji.....	38
Tabel 4. 18 Nilai Probabilitas Kelas Positif.....	39
Tabel 4. 19 Nilai Probabilitas Kelas Negatif	39
Tabel 4. 20 Confusion Matrix Model	41

DAFTAR GAMBAR

Gambar 3. 1 Tahapan Penelitian	12
Gambar 4. 1 Hasil Crawling Data Twitter	20
Gambar 4. 2 Wordcloud Sentimen Positif	43
Gambar 4. 3 Wordcloud Sentimen Negatif.....	44

DAFTAR SIMBOL

Simbol Flowchart			
NO	GAMBAR	NAMA	KETERANGAN
1		<i>Terminal (start, end)</i>	MengGambarkan bagaimana kegiatan dimulai atau kegiatan berakhir.
2		<i>Flow Direction</i>	MengGambarkan hubungan antar simbol yang menyatakan suatu jalannya proses dalam sistem.
3		<i>Process</i>	MengGambarkan deskripsi dari proses yang dijalankan.
4		<i>Document</i>	MengGambarkan bahwa masukkan (<i>input</i>) berasal dari sebuah data dokumen yang dapat berupa kertas atau keluaran (<i>output</i>) yang dicetak ke kertas.
5		<i>Predefine process</i>	MengGambarkan pelaksanaan dari sebuah proses atau disebut sebagai subprogram.