

**ANALISIS PENGARUH PENERAPAN *STOPWORD REMOVAL*
DAN *GRID SEARCH (TUNING HYPERPARAMETER)* PADA
PERFORMA KLASIFIKASI SENTIMEN *TWEET* BAHASA
INDONESIA**

SKRIPSI

Program Studi Informatika

Jurusan Informatika

Oleh:

SHERREN JESSICA ANGELINA

NIM D1041161003



**JURUSAN INFORMATIKA
FAKULTAS TEKNIK
UNIVERSITAS TANJUNGPURA
PONTIANAK
2023**

HALAMAN PERNYATAAN

Yang bertanda tangan dibawah ini:

Nama : Sherren Jessica Angelina

NIM : D1041161003

menyatakan bahwa dalam skripsi yang berjudul “Analisis Pengaruh Penerapan *Stopword Removal Dan Grid Search (Tuning Hyperparameter)* Pada Performa Klasifikasi Sentimen *Tweet Bahasa Indonesia*” tidak terdapat karya yang pernah diajukan untuk memperoleh gelar sarjana di suatu perguruan tinggi manapun. Sepanjang pengetahuan Saya, 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.

Demikian pernyataan ini dibuat dengan sebenar-benarnya. Saya sanggup menerima konsekuensi akademis dan hukum di kemudian hari apabila pernyataan yang dibuat ini tidak benar.

Pontianak, 7 Juni 2023

Sherren Jessica Angelina

NIM D1041161003



KEMENTERIAN PENDIDIKAN, KEBUDAYAAN,
RISET, DAN TEKNOLOGI
UNIVERSITAS TANJUNGPURA
FAKULTAS TEKNIK

Jalan Prof. Dr. H. Hadari Nawawi Pontianak 78124
Telepon (0561) 740186 Email : ft@untan.ac.id Website : http://teknik.untan.ac.id

HALAMAN PENGESAHAN

**ANALISIS PENGARUH PENERAPAN *STOPWORD REMOVAL* DAN
GRID SEARCH (TUNING HYPERPARAMETER) PADA PERFORMA
KLASIFIKASI SENTIMEN *TWEET* BAHASA INDONESIA**

Program Studi Sarjana Informatika
Jurusan Informatika

Oleh:

Sherren Jessica Angelina
NIM. D1041161003

Telah dipertahankan di depan Pengaji Skripsi pada tanggal 7 Juni 2023 dalam sidang dan diterima sebagai salah satu persyaratan untuk memperoleh gelar sarjana.

Susunan Pengaji Skripsi

Ketua,

Dr. Arif Bijaksana Putra Negara, S.T., M.T.
NIP 197208081998021002

Pengaji Utama,

Prof. Dr. Henry Sujaini, S.T., M.T.
NIP 196806291997021001

Sekretaris,

Hafiz Muhamadi, S.T., M.Kom.
NIP. 0023079006

Pengaji Pendamping,

Rina Septiriana, S.T, M.Cs.
NIP 198709232020122001



Pontianak, 7 Juni 2023
Dekan,

Dr.-Ing. Ir. Slamet Widodo, M.T.,IPM.
NIP.196712231992031002

HALAMAN PERSEMPAHAN

Skripsi ini saya persembahkan kepada kedua orang tua, saudara, keluarga, teman, sahabat, dan semua pihak yang telah senantiasa memberikan semangat, masukkan, serta arahan dalam upaya mendukung penyelesaian skripsi ini.

KATA PENGANTAR

Puji syukur saya panjatkan kepada Tuhan Yang Maha Esa atas berkat dan rahmat-Nya yang besar, sehingga saya dapat menyelesaikan penggerjaan skripsi yang berjudul "*Analisis Pengaruh Penerapan Stopword Removal Dan Grid Search (Tuning Hyperparameter) Dalam Performa Klasifikasi Sentimen Tweet Bahasa Indonesia*". Tujuan dari penyusunan skripsi ini dibuat untuk memenuhi syarat dalam menyelesaikan Program Studi Sarjana Informatika Universitas Tanjungpura.

Penyelesaian penyusunan skripsi ini tidak lepas dari bantuan, arahan dan bimbingan banyak pihak. Oleh sebab itu penulis ingin menyampaikan ucapan terima kasih kepada Bapak Dr. Arif Bijaksana Putra Negara, ST., MT., dan Bapak Hafiz Muhardi, S.T., M.Kom. selaku pembimbing skripsi, kepada Bapak Prof. Dr. Herry Sujaini, S.T., M.T. dan Ibu Rina Septiriana, S.T, M.Cs selaku penguji skripsi yang telah memberikan banyak arahan, masukan, serta motivasi dalam membimbing penulis untuk dapat menyelesaikan skripsi ini dengan baik. Segenap dosen jurusan informatika atas segala ilmu dan bimbingannya. Kedua orang tua, keluarga, saudara/saudari saya yang sangat saya sayangi serta teman, sahabat, dan semua pihak yang tidak dapat disebutkan satu persatu.

Meski demikian, penulis merasa masih banyak kesalahan dalam penyusunan skripsi ini. Oleh sebab ini penulis sangat terbuka dalam menerima kritik dan saran yang membangun untuk dijadikan sebagai bahan evaluasi. Akhir kata, penulis mengharapkan semoga tujuan dari pembuatan skripsi ini dapat tercapai sesuai dengan yang diharapkan.

Pontianak, 7 Juni 2023

Penulis,

Sherren Jessica Angelina

ABSTRAK

Melalui *tweet* pada Twitter didapatkanlah berbagai informasi, salah satunya sentimen / pendapat yang bisa dijadikan acuan timbal balik respon masyarakat terkait suatu hal. Sentimen atau pendapat dapat berupa sentimen bersifat positif, netral dan negatif, yang bisa didapatkan melalui analisis sentimen atau *opinion mining*, yaitu sebuah metode penganalisa teks berbasis komputasi. Tujuan penelitian ini adalah menghasilkan model klasifikasi yang memiliki performa terbaik dalam mengklasifikasikan sentimen pada *tweet* Bahasa Indonesia dan mengetahui pengaruh penerapan *Stopword Removal*, SMOTE dan *Grid Search (Tuning Hyperparameter)* dalam membangun model klasifikasi sentimen analisis. Algoritma yang digunakan pada penelitian ini adalah *Logistic Regression* dan *Random Forest*. Berdasarkan hasil evaluasi yang dilakukan, diketahui bahwa pengimplementasian *Tuning Hyperparameter (Grid Search)* dan SMOTE pada algoritma *Logistic Regression* dan *Random Forest*, menghasilkan nilai *f1-score* dan peningkatan nilai *f1-score* setiap skenario terhadap skenario *default* tertinggi, yaitu sebesar 72.70% dan +1.20% untuk *Logistic Regression* dan sebesar 75.03% dan +6.77 untuk *Random Forest*. Oleh karena itu model klasifikasi terbaik pada penelitian ini adalah pada pengimplementasian algoritma *Random Forest* disertai *Tuning Hyperparameter (Grid Search)* dan SMOTE dengan nilai *f1-score* sebesar 75.03%. Kemudian untuk nilai *f1-score* dan peningkatan nilai *f1-score* skenario terhadap skenario *default* terendah, pada algoritma *Logistic Regression* adalah dengan pengimplementasian *Stopword Removal* dan SMOTE yaitu sebesar 68.60% dan -2.90% dan pada algoritma *Random Forest* adalah dengan pengimplementasian *Stopword Removal* yaitu sebesar 68.09% dan -0.17%. Oleh karena itu model klasifikasi terburuk pada penelitian ini adalah pada pengimplementasian algoritma *Random Forest* disertai *Stopword Removal* dengan nilai *f1-score* sebesar 68.09%. Penerapan *Stopword Removal* pada kedua algoritma, memberikan pengaruh kurang baik berupa penurunan nilai *f1-score* yang dihasilkan. Penurunan nilai *f1-score* ini dikarena *Stopword Removal* dapat mengurangi informasi dan mengubah makna tweet yang diolah sehingga tweet tersebut kehilangan sentimennya. Selain itu penerapan *stoplist NLTK* yang digunakan untuk melakukan *Stopword Removal* pada penelitian ini lebih bekerja optimal pada pengklasifikasian dokumen dibandingkan sentimen. Untuk penerapan *Tuning Hyperparameter (Grid Search)* dan SMOTE pada kedua algoritma, memberikan pengaruh yang baik berupa peningkatan nilai *f1-score* yang dihasilkan. Peningkatan nilai *f1-score* terjadi dikarenakan telah dioptimalkannya *hyperparameter* yang digunakan dan diseimbangkannya jumlah data antar kelas dalam *dataset*.

Kata kunci: Twitter, Sentimen Analisis, *Stopword Removal*, *Tuning Hyperparameter*, *Grid Search*, *Smote*, *Logistic Regression*, *Random Forest*.

DAFTAR ISI

HALAMAN PERNYATAAN.....	ii
HALAMAN PENGESAHAN.....	iii
HALAMAN PERSEMBAHAN	iv
KATA PENGANTAR.....	v
ABSTRAK	vi
DAFTAR ISI.....	vii
DAFTAR TABEL	xi
DAFTAR GAMBAR.....	xii
DAFTAR KODE PROGRAM	xiv
DAFTAR LAMPIRAN	xvi
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	4
1.3 Tujuan Penelitian.....	5
1.4 Batasan Masalah.....	6
1.5 Sistematika Penulisan Skripsi	7
BAB II TINJAUAN PUSTAKA.....	8
2.1 Kajian Terkait.....	8
2.2 Sentimen / Pendapat	14
2.3 Twitter	14
2.4 Text Mining	15
2.5 Sentimen Analisis (<i>Opinion Mining</i>)	15
2.6 Text Pre-processing	16
2.6.1 Cleaning	16

2.6.2	Case Folding	16
2.6.3	Normalization.....	16
2.6.4	Tokenizing	17
2.6.5	Stopword Removal.....	17
2.6.6	Stemming	18
2.6.7	Rejoin	18
2.7	Term Frequency – Inverse Document Frequency (TF-IDF)	18
2.8	Synthetic Minority Oversampling Technique (SMOTE)	20
2.9	Tuning Hyperparameter.....	20
2.9.1	Grid Search	21
2.10	Logistic Regression	21
2.11	Random Forest	22
2.12	K- Fold Cross Validation	22
2.13	Confusion Matrix	23
BAB III METODOLOGI PENELITIAN	25
3.1	Alat Penelitian	25
3.1.1	Perangkat Keras	25
3.1.2	Perangkat Lunak	25
3.2	Data Penelitian	25
3.3	Langkah Penelitian	26
3.3.1	Pengumpulan Data (<i>Data Gathering</i>)	28
3.3.2	Pembersihan Data (Text Pre-processing).....	28
3.3.3	Pembobotan (Vektorisasi Data)	28
3.3.4	Pembagian Dataset (Spliting Dataset).....	29
3.3.5	Penyeimbangan Dataset (Balancing Dataset)	29
3.3.6	Pemodelan (<i>Modelling</i>).....	29

3.3.7	Skenario Pengujian	30
3.3.8	Validasi (<i>Validation</i>).....	36
3.3.9	Evaluasi (<i>Evaluation</i>).....	36
BAB IV IMPLEMENTASI DAN HASIL.....	37	
4.1	Implementasi	37
4.1.1	Pengumpulan Data (<i>Data Gathering</i>)	37
4.1.2	Pembersihan Data (<i>Text Pre-processing</i>).....	38
4.1.3	Distribusi Frekuensi	44
4.1.4	Pembobotan (Vektorisasi Data)	44
4.1.5	Penyeimbangan Dataset (<i>Balancing Dataset</i>)	45
4.1.6	Pembagian Dataset (Splitting Dataset).....	45
4.1.7	Pemodelan (<i>Modelling</i>).....	46
4.1.8	Validasi (<i>Validation</i>).....	56
4.1.9	Evaluasi (<i>Evaluation</i>).....	56
4.1.10	Workflow Sistem	57
4.2	Hasil.....	58
4.2.1	Pengumpulan Data (<i>Data Gathering</i>)	58
4.2.2	Pembersihan Data (<i>Text Pre-processing</i>).....	60
4.2.3	Distribusi Frekuensi	62
4.2.4	Pembobotan (Vektorisasi Data)	66
4.2.5	Penyeimbangan Dataset (<i>Balancing Dataset</i>)	67
4.2.6	Pembagian Dataset (Splitting Dataset).....	69
4.2.7	Tuning Hyperparameter	70
4.2.8	Validasi (<i>Validation</i>).....	71
4.2.9	Evaluasi (<i>Evaluation</i>).....	80

4.2.10	Analisis Perbandingan Evaluasi (Evaluation) Logistic Regression dan Random Forest.....	111
4.2.11	Analisis Validasi (<i>Validation</i>) dan Evaluasi (<i>Evaluation</i>)	123
BAB V PENUTUP	127
5.1	Kesimpulan.....	127
5.2	Saran.....	128
DAFTAR PUSTAKA	129
LAMPIRAN	132

DAFTAR TABEL

Tabel 2. 1 Rangkuman Kajian Terkait	10
Tabel 2. 2 Confusion Matrix	23
Tabel 3. 1 Sampel Dataset Indonesia Untuk Analisis Sentimen	26
Tabel 3. 2 Skenario Pengujian Logistic Regression.....	31
Tabel 3. 3 Skenario Pengujian Random Forest	31
Tabel 4. 1 Jumlah Dataset	59
Tabel 4. 2 Hasil Text Pre-processing	61
Tabel 4. 3 Distribusi Frekuensi Text Pre-Processing	62
Tabel 4. 4 Distribusi Frekuensi Text Pre-Processing + Stopword Removal.....	63
Tabel 4. 5 Kalimat Sintetik SMOTE	69
Tabel 4. 6 Tuning Hyperparameter Logistic Regression.....	70
Tabel 4. 7 Tuning Hyperparameter Random Forest	71
Tabel 4. 8 Hasil Validasi Setiap Skenario Logistic Regression	72
Tabel 4. 9 Hasil Validasi Setiap Skenario Random Forest	76
Tabel 4. 10 Hasil Evaluasi, Perbandingan Dan Perubahan Akurasi Logistic Regression	86
Tabel 4. 11 Hasil Evaluasi, Perbandingan Dan Perubahan Akurasi Random Forest	101
Tabel 4. 12 Perbandingan Evaluasi dan Jarak Perubahan Nilai F1-Score Logistic Regression Dan Random Forest.....	111
Tabel 4. 13 Perubahan Makna Kalimat Dikarenakan Stopword Removal.....	122
Tabel 4. 14 Perbandingan Validasi Dan Evaluasi Logistic Regression	123
Tabel 4. 15 Perbandingan Validasi Dan Evaluasi Random Forest.....	125

DAFTAR GAMBAR

Gambar 2. 1 Ilustrasi TF-IDF.....	20
Gambar 2. 2 10-Fold Cross Validation.....	23
Gambar 3. 1 Metodologi Penelitian.....	27
Gambar 3. 2 Flowchart Skenario Pengujian Pertama.....	32
Gambar 3. 3 Flowchart Skenario Pengujian Kedua	32
Gambar 3. 4 Flowchart Skenario Pengujian Ketiga	33
Gambar 3. 5 Flowchart Skenario Pengujian Keempat	33
Gambar 3. 6 Flowchart Skenario Pengujian Kelima	34
Gambar 3. 7 Flowchart Skenario Pengujian Keenam.....	34
Gambar 3. 8 Flowchart Skenario Pengujian Ketujuh	35
Gambar 3. 9 Flowchart Skenario Pengujian Kedelapan.....	35
Gambar 4. 1 Inisialisasi Data.....	38
Gambar 4. 2 Cleaning	39
Gambar 4. 3 Case Folding	39
Gambar 4. 4 Normalisasi	40
Gambar 4. 5 Kamus Normalisasi.....	40
Gambar 4. 6 Tokenisasi	41
Gambar 4. 7 Stopword Removal	42
Gambar 4. 8 Stoplist NLTK	42
Gambar 4. 9 Stemming	43
Gambar 4. 10 Rejoin.....	43
Gambar 4. 11 Jumlah Data Sebelum SMOTE.....	45
Gambar 4. 12 Jumlah Data Setelah SMOTE.....	45
Gambar 4. 13 Workflow Sistem	57
Gambar 4. 14 Dataset Indonesia Untuk Analisis Sentimen.....	59
Gambar 4. 15 Jumlah Dataset Indonesia Untuk Analisis Sentimen	59
Gambar 4. 16 Hasil Text Pre-processing	60
Gambar 4. 17 Hasil Text Pre-processing + Stopword.....	60
Gambar 4. 18 Distribusi Frekuensi Text Pre-Processing	65
Gambar 4. 19 Distribusi Frekuensi Text Pre-Processing + Stopword Removal .	65

Gambar 4. 20 TF-IDF Text Pre-Processing	66
Gambar 4. 21 Vocabulary Content Text Pre-Processing.....	66
Gambar 4. 22 TF-IDF Text Pre-Processing + Stopword Removal	67
Gambar 4. 23 Vocabulary Content Text Pre-Processing + Stopword Removal .	67
Gambar 4. 24 Data Sentimen Sebelum SMOTE	68
Gambar 4. 25 Data Sentimen Sesudah SMOTE.....	68
Gambar 4. 26 Pembagian Dataset	70
Gambar 4. 27 Pembagian Dataset SMOTE.....	70
Gambar 4. 28 Hasil Validasi Setiap Skenario Logistic Regression	72
Gambar 4. 29 Hasil Validasi Setiap Skenario Random Forest.....	77
Gambar 4. 30 Logistic Regression Confusion Matrix Skenario 1.....	80
Gambar 4. 31 Logistic Regression Confusion Matrix Skenario 2.....	81
Gambar 4. 32 Logistic Regression Confusion Matrix Skenario 3.....	82
Gambar 4. 33 Logistic Regression Confusion Matrix Skenario 4.....	82
Gambar 4. 34 Logistic Regression Confusion Matrix Skenario 5.....	83
Gambar 4. 35 Logistic Regression Confusion Matrix Skenario 6.....	84
Gambar 4. 36 Logistic Regression Confusion Matrix Skenario 7.....	84
Gambar 4. 37 Logistic Regression Confusion Matrix Skenario 8.....	85
Gambar 4. 38 Random Forest Confusion Matrix Skenario 1	96
Gambar 4. 39 Random Forest Confusion Matrix Skenario 2	96
Gambar 4. 40 Random Forest Confusion Matrix Skenario 3	97
Gambar 4. 41 Random Forest Confusion Matrix Skenario 4	98
Gambar 4. 42 Random Forest Confusion Matrix Skenario 5	98
Gambar 4. 43 Random Forest Confusion Matrix Skenario 6	99
Gambar 4. 44 Random Forest Confusion Matrix Skenario 7	99
Gambar 4. 45 Random Forest Confusion Matrix Skenario 8	100
Gambar 4. 46 Perbandingan Evaluasi Logistic Regression Dan Random Forest	112
Gambar 4. 47 Perbandingan Validasi dan Evaluasi Logistic Regression	123
Gambar 4. 48 Perbandingan Validasi dan Evaluasi Random Forest.....	125

DAFTAR KODE PROGRAM

Kode Program 4. 1 Inisialisasi Data.....	38
Kode Program 4. 2 Cleaning.....	39
Kode Program 4. 3 Case Folding	39
Kode Program 4. 4 Kamus Normalisasi Bahasa Indonesia	40
Kode Program 4. 5 Normalisasi	41
Kode Program 4. 6 Tokenisasi	41
Kode Program 4. 7 Stoplist NLTK	42
Kode Program 4. 8 Stopword Removal	42
Kode Program 4. 9 Stemming	43
Kode Program 4. 10 Rejoin.....	44
Kode Program 4. 11 Distribusi Frekuensi.....	44
Kode Program 4. 12 TF-IDF.....	44
Kode Program 4. 13 SMOTE.....	45
Kode Program 4. 14 Pembagian Dataset.....	46
Kode Program 4. 15 Pembagian Dataset SMOTE	46
Kode Program 4. 16 Grid Search Logistic Regression Dataset	47
Kode Program 4. 17 Grid Search Logistic Regression Dataset Stopword.....	47
Kode Program 4. 18 Grid Search Logistic Regression Dataset + SMOTE	47
Kode Program 4. 19 Grid Search Logistic Regression Dataset Stopword + SMOTE	47
Kode Program 4. 20 Grid Search Random Forest Dataset.....	48
Kode Program 4. 21 Grid Search Random Forest Dataset Stopword	48
Kode Program 4. 22 Grid Search Random Forest Dataset + SMOTE.....	49
Kode Program 4. 23 Grid Search Random Forest Dataset Stopword + SMOTE	49
Kode Program 4. 24 Logistic Regression Skenario 1	49
Kode Program 4. 25 Logistic Regression Skenario 2	50
Kode Program 4. 26 Logistic Regression Skenario 3	50
Kode Program 4. 27 Logistic Regression Skenario 4	50
Kode Program 4. 28 Logistic Regression Skenario 5	51
Kode Program 4. 29 Logistic Regression Skenario 6	51

Kode Program 4. 30	Logistic Regression Skenario 7	52
Kode Program 4. 31	Logistic Regression Skenario 8	52
Kode Program 4. 32	Random Forest Skenario 1	53
Kode Program 4. 33	Random Forest Skenario 2	53
Kode Program 4. 34	Random Forest Skenario 3	53
Kode Program 4. 35	Random Forest Skenario 4	54
Kode Program 4. 36	Random Forest Skenario 5	54
Kode Program 4. 37	Random Forest Skenario 6	55
Kode Program 4. 38	Random Forest Skenario 7	55
Kode Program 4. 39	Random Forest Skenario 8	56
Kode Program 4. 40	Validasi Setiap Skenario	56
Kode Program 4. 41	Evaluasi Tiap Skenario	56

DAFTAR LAMPIRAN

Lampiran 1 Kamus Normalisasi KBBI_Final	132
Lampiran 2 Stoplist NLTK	135
Lampiran 3 Dataset Indonesia Untuk Analisis Sentimen.....	136
Lampiran 4 Hasil Preprocessing (Cleaning)	137
Lampiran 5 Hasil Preprocessing (Case Folding).....	138
Lampiran 6 Hasil Preprocessing (Normalisasi)	139
Lampiran 7 Hasil Preprocessing (Tokenisasi)	140
Lampiran 8 Hasil Preprocessing (Stemming)	141
Lampiran 9 Hasil Preprocessing (Rejoin)	142
Lampiran 10 Hasil Preprocessing (Stopword Removal).....	143
Lampiran 11 Hasil Preprocessing (Stemming + Stopword Removal)	144
Lampiran 12 Hasil Preprocessing (Rejoin + Stopword Removal)	145
Lampiran 13 Classification Report Logistic Regression.....	146
Lampiran 14 Classification Report Random Forest.....	149

BAB I

PENDAHULUAN

1.1 Latar Belakang

Data dan informasi merupakan dua hal yang saling terkait dan terikat. Dikatakan demikian karena untuk mendapatkan suatu informasi diperlukanlah sekumpulan data yang faktual dan aktual untuk diolah. Data merupakan bahan mentah dari informasi, yang dirumuskan sebagai sekelompok lambang-lambang tidak acak yang menunjukkan jumlah atau tindakan atau hal-hal lain (Darmoyo, 2020). Sedangkan informasi dapat dijabarkan sebagai data yang telah diproses sedemikian rupa sehingga meningkatkan pengetahuan seseorang yang menggunakan data tersebut (Kadir, 2008).

Di era digital seperti sekarang, dimana hampir semua kalangan sudah memiliki gawai / *gadget*, penyebaran informasi dapat dikatakan sangat cepat. Hal ini dikarenakan banyaknya wadah / *platform* yang dapat diakses untuk menguraikan pokok pikiran dan pendapat terkait kebijakan maupun isu-isu yang berkembang. *Website*, media sosial dan *e-commers* merupakan contoh-contoh dari wadah / *platform digital*. Dimana dari ketiga wadah / *platform digital* tersebut, media sosial adalah yang terbanyak penggunanya, sehingga penyebaran informasi melalui media sosial terasa lebih cepat dan efisien.

Di Indonesia, Twitter merupakan salah satu media sosial populer dan favorit yang digunakan oleh berbagai kalangan lapisan masyarakat bahkan pemerintah. Dengan berbentuk pesan singkat yang berisikan cuitan / kicauan yang lebih dikenal dengan “*tweet*” sebagai sarana penggunaannya, masyarakat lebih dapat bereksresi dalam mengemukakan isi pikiran dan pendapat, ditambah twitter sendiri tidak mengharuskan para penggunanya saling berteman dahulu untuk dapat berinteraksi / berdiskusi, hanya cukup dengan adanya bahasan yang sama maka interaksi akan terjadi.

Diskusi / interaksi ini akan menghasilkan data, yaitu berupa cuitan / kicauan / *tweet*, yang jika diolah, data tersebut akan menghasilkan informasi, salah satunya berupa sentimen / pendapat masyarakat yang bisa dijadikan acuan timbal balik

respon masyarakat. Sentimen / pendapat masyarakat dapat diklasifikasikan menjadi 3, yaitu sentimen / pendapat bersifat positif, netral dan negatif. Data sentimen / pendapat ini umumnya dapat dikumpulkan secara manual, salah satunya melalui penyebaran kuesioner. Namun dalam penerapannya akan cukup menguras waktu dan tenaga, maka penarikan data tweet terasa lebih efisien untuk digunakan. Untuk melakukan sentimen analisis diperlukanlah penerapan metode penganalisa *text* berbasis komputasi yang dikenal dengan analisis sentimen atau *opinion mining*.

Secara umum terdapat beberapa tahapan yang akan dilakukan dalam analisis sentimen yaitu tahapan Pembersihan Data (*Text Pre-processing*), Pembobotan (Vektorisasi Data), Pemodelan, dan Evaluasi. Berfokus pada tahapan Pembersihan Data (*Text Pre-processing*) dan Pemodelan. Dimana keduanya mengambil peran penting dalam analisis sentimen, karena kedua tahapan tersebut merupakan tahapan yang akan bertugas dalam menghasilkan *dataset* dan pemodelan mesin klasifikasi yang nantinya akan berpengaruh pada performa nilai *f1-score* klasifikasi sentimen yang dihasilkan.

Pada dasarnya data twitter “*tweet*” yang akan digunakan sebagai *dataset*, pastinya terdapat derau atau *noise*, seperti banyak penggunaan kata yang tidak baku, kata singkatan bahkan sampai penggunaan kata dalam bahasa gaul. Hal inilah yang menyebabkan perlunya dilakukan tahapan *Text Pre-processing*. *Stopword Removal* merupakan salah satu metode yang bisa digunakan pada tahapan *Text Pre-processing*. Dimana *Stopword Removal* akan melakukan penghapusan kata-kata yang cukup umum dan sering muncul namun tidak mempunyai pengaruh yang signifikan terhadap makna suatu teks atau kalimat. Sehingga diharapkan dapat menghasilkan *dataset* yang lebih baik, karena kualitas *dataset* akan berpengaruh pada performa nilai *f1-score* klasifikasi sentimen.

Selanjutnya dalam tahapan pemodelan, sebelum melakukan Pemodelan Mesin Klasifikasi alangkah lebih baik jika dilakukan penerapan *Tuning Hyperparameter* terlebih dahulu. Sebenarnya kedua kegiatan ini saling berhubungan, dimana dengan adanya penerapan *Tuning Hyperparameter* akan berpengaruh pada optimalisasi pemodelan mesin klasifikasi. *Grid Search* merupakan salah satu algoritma yang dapat digunakan untuk memilih kombinasi *hyperparameter*, dengan cara mengkombinasikan *hyperparamater*-

hyperparameter yang dimasukkan dan mencari kombinasi *hyperparameter* dengan nilai *f1-score* yang paling tinggi untuk diimplementasikan pada model klasifikasi yang dibuat. Untuk pemodelan mesin klasifikasi, dapat dilakukan dengan pendekatan algoritma klasifikasi, seperti Algoritma *Logistic Regression* dan *Random Forest*. Hal ini dikarenakan, pada penerapannya Algoritma *Logistic Regression* dan *Random Forest* memiliki kemampuan untuk mengolah data yang besar dan menghasilkan nilai *f1-score* yang baik.

Pada umumnya performa metode dan algoritma yang diimplementasikan dalam melakukan klasifikasi sentimen akan kurang optimal jika *dataset* yang digunakan mengalami ketidakseimbangan data atau *imbalanced data*. Untuk mengatasinya dapat dengan dilakukannya pengimplementasian SMOTE. SMOTE atau *Synthetic Minority Oversampling Technique* merupakan teknik untuk mengatasi masalah *imbalanced data*, dengan cara kerja yaitu membuat data sintetik baru dari kelas minoritas pada *dataset* sehingga terjadi keterseimbangan data antar kelas (Cahyaningtyas et al., 2021). Dengan adanya SMOTE maka *dataset* yang digunakan tidak akan bias terhadap kelas mayoritas, sehingga diharapkan dapat mengoptimalkan kinerja dari metode dan algoritma klasifikasi.

Adapun kajian terkait penggunaan *Stopword Removal* pada tahapan *Text Pre-processing* telah dilakukan oleh beberapa pihak, seperti pada penelitian yang dilakukan oleh (Santosa et al., 2022) dimana dilakukan uji coba pengaruh penggunaan *Stopword Removal* dan *Stemming* dengan Algoritma LSM, didapatkan hasil dimana nilai keakurasaan yang paling baik adalah tanpa digunakannya *Stopword Removal* dan *Stemming*, dengan nilai akurasi sebesar 0.82, *loss* 0.4, *precision* 0.83, *recall* 0.81, dan *f1-score* 0.82. Sedangkan pada penelitian yang dilakukan (Hidayatullah, 2016) dimana dilakukan uji coba pengaruh penggunaan *Stopword Removal* dengan menggunakan Algoritma KNN, SVM dan *Naïve Bayes*, didapatkanlah hasil bahwa nilai akurasi data jika digunakannya *Stopword Removal* dengan penerapan Algoritma SVM sebagai yang terbaik dengan nilai akurasi sebesar 88,59%.

Selain itu, untuk penggunaan *Grid Search* sebagai *Tuning Hyperparameter* juga sudah pernah dilakukan, seperti penelitian yang dilakukan (Negara et al., 2021) dimana dilakukan perbandingan Algoritma KNN dan *Logistic Regression* untuk

pengklasifikasian emosi *tweet* Bahasa Indonesia dengan penerapan *TF-IDF* pada tahapan vektorisasi dan *Grid Search* sebagai *Tuning Hyperparameter*, dimana di hasilkan informasi bahwa nilai akurasi dan *f1-score* terbaik yang didapatkan adalah dengan pengimplementasian Algoritma *Logistic Regression* dengan penerapan *TF-IDF* dan *Grid Search* sebagai *Tuning Hyperparameter*, dengan nilai *accuracy* sebesar 65% dan *f1-score* sebesar 66%. Sementara untuk penelitian yang dilakukan (Andreyestha & Azizah, 2022) terkait penggunaan SMOTE dengan pengimplementasian algoritma *Naïve Bayes* dan *Random Forest* dalam mengklasifikasikan analisis sentimen kicauan twitter tokopedia, dihasilkan informasi bahwa penggunaan SMOTE dapat meningkatkan nilai akurasi sebesar 3.4% pada Algoritma *Naïve Bayes* dan 1.55% pada Algoritma *Random Forest*. Nilai akurasi dan *f1-score* terbaik didapatkan oleh Algoritma *Random Rorest* dengan pengimplementasian SMOTE yaitu sebesar 88.44% dan 88.30%.

Berdasarkan uraian yang telah dipaparkan, maka penelitian ini akan merujuk pada analisis bagaimana pengaruh penerapan *Stopword Removal* pada tahapan *Text Pre-processing* dan *Grid Search* sebagai *Tuning Hyperparameter*, terhadap nilai *f1-score* yang dihasilkan. Penelitian ini juga menerapkan metode SMOTE dalam upaya menyeimbangkan *dataset* yang digunakan. Selain itu penelitian ini menggunakan 2 algoritma klasifikasi untuk pemodelan mesin analisis yaitu *Logistic Regression* dan *Random Forest* sehingga akan dilakukan perbandingan nilai *f1-score* yang dihasilkan keduanya untuk melihat algoritma mana yang pengimplementasiannya lebih optimal dalam melakukan analisis sentimen *tweet* berbahasa Indonesia.

1.2 Rumusan Masalah

Berdasarkan pada permasalahan yang telah dipaparkan dalam latar belakang di atas, maka diperoleh rumusan masalah sebagai berikut, yaitu :

1. Bagaimana pengaruh penerapan *Stopword Removal* pada tahapan *Text Pre-processing* dalam performa klasifikasi sentimen analisis *tweet* Bahasa Indonesia dengan penerapan Algoritma *Logistic Regression* dan *Random Forest* ?

2. Bagaimana pengaruh penerapan *Grid Search* sebagai *Tuning Hyperparameter* dalam performa klasifikasi sentimen analisis *tweet* Bahasa Indonesia dengan penerapan Algoritma *Logistic Regression* dan *Random Forest* ?
3. Apakah penerapan metode SMOTE dalam menyeimbangkan *dataset* mampu mengoptimalkan performa klasifikasi sentimen analisis *tweet* Bahasa Indonesia dengan penerapan Algoritma *Logistic Regression* dan *Random Forest* ?
4. Algoritma manakah yang memiliki performa paling optimal antara *Logistic Regression* dan *Random Forest* dalam melakukan analisis sentimen *tweet* berbahasa Indonesia ?

1.3 Tujuan Penelitian

Berdasarkan rumusan masalah yang telah disampaikan, maka diperoleh tujuan yang akan dicapai adalah sebagai berikut :

1. Untuk mengetahui perbandingan nilai *f1-score* sentimen yang dihasilkan oleh Algoritma *Logistic Regression* dan *Random Forest* dengan atau tanpa diterapkannya *Stopword Removal* pada tahapan *Text Pre-processing* dalam mengklasifikasi sentimen *tweet* Bahasa Indonesia.
2. Untuk mengetahui perbandingan nilai *f1-score* sentimen yang dihasilkan oleh Algoritma *Logistic Regression* dan *Random Forest* dengan atau tanpa diterapkannya *Grid Search* sebagai *Tuning Hyperparameter* dalam mengklasifikasi sentimen *tweet* Bahasa Indonesia.
3. Untuk mengetahui performa metode SMOTE dalam mengoptimalkan nilai *f1-score* sentimen yang dihasilkan oleh Algoritma *Logistic Regression* dan *Random Forest* dalam mengklasifikasi sentimen *tweet* Bahasa Indonesia.
4. Untuk mengetahui algoritma yang paling optimal antara *Logistic Regression* dan *Random Forest* dalam melakukan analisis sentimen *tweet* Bahasa Indonesia.

1.4 Batasan Masalah

Adapun beberapa hal yang menjadi batasan permasalahan untuk menghindari adanya penyimpangan dan pelebaran inti permasalahan dalam penelitian ini adalah sebagai berikut :

1. Data yang digunakan sebagai *dataset* dalam penelitian ini berupa data berasal dari penelitian Ridi Ferdiana, Fahim Jatmiko, Desi Dwi Purwanti, Artmita Sekar Tri Ayu, Wiliam Fajar Dicka pada tahun 2019 yang berjudul “*Dataset Indonesia untuk Analisis Sentimen*” dan berbentukan teks berbahasa Indonesia yang diambil dari twitter.
2. Penelitian ini akan melakukan pengklasifikasian sentimen / pendapat menjadi 3 kelas yaitu sentimen / pendapat bersifat positif, netral dan negatif.
3. Penerapan *Stopword Removal* pada penelitian ini menggunakan pengimplementasian *library python* yaitu *library NLTK (Natural Language Toolkit)*.
4. Penerapan *Tuning Hyperparameter* pada penelitian ini menggunakan pengimplementasian metode *Grid Search* dengan pemanfaatan *library python* yaitu *library Scikit Learn*.
5. Penerapan *SMOTE* pada penelitian ini menggunakan pengimplementasian *library python* yaitu *library Imblearn*.
6. Algoritma klasifikasi yang digunakan adalah Algoritma *Logistic Regression* dan *Random Forest*.
7. *Hyperparameter* algoritma klasifikasi yang akan dilakukan pengimplementasian *Tuning Hyperparameter* pada penelitian ini berupa nilai *C* pada *Logistic Regression* dan nilai *Estimators*, *Max_depth*, *Max_features* dan *Criterion* pada *Random Forest*.
8. Pembagian *dataset* pada penelitian ini terdiri atas 2 bagian, yaitu data latih (*training*) dan uji (*testing*). Dengan rincian yaitu, 90% untuk data latih (*training*) dan 10% untuk data uji (*testing*).
9. Penelitian ini akan berfokus pada penganalisaan pengaruh penerapan *Stopword Removal* pada tahapan *Text Pre-processing* dan *Grid Search*

sebagai *Tuning Hyperparameter* dalam melakukan sentimen analisis *tweet* berbahasa Indonesia terhadap nilai *f1-score* yang dihasilkan.

1.5 Sistematika Penulisan Skripsi

Sistematika penulisan untuk tugas akhir ini disusun ke dalam lima bab yang terdiri dari:

1. Bab I Pendahuluan adalah bab yang memaparkan mengenai latar belakang penelitian, perumusan masalah, tujuan penelitian, pembatasan masalah, dan sistematika penulisan.
2. Bab II Tinjauan Pustaka adalah bab yang berisi mengenai landasan teori, prinsip-prinsip, serta kajian-kajian yang berkaitan dengan penelitian yang akan dilakukan.
3. Bab III Metodologi Penelitian adalah bab yang didalamnya berisikan tentang data penelitian, alat yang digunakan, metode penelitian serta diagram alir penelitian yang meliputi tahapan Pembersihan Data (*Text Pre-processing*), Pembobotan (Vektorisasi Data), Penyeimbangan Dataset (*Balancing Data*), Pemodelan (*Modelling*), Validasi (*Validation*) dan Evaluasi (*Evaluation*).
4. Bab IV Hasil dan Analisis adalah bab yang menjabarkan mengenai penelitian yang dilakukan sesuai dengan metodologi penelitian. Dimana kemudian hasil dari penelitian akan dianalisis dan dipaparkan menjadi suatu narasi yang mengacu pada penarikan suatu kesimpulan, dapat dilengkapi dengan tabel, gambar maupun grafik.
5. Bab V Penutup adalah bab yang berisi kesimpulan dari penelitian yang telah dilakukan, dilengkapi dengan saran serta rekomendasi untuk perbaikan, pengembangan, penyempurnaan dan pelengkapan penelitian yang telah dilakukan.