

**ANALISIS SENTIMEN PENGGUNA TWITTER TERHADAP
KEBOCORAN DATA REGISTRASI SIM MENGGUNAKAN
METODE NAÏVE BAYES**

SKRIPSI

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi S1 Informatika



diajukan oleh

RAIHAN ALFAIN SHUBHIY

18.11.2333

Kepada

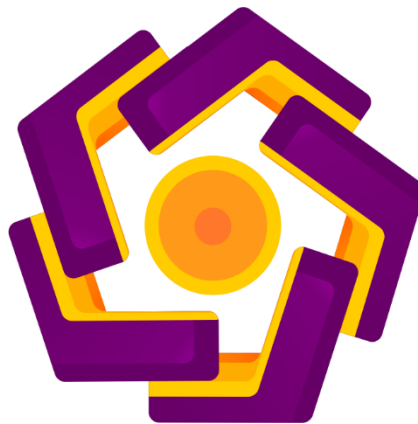
**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2023

**ANALISIS SENTIMEN PENGGUNA TWITTER TERHADAP
KEBOCORAN DATA REGISTRASI SIM MENGGUNAKAN
METODE NAÏVE BAYES**

SKRIPSI

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi S1 Informatika



diajukan oleh

RAIHAN ALFAIN SHUBHIY

18.11.2333

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2023

HALAMAN PERSETUJUAN

SKRIPSI

ANALISIS SENTIMEN PENGGUNA TWITTER TERHADAP

KEBOCORAN DATA REGISTRASI SIM MENGGUNAKAN

METODE NAÏVE BAYES

yang disusun dan diajukan oleh

RAIHAN ALFAIN SHUBHIY

18.11.2333

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 4 April 2023

Dosen Pembimbing,



Theopilus Bayu Sasongko, S.Kom, M.Eng.

NIK. 190302375

HALAMAN PENGESAHAN

SKRIPSI

**ANALISIS SENTIMEN PENGGUNA TWITTER TERHADAP
KEBOCORAN DATA REGISTRASI SIM MENGGUNAKAN
METODE NAÏVE BAYES**

yang disusun dan diajukan oleh
RAIHAN ALFAIN SHUBHIY

18.11.2333

Telah dipertahankan di depan Dewan Penguji
pada tanggal 4 April 2023

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Arif Dwi Laksito, M.Kom.
NIK. 190302150



Supriatin, M.Kom.
NIK. 190302239



Theopilus Bayu Sasongko, S.Kom., M.Eng.
NIK. 190302375



Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 4 April 2023

DEKAN FAKULTAS ILMU KOMPUTER



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

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

Nama mahasiswa : Raihan Alfain Shubhiy
NIM : 18.11.2333

Menyatakan bahwa Skripsi dengan judul berikut:

Analisis Sentimen Pengguna Twitter Terhadap Kebocoran Data Registrasi SIM Menggunakan Metode Naïve Bayes

Dosen Pembimbing : Theopilus Bayu Sasongko, S.Kom, M.Eng.

1. Karya tulis ini adalah benar-benar ASLI dan BELUM PERNAH diajukan untuk mendapatkan gelar akademik, baik di Universitas AMIKOM Yogyakarta maupun di Perguruan Tinggi lainnya.
2. Karya tulis ini merupakan gagasan, rumusan dan penelitian SAYA sendiri, tanpa bantuan pihak lain kecuali arahan dari Dosen Pembimbing.
3. Dalam karya tulis ini tidak terdapat karya atau pendapat orang lain, kecuali secara tertulis dengan jelas dicantumkan sebagai acuan dalam naskah dengan disebutkan nama pengarang dan disebutkan dalam Daftar Pustaka pada karya tulis ini.
4. Perangkat lunak yang digunakan dalam penelitian ini sepenuhnya menjadi tanggung jawab SAYA, bukan tanggung jawab Universitas AMIKOM Yogyakarta.
5. Pernyataan ini SAYA buat dengan sesungguhnya, apabila di kemudian hari terdapat penyimpangan dan ketidakbenaran dalam pernyataan ini, maka SAYA bersedia menerima SANKSI AKADEMIK dengan pencabutan gelar yang sudah diperoleh, serta sanksi lainnya sesuai dengan norma yang berlaku di Perguruan Tinggi.

Yogyakarta, 4 April 2023

Yang Menyatakan,



Raihan Alfain Shubhiy

HALAMAN PERSEMBAHAN

Alhamdulillah, skripsi ini telah saya selesaikan dengan baik dan maksimal. Hal ini tentunya tak lepas dari karunia, rahmat serta hidayah yang telah diberikan oleh Allah swt. Sehingga saya mendapat kemudahan dan kelancaran dalam mengerjakan skripsi ini. Selain itu ada orang-orang hebat yang selalu memberikan dukungan, motivasi, dan bantuan baik secara langsung maupun tidak langsung, antar lain.

1. Orang tua serta keluarga yang selalu memberikan dukungan, motivasi dan doa restu tanpa henti.
2. Bapak Theopilus Bayu Sasongko, S.Kom, M.Eng yang telah membimbing saya dalam peroses pengerjaan skripsi ini.
3. Teman-Teman dari Himpunan Mahasiswa Informatika yang membantu dikala saya mendapati kesulitan dalam pengerjaan skripsi.
4. Shabrina Azzahra, Amd. Keb yang selalu menemani dan membantu saya dalam proses pengerjaan skripsi hingga selesai.
5. Teman-Teman yang tidak bisa saya sebutkan namanya satu-persatu yang telah membantu dan *men-support* baik secara langsung maupun tidak langsung.

KATA PENGANTAR

Puji syukur kehadiran Allah swt. yang telah memberikan karunia, rahmat serta hidayah-nya sehingga penulis dapat menyelesaikan skripsi ini dengan baik dan maksimal. Skripsi yang berjudul **“Analisis Sentimen Pengguna Twitter Terhadap Kebocoran Data Registrasi SIM Menggunakan Metode Naïve Bayes”** ini disusun sebagai salah satu syarat dalam menyelesaikan masa studi program sarjana di Universitas AMIKOM Yogyakarta.

Dalam penulisan skripsi ini, penulis mengucapkan rasa terimakasih atas motivasi, bimbingan, saran dan masukan dari berbagai pihak secara moral maupun spiritual. Pada kesempatan kali ini penulis mengucapkan terimakasih kepada:

1. Bapak Prof. Dr. M. Suyanto, M.M selaku Rektor Universitas AMIKOM Yogyakarta.
2. Bapak Hanif Al Fatta, M.Kom selaku Dekan Fakultas Ilmu Komputer Universitas Amikom Yogyakarta.
3. Ibu Windha Mega Pradya D., M.Kom selaku Ketua Program Studi S1 Informatika Universitas Amikom Yogyakarta.
4. Bapak Theopilus Bayu Sasongko, S.Kom, M.Eng selaku Dosen Pembimbing yang selalu memberikan saran dan masukan dalam proses penulisan skripsi.
5. Bapak Arif Dwi Laksito, M.Kom dan Ibu Supriatin, M.Kom selaku Dosen Penguji yang telah memberikan evaluasi dan saran agar penelitian ini menjadi lebih baik.

Akhir kata, semoga penyusunan skripsi ini dapat bermanfaat bagi pembaca dalam menambah wawasan dan pengetahuan khususnya dalam bidang informatika.

Yogyakarta, 4 April 2023

Penulis

DAFTAR ISI

HALAMAN JUDUL	ii
HALAMAN PERSETUJUAN.....	iii
HALAMAN PERNYATAAN KEASLIAN SKRIPSI	v
HALAMAN PERSEMBAHAN	vi
KATA PENGANTAR	vii
DAFTAR ISI.....	viii
DAFTAR GAMBAR	xi
INTISARI	xii
ABSTRACT.....	xiii
BAB I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah.....	2
1.3 Batasan Masalah	2
1.4 Tujuan Penelitian	2
1.5 Manfaat Penelitian	3
1.6 Metode Penelitian	3
1.6.1 Metode Pengumpulan Data.....	3
1.6.2 Metode Analisis Data	3
1.6.3 Pengujian	3
1.7 Sistematika Penulisan	4
BAB II TINJAUAN PUSTAKA	5
2.1 Studi Literatur	5
2.2 Dasar Teori	12
2.2.1 Analisis Sentimen	12
2.2.2 <i>Machine Learning</i>	12
2.2.3 <i>Twitter</i>	12
2.2.4 <i>Scraping</i>	13
2.2.5 <i>Text Mining</i>	13
2.2.6 <i>Pre-processing</i>	13
2.2.7 Pelabelan Data	14
2.2.8 TF-IDF.....	15
2.2.9 SMOTE.....	16
2.2.10 Naïve Bayes	16
2.2.11 <i>Confusion Matrix</i>	18
2.2.12 K-Fold Cross Validation	20
2.2.13 Python	21
2.2.14 <i>Wordcloud</i>	21

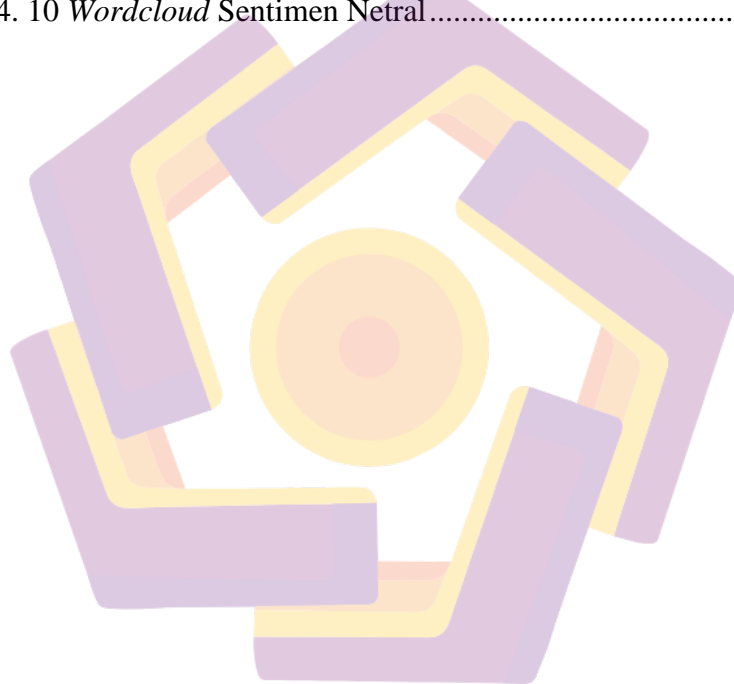
BAB III METODE PENELITIAN	22
3.1 Prosedur Penelitian	22
3.2 Tahapan Penelitian.....	23
3.3 Alat Penelitian.....	24
3.4 Data Penelitian.....	24
3.5 Instrumen Pengambilan Data.....	25
3.6 <i>Pre-Processing</i>	25
3.6.1 <i>Case Folding</i>	25
3.6.2 <i>Cleansing</i>	25
3.6.3 <i>Tokenizing</i>	26
3.6.4 <i>Normalization</i>	26
3.6.5 <i>Stopword Removal</i>	27
3.6.6 <i>Stemming</i>	27
3.7 <i>Labelling</i>	28
3.8 TF-IDF	28
3.9 Multinomial Naïve Bayes	29
3.10 Confusion Matrix	30
BAB IV HASIL DAN PEMBAHASAN	32
4.1 Pengambilan Data Twitter	32
4.2 <i>Preprocessing Data</i>	33
4.2.1 <i>Case Folding</i>	33
4.2.2 <i>Cleansing Data</i>	34
4.2.3 <i>Tokenizing</i>	36
4.2.4 <i>Normalization</i>	37
4.2.5 <i>Stopword Removal</i>	38
4.2.6 <i>Stemming</i>	40
4.3 <i>Labelling</i>	41
4.4 Pembobotan TF-IDF	43
4.5 Pembagian Data	47
4.6 SMOTE.....	48
4.7 Hasil Klasifikasi.....	49
4.7.1 <i>Confusion Matrix</i>	50
4.7.2 <i>K-fold Cross Validation</i>	52
4.8 Hasil Visualisasi	54
4.8.1 Sentimen Positif.....	55
4.8.2 Sentimen Negatif	56
4.8.3 Sentimen Netral	57
BAB V PENUTUP	59
5.1 Kesimpulan	59
5.2 Saran	59
REFERENSI	61

DAFTAR TABEL

Tabel 2. 1 Keaslian Penelitian	8
Tabel 2. 2 <i>Confusion Matrix</i>	19
Tabel 2. 3 <i>K-Fold Cross Validation</i>	20
Tabel 3.1 Alat Penelitian.....	24
Tabel 3.2 Contoh <i>Case Folding</i>	25
Tabel 3.3 Contoh <i>Cleansing</i>	26
Tabel 3.4 Contoh <i>Tokenizing</i>	26
Tabel 3. 5 Contoh <i>Normalization</i>	27
Tabel 3. 6 Contoh <i>Stopword Removal</i>	27
Tabel 3. 7 Contoh <i>Stemming</i>	28
Tabel 3. 8 <i>Confusion Matrix 3 Kelas</i>	31
Tabel 4. 1 Proses <i>Scrapping</i>	32
Tabel 4. 2 Hasil Pengambilan Data.....	33
Tabel 4. 3 Proses <i>Case Folding</i>	34
Tabel 4. 4 Hasil <i>Case Folding</i>	34
Tabel 4. 5 Proses <i>Cleansing Data</i>	35
Tabel 4. 6 Hasil <i>Cleansing Data</i>	36
Tabel 4. 7 Proses <i>Tokenizing</i>	36
Tabel 4. 8 Hasil <i>Tokenizing</i>	36
Tabel 4. 9 Proses <i>Normalization</i>	37
Tabel 4. 10 Hasil <i>Normalization</i>	38
Tabel 4. 11 Proses <i>Stopward Removal</i>	38
Tabel 4. 12 Hasil <i>Stopward Removal</i>	39
Tabel 4. 13 Proses <i>Stemming</i>	40
Tabel 4. 14 Hasil <i>Stemming</i>	41
Tabel 4. 15 Proses <i>Labelling</i>	41
Tabel 4. 16 Hasil <i>Labelling</i> dengan <i>InSet Lexicon</i>	43
Tabel 4. 17 Perhitungan TF	44
Tabel 4. 18 Perhitungan IDF.....	44
Tabel 4. 19 Perhitungan TF-IDF.....	45
Tabel 4. 20 Kode Proses TF-IDF	46
Tabel 4. 21 Proses Pembagian Data.....	47
Tabel 4. 22 Proses <i>oversampling</i> menggunakan SMOTE	48
Tabel 4. 23 Hasil <i>oversampling</i> menggunakan SMOTE	48
Tabel 4. 24 Proses Klasifikasi.....	49
Tabel 4. 25 Hasil Klasifikasi.....	49
Tabel 4. 26 Proses <i>Confusion Matrix</i>	50
Tabel 4. 27 Hasil <i>Confusion Matrix</i>	51
Tabel 4. 28 Proses <i>K-fold Cross Validation</i>	52
Tabel 4. 29 Hasil <i>K-fold Cross Validation</i>	53
Tabel 4. 30 Proses Visualisasi <i>Barplot</i>	54
Tabel 4. 31 Proses Visualisasi <i>Wordcloud</i>	54

DAFTAR GAMBAR

Gambar 3. 1 Tahapan Penelitian	23
Gambar 4. 1 Hasil TF-IDF	47
Gambar 4. 2 Hasil <i>Confusion Matrix</i>	51
Gambar 4. 3 Frekuensi Kata pada Data <i>Tweet</i>	54
Gambar 4. 4 Hasil Visualiasi <i>Wordcloud</i>	55
Gambar 4. 5 Frekuensi Kata Sentimen Positif	56
Gambar 4. 6 <i>Wordcloud</i> Sentimen Positif	56
Gambar 4. 7 Frekunesi Kata Sentimen Negatif	57
Gambar 4. 8 <i>Wordcloud</i> Sentimen Negatif	57
Gambar 4. 9 Frekunesi Kata Sentimen Netral	58
Gambar 4. 10 <i>Wordcloud</i> Sentimen Netral	58



INTISARI

Perkembangan teknologi digital di Indonesia yang semakin masif membuat masyarakat tidak bisa menolak dampak dari digitalisasi yang dapat membawa ancaman besar seperti ancaman keamanan data pribadi. Pada bulan September 2022, terdapat insiden kebocoran data registrasi kartu sim berisikan data pribadi pengguna yang disebarakan melalui situs gelap. Sebagai salah satu media sosial paling populer di Indonesia, Twitter dijadikan tempat oleh masyarakat untuk menyuarakan opininya terkait isu kebocoran data registrasi sim. Penelitian ini bertujuan untuk menganalisis sentimen dan sebaran kata dari opini pengguna Twitter terkait dengan isu tersebut. Analisis sentimen dilakukan menggunakan pendekatan *machine learning* dengan metode klasifikasi *Naïve Bayes*. Penelitian ini menggunakan 901 data *tweet* yang sudah diberi label. Terdapat ketidakseimbangan antar kelas sentimen pada *dataset* yang digunakan dimana data sentimen positif memiliki jumlah yang jauh lebih sedikit dibanding dengan sentimen negatif dan netral. Sehingga dilakukan teknik *oversampling* menggunakan SMOTE pada data latih untuk membantu algoritma dalam membentuk model klasifikasi. Hasil dari model yang dibangun menggunakan algoritma *Naïve Bayes* mendapatkan nilai *accuracy* sebesar 71%, *precision* 62%, *recall* 74%, dan *f1-score* sebesar 62%. Ketidakseimbangan kelas sentimen membuat algoritma *Naïve Bayes* memiliki performa yang rendah.

Kata kunci: *Naïve Bayes*, SMOTE, Kebocoran Data SIM, Analisis Sentimen, Twitter.

ABSTRACT

The increasingly massive development of digital technology in Indonesia has made people unable to resist the impact of digitalization which can bring great threats such as threats to personal data security. In September 2022, there was an incident of sim card registration data leak containing users' personal data that was spread through dark sites. As one of the most popular social media in Indonesia, Twitter is used by the public to voice their opinions regarding the issue of sim registration data leakage. This study aims to analyze the sentiment and word distribution of Twitter users' opinions related to the issue. Sentiment analysis is conducted using a machine learning approach with the Naïve Bayes classification method. This research uses 901 tweet data that has been labeled. There is an imbalance between sentiment classes in the dataset used where positive sentiment data has a much smaller amount than negative and neutral sentiment. So an oversampling technique using SMOTE is carried out on the training data to help the algorithm in forming a classification model. The results of the model built using the Naïve Bayes algorithm get an accuracy value of 71%, precision 62%, recall 74%, and f1-score of 62%. The imbalance of sentiment classes makes the Naïve Bayes algorithm have low performance.

Keyword: *Naïve Bayes, SMOTE, SIM data leak, sentiment analysis, Twitter.*