

**PENGARUH STEMMING DAN STOPWORDS TERHADAP  
AKURASI ANALISIS SENTIMEN MENGGUNAKAN  
METODE SUPPORT VECTOR MACHINE**

**SKRIPSI**



disusun oleh

**Aditya Wiha Pradana**

**15.11.9295**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2019**

**PENGARUH STEMMING DAN STOPWORDS TERHADAP  
AKURASI ANALISIS SENTIMEN MENGGUNAKAN  
METODE SUPPORT VECTOR MACHINE**

**SKRIPSI**

untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Aditya Wiha Pradana**

**15.11.9295**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2019**

**PERSETUJUAN**

**SKRIPSI**

**PENGARUH STEMMING DAN STOPWORDS TERHADAP  
AKURASI ANALISIS SENTIMEN MENGGUNAKAN  
METODE SUPPORT VECTOR MACHINE**

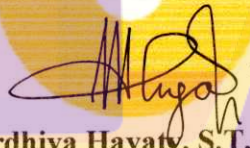
yang dipersiapkan dan disusun oleh

**Aditya Wiha Pradana**

**15.11.9295**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 23 Juli 2019

**Dosen Pembimbing,**



**Mardhiya Hayaty, S.T, M.Kom**

**NIK. 190302108**

## PENGESAHAN

### SKRIPSI

#### PENGARUH STEMMING DAN STOPWORDS TERHADAP AKURASI ANALISIS SENTIMEN MENGGUNAKAN METODE SUPPORT VECTOR MACHINE

yang dipersiapkan dan disusun oleh  
**Aditya Wiha Pradana**

15.11.9295

telah dipertahankan di depan Dewan Penguji  
pada tanggal 23 Agustus 2019

#### Susunan Dewan Penguji

**Nama Penguji**

**Tanda Tangan**

Sri Ngudi Wahyuni, S.T, M.Kom  
NIK. 190302060



Rumini, M.Kom  
NIK. 190302246



Mardhiya Hayaty, S.T, M.Kom  
NIK. 190302108



Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 5 September 2019

**DEKAN FAKULTAS ILMU KOMPUTER**



  
Krisnawati, S.Si, M.T  
NIK. 190302038



## 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 tanggung jawab saya pribadi.

Yogyakarta, 23 Agustus 2019



Aditya Wiha Pradana

NIM. 15.11.9295

## MOTTO

**“***Door Duisternis Tot Licht***”**

- R.A. Kartini

**“**それが悪魔であろうと聖人であろうと、私の名前は世界中で聞かれます!**”**

- Zorojuro

**“***The fool doth think he is wise, but the wise man knows himself to be a fool.***”**

- William Shakespeare

**“***I don't care that they stole my idea ... I care that they don't have any of their own.***”**

- Nikola Tesla

## PERSEMBAHAN

Alhamdulillah segala puji syukur atas kehadiran Allah SWT yang telah memberikan limpahan rahmat dan karunia-Nya sehingga penelitian ini dapat dilakukan dan diselesaikan dengan sebaik-baiknya. Saya juga ucapkan terimakasih untuk dukungan dan bantuan semua pihak yang membantu selesainya penelitian ini.

Ucapan terimakasih saya persembahkan kepada :

1. Kedua orang tua dan semua keluarga yang senantiasa memberikan dukungan dan do'a kepada saya.
2. Ibu Mardhiya Hayaty, S.T, M.Kom selaku dosen pembimbing yang telah banyak memberikan bimbingan, pelajaran, serta ilmu yang sangat bermanfaat.
3. Teman-teman kelas 15-S1IF-12 yang menjadi teman seperjuangan dari awal perkuliahan hingga saat ini.
4. Sahabat-sahabat Nongki, yang banyak memberikan dukungan, kritik dan saran yang membantu proses penulisan skripsi ini.
5. Sobat Rodi, yang telah berjuang bersama dalam proses penulisan skripsi dari awal sampai selesai.
6. Teman-teman kos ijo, yang sudah menemani dan bersenang-senang bersama penulis selama perkuliahan.

## KATA PENGANTAR

Assalamualaikum Wr.Wb.

Puji syukur kepada Allah SWT yang telah memberikan rahmat hidayah, karunia dan kesehatan, sehingga penulis dapat menyelesaikan laporan skripsi yang berjudul “Pengaruh Stemming dan Stopwords Terhadap Akurasi Analisis Sentimen Menggunakan Metode Support Vector Machine”. Keberhasilan dalam menyelesaikan pembuatan laporan skripsi ini adalah berkat bantuan dan dukungan dari berbagai pihak. Maka dari itu pada kesempatan kali ini penulis ingin mengucapkan terimakasih yang sebesar-besarnya kepada:

1. Bapak Prof. Dr. M. Suyanto, M.M. selaku ketua yayasan Universitas Amikom Yogyakarta
2. Ibu Krisnawati, S.Si., M.T. selaku Dekan Fakultas Ilmu Komputer dan Bapak Sudarmawan, S.T., M.T. selaku Kepala Program Studi S1 Informatika.
3. Ibu Mardhiya Hayaty, S.T, M.Kom Selaku dosen pembimbing yang telah memberikan pengarahan, bimbingan dan motivasi selama proses penyusunan skripsi hingga selesai.
4. Segenap Dosen dan Karyawan Universitas Amikom Yogyakarta yang telah memberikan ilmu pengetahuan dan pengalamannya.
5. Kepada kedua orang tua penulis yang telah memberikan dukungan terbaiknya selama kuliah.
6. Kepada Teman-teman angkatan 2015 khususnya kelas 11-S1-TI12 yang telah berjuang bersama.

Penulis menyadari sepenuhnya bahwa laporan skripsi ini masih sangat jauh dari kesempurnaan, itu semua tidak lepas dari keterbatasan pengetahuan dan kemampuan dari penulis sendiri. Untuk itu, penulis mengharapkan kritik dan saran

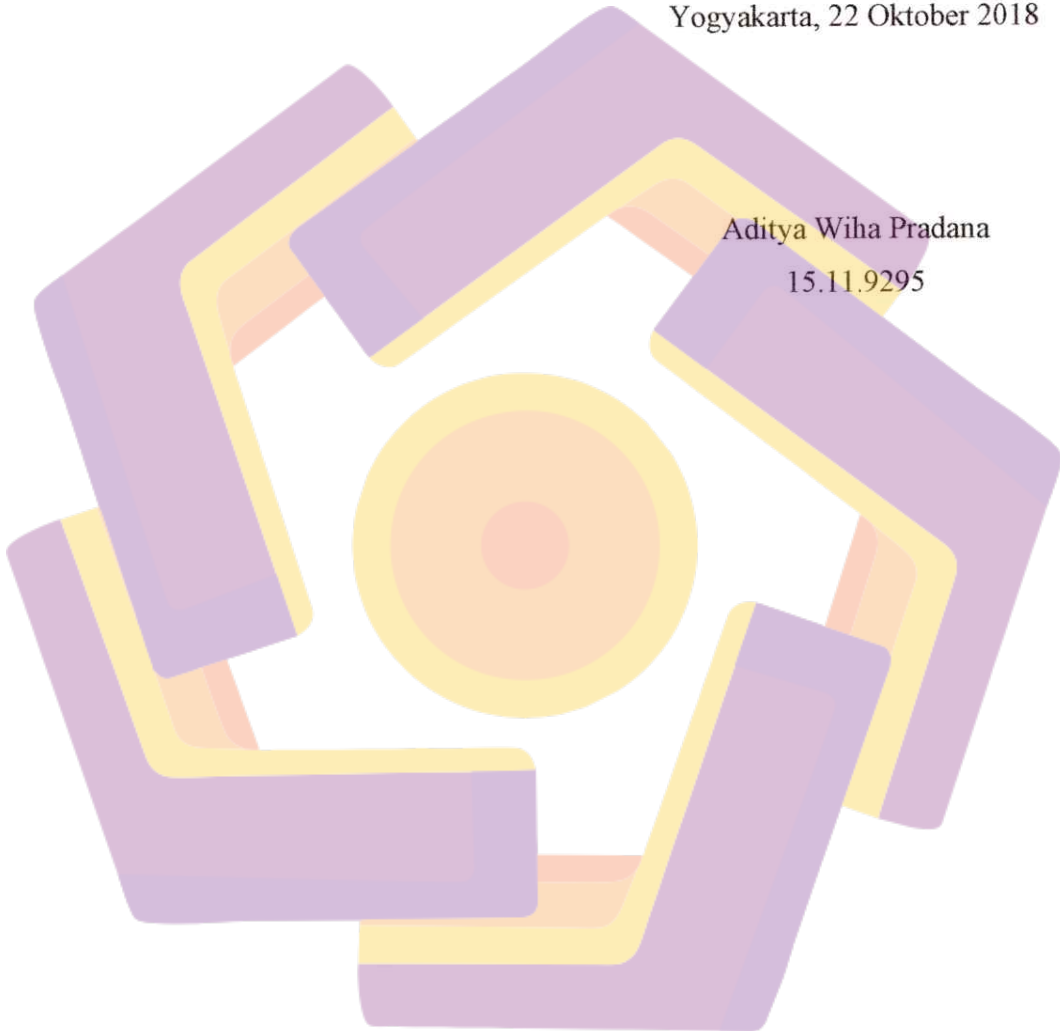


yang bersifat membangun guna mencapai kesempurnaan yang selalu penulis harapkan sehingga dapat bermanfaat bagi penulis, serta pihak-pihak yang membutuhkan.

Yogyakarta, 22 Oktober 2018

Aditya Wiha Pradana

15.11.9295



## DAFTAR ISI

JUDUL .....	ii
PERSETUJUAN .....	iii
PENGESAHAN .....	iv
PERNYATAAN.....	v
MOTTO .....	vi
PERSEMBAHAN.....	vii
KATA PENGANTAR .....	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR .....	xiv
INTISARI.....	xvi
ABSTRACT.....	xvii
BAB I PENDAHULUAN.....	1
1.1. Latar Belakang .....	1
1.2. Rumusan Masalah .....	2
1.3. Batasan Masalah.....	3
1.4. Tujuan Penelitian.....	3
1.5. Manfaat Penelitian.....	4
1.6. Hipotesis Penelitian.....	4
1.7. Metode Penelitian.....	4
1.8. Sistematika Penulisan.....	5
BAB II LANDASAN TEORI .....	7
2.1 Tinjauan Pustaka .....	7
2.2 Landasan Teori.....	10
2.2.1 Text Mining.....	10
2.2.2 Analisis Sentimen .....	11

2.2.1	Text Mining.....	10
2.2.2	Analisis Sentimen .....	11
2.2.3	Klasifikasi .....	11
2.2.4	Preprocessing Data.....	12
2.2.5	Term Frequency – Inverse Document Frequency .....	13
2.2.6	Support Vector Machine .....	15
2.2.7	Confusion Matrix .....	20
2.2.8	K-Fold Cross Validation .....	22
<b>BAB III METODE PENELITIAN.....</b>		<b>23</b>
3.1	Alat dan Bahan .....	23
3.1.1	Alat Penelitian.....	23
3.1.2	Bahan Penelitian.....	24
3.2	Alur Penelitian.....	24
3.2.1	Pengumpulan Data .....	27
3.2.2	Preprocessing Data.....	28
3.2.3	Pembagian Data .....	31
3.2.4	TF-IDF .....	31
3.2.5	Klasifikasi .....	31
3.2.6	Evaluasi .....	33
3.2.7	Validasi .....	33
3.2.8	Analisis Hasil Evaluasi dan Validasi .....	34
<b>BAB IV HASIL DAN PEMBAHASAN .....</b>		<b>35</b>
4.1	Eksperimen .....	35
4.1.1	Pengumpulan Data .....	35
4.1.2	Preprocessing Data.....	37

4.1.3	Pembagian Data .....	47
4.1.4	Pembobotan TF-IDF .....	48
4.1.5	Pelatihan dan Pengujian SVM .....	52
4.1.6	Evaluasi .....	60
4.1.7	Validasi .....	61
4.2	Hasil dan Pembahasan .....	63
4.2.1	Hasil Evaluasi .....	63
4.2.2	Hasil Validasi .....	67
4.2.3	Analisis Hasil .....	71
BAB V PENUTUP .....		75
5.1	Kesimpulan .....	75
5.2	Saran .....	76
DAFTAR PUSTAKA .....		77

## DAFTAR TABEL

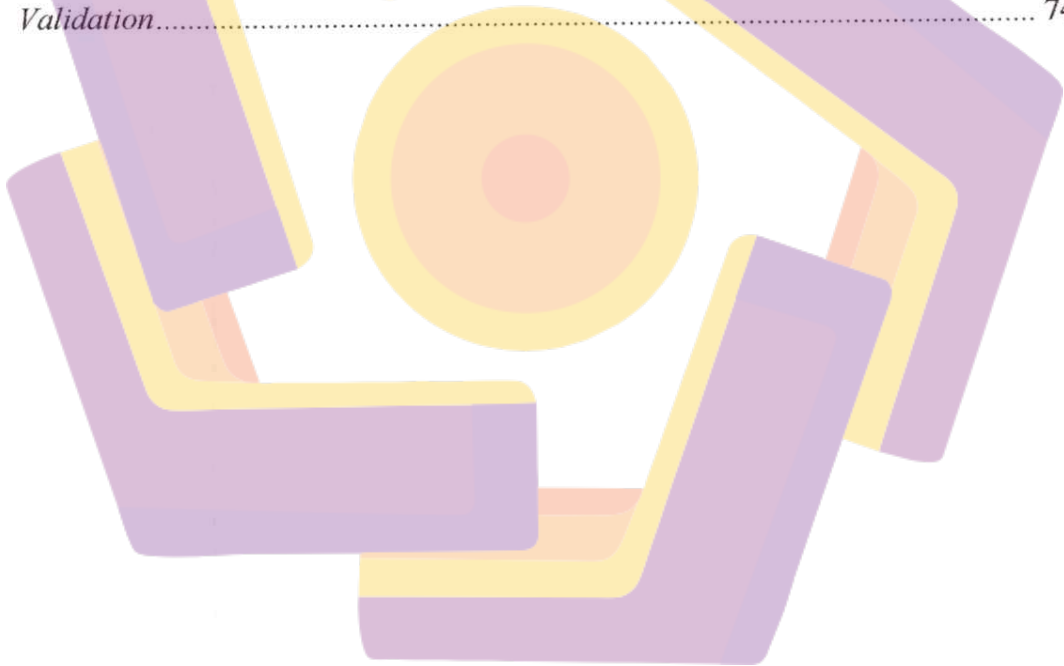
Tabel 2.1 Perbedaan Penelitian .....	9
Tabel 3.1 Spesifikasi Perangkat keras dan Perangkat Lunak.....	23
Tabel 3.2 Detail Dataset Penelitian.....	24
Tabel 3.3 Perbandingan Model <i>Preprocessing</i> .....	28
Tabel 4.1 Contoh pembobotan Term Frequency.....	50
Tabel 4.2 Contoh pembobotan Inverse Document Frequency .....	50
Tabel 4.3 Contoh hasil pembobotan TF-IDF .....	51
Tabel 4.4 Data latih dalam format Support Vector.....	53
Tabel 4.5 Hasil kernelisasi matrix $x_i^T x$ .....	58
Tabel 4.6 Hasil Evaluasi Model 1 .....	63
Tabel 4.7 Hasil Evaluasi Model 2 .....	65
Tabel 4.8 Hasil Validasi Model 1 .....	67
Tabel 4.9 Hasil Validasi Model 2 .....	69



## DAFTAR GAMBAR

Gambar 2.1 (a) Pencarian Hyperplane (b) Hyperplane Terbaik .....	15
Gambar 2.2 Fungsi $\Phi$ memetakan data ke ruang vector yang berdimensi tinggi. ....	19
Gambar 2.3 Ilustrasi Confusion Matrix .....	20
Gambar 2.4 Ilustrasi K-Fold Cross Validation .....	22
Gambar 3.1 Gambaran Umum Alur Penelitian.....	25
Gambar 3.2 Diagram alur tahapan pre-processing model 1.....	29
Gambar 3.3 Diagram Alur Tahapan <i>Preprocessing</i> Model 2 .....	29
Gambar 3.4 Diagram alur tahapan SVM.....	33
Gambar 4.1 Instalasi <i>GetOldTweets-python</i> .....	35
Gambar 4.2 Mengubah <i>OutputFile</i> di <i>Exporter.py</i> .....	36
Gambar 4.3 Pengambilan Data dengan <i>GetOldTweets-python</i> .....	36
Gambar 4.4 <i>Import</i> Dataset Menggunakan <i>Library Pandas</i> .....	37
Gambar 4.5 Implementasi <i>Case Folding</i> .....	37
Gambar 4.6 Hasil Implementasi <i>Case Folding</i> .....	38
Gambar 4.7 Implementasi <i>Cleaning</i> .....	39
Gambar 4.8 Hasil Implementasi <i>Cleaning</i> .....	40
Gambar 4.9 Implementasi Penggantian Kata <i>Slang</i> 1.....	40
Gambar 4.10 Implementasi Penggantian Kata <i>Slang</i> 2.....	41
Gambar 4.11 Hasil Implementasi Penggantian Kata <i>Slang</i> .....	42
Gambar 4.12 Implementasi Penghapusan <i>Stopwords</i> .....	42
Gambar 4.13 Hasil Implementasi Penghapusan <i>Stopwords</i> .....	43
Gambar 4.14 Implementasi <i>Stemming</i> .....	44
Gambar 4.15 Hasil Implementasi <i>Stemming</i> .....	44
Gambar 4.16 Implementasi <i>Filtering</i> .....	45
Gambar 4.17 Hasil Implementasi <i>Filtering</i> .....	45
Gambar 4.18 Implementasi <i>Tokenization</i> .....	46
Gambar 4.19 Hasil Implementasi <i>Tokenization</i> .....	47
Gambar 4.20 Implementasi Pembagian Data.....	47
Gambar 4.21 Implementasi TF-IDF.....	48

Gambar 4.22 Implementasi Pelatihan dan Pengujian SVM.....	52
Gambar 4.23 Hasil Implementasi Pelatihan dan Pengujian SVM .....	60
Gambar 4.24 Implementasi Evaluasi Model.....	61
Gambar 4.25 Implementasi Evaluasi Model.....	62
Gambar 4.26 Grafik Garis Hasil Akurasi dengan <i>Confusion Matrix</i> Model 1 .....	64
Gambar 4.27 Grafik Garis Hasil Akurasi <i>Confusion Matrix</i> Model 2.....	66
Gambar 4.28 Grafik Garis Hasil Akurasi <i>K-Fold Cross Validation</i> Model 1 .....	68
Gambar 4.29 Grafik Garis Hasil Akurasi <i>K-Fold Cross Validation</i> Model 2 .....	70
Gambar 4.30 Grafik Garis Perbandingan Akurasi <i>Confusion Matrix</i> .....	71
Gambar 4.31 Grafik Batang Perbandingan Rata-rata Akurasi <i>Confusion Matrix</i>	72
Gambar 4.32 Grafik Garis Perbandingan Akurasi <i>K-Fold Cross Validation</i> .....	73
Gambar 4.33 Grafik Batang Perbandingan Rata-rata Akurasi <i>K-Fold Cross Validation</i> .....	74



## INTISARI

*Preprocessing* adalah tugas penting untuk analisis sentimen karena informasi tekstual membawa banyak data yang tidak terstruktur dan mengandung *noise*. *Stemming* maupun penghapusan *stopword* adalah teknik *preprocessing* yang cukup populer untuk klasifikasi teks. Namun, penelitian sebelumnya memberikan hasil yang berbeda mengenai pengaruh kedua metode tersebut terhadap akurasi pada klasifikasi sentimen.

Oleh karena itu, penelitian ini melakukan investigasi lebih lanjut tentang pengaruh *stemming* dan penghapusan *stopword* pada analisis sentimen bahasa Indonesia. Selain itu, penelitian ini mengusulkan dua kondisi *preprocessing* yang menggunakan *stemming* dan penghapusan *stopword* dan tanpa menggunakan keduanya. *Support Vector Machine* digunakan untuk algoritma klasifikasi dan TF-IDF sebagai skema pembobotan. Hasilnya dievaluasi menggunakan *confusion matrix* dan kemudian divalidasi menggunakan metode *k-fold cross validation*.

Hasil eksperimen menunjukkan bahwa semua akurasi tidak meningkat dan cenderung menurun ketika melakukan skenario *stemming* dan penghapusan *stopword*. Karya ini menyimpulkan bahwa penerapan teknik *stemming* dan penghapusan *stopword* tidak berpengaruh signifikan terhadap akurasi analisis sentimen pada dokumen teks bahasa Indonesia.

**Kata Kunci:** *stemming*, penghapusan *stopword*, *preprocessing*, analisis sentimen, *svm*

## **ABSTRACT**

*Preprocessing is an essential task for sentiment analysis since textual information carries a lot of noisy and unstructured data. Both stemming and stopword removal are pretty popular preprocessing techniques for text classification. However, the prior research gives different results concerning the influence of both methods toward accuracy on sentiment classification.*

*Therefore, this research conducts further investigations about the effect of stemming and stopword removal on Indonesian language sentiment analysis. Furthermore, we propose two preprocessing conditions which are with using both stemming and stopword removal and without using both. Support Vector Machine was used for the classification algorithm and TF-IDF as a weighting scheme. The result was evaluated using the confusion matrix and then validated using k-fold cross-validation methods.*

*The experiments result show that all accuracy did not improve and tends to decrease when performing stemming or stopword removal scenarios. This work concludes that the application of stemming and stopword removal technique does not significantly affect the accuracy of sentiment analysis in Indonesian text documents.*

**Keywords:** *stemming, stopword removal, preprocessing, sentiment analysis, svm*