

**PENGARUH PENGGUNAAN STEMMING DAN LEMMATIZATION
TERHADAP AKURASI ANALISIS SENTIMEN**

SKRIPSI



disusun oleh

Fadel Maulana Ichsan

16.11.0587

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2020**

**PENGARUH PENGGUNAAN STEMMING DAN LEMMATIZATION
TERHADAP AKURASI ANALISIS SENTIMEN**

SKRIPSI

untuk memenuhi sebagian persyaratan
mencapai gelar Sarjana S1
pada jurusan Informatika



disusun oleh

Fadel Maulana Ichsan

16.11.0587

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2020**

PERSETUJUAN

SKRIPSI

**PENGARUH PENGGUNAAN STEMMING DAN LEMMATIZATION
TERHADAP AKURASI ANALISIS SENTIMEN**

yang disusun oleh

Fadel Maulana Ichsan

16.11.0587

telah disetujui oleh Dosen Pembimbing Skripsi

pada tanggal 4 Juli 2020

Dosen Pembimbing,

Mardhiya Hayaty. S.T, M.Kom

NIK. 190302108

PENGESAHAN**SKRIPSI****PENGARUH PENGGUNAAN STEMMING DAN LEMMATIZATION****TERHADAP AKURASI ANALISIS SENTIMEN**

yang disusun oleh

Fadel Maulana Ichsan

16.11.0587

telah dipertahankan di depan Dewan Penguji
pada tanggal 20 Juli 2020

Susunan Dewan Penguji

Nama Penguji

Ainul Yaqin, M.Kom
NIK. 190302255

Majid Rahardi, S.Kom., M.Eng.
NIK. 190302393

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

Tanda Tangan

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

DEKAN FAKULTAS ILMU KOMPUTER

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

PERNYATAAN

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, 20 Juli 2020



NIM. 16.11.0587

MOTTO

“Bermimpilah setinggi langit, karena bila kau jatuh. Kau akan jatuh diantara bintang-bintang”

-Ir Soekarno

"Ask not what your country can do for you, but what you can do for your country."

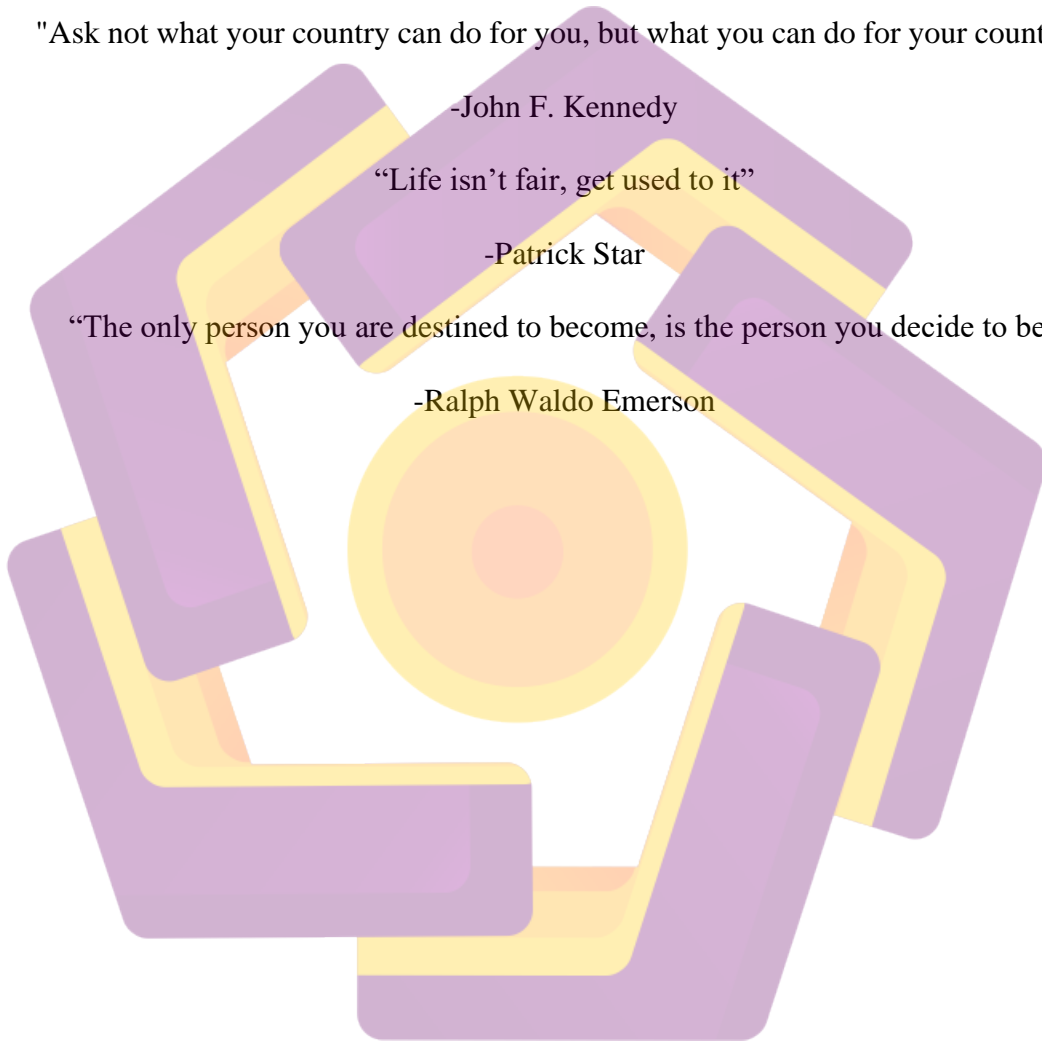
-John F. Kennedy

“Life isn’t fair, get used to it”

-Patrick Star

“The only person you are destined to become, is the person you decide to be”

-Ralph Waldo Emerson



PERSEMBAHAN

Puji Syukur saya panjatkan kepada Allah SWT, karena berkat rahmat dan karunia-Nya sehingga saya dapat menyelesaikan skripsi ini. Untuk itu skripsi ini saya persembahkan untuk:

1. Kedua orang tua yang selalu memberikan motivasi, dorongan baik berupa materi maupun non materi. Berkat doa dari orang tua saya lah skripsi ini dapat terselesaikan.
2. Ibu Mardhiya Hayaty selaku dosen pembimbing yang senantiasa sabar dalam membimbing saya dan memberikan masukan serta saran terhadap skripsi yang saya kerjakan agar dapat terselesaikan dengan baik.
3. Seluruh Dosen Universitas Amikom Yogyakarta yang telah banyak memberikan ilmu pengetahuan serta informasi kepada saya yang sangat berguna dalam penyusunan skripsi ini.
4. Keluarga besar yang rutin menanyakan “kapan lulus”, terimakasih untuk support dan motivasinya sehingga saya bisa menyelesaikan skripsi ini.
5. Teman-teman grup “ Daily New Jokes” yang berisikan sobat sambat yang senantiasa menjadi ice breaker kala jenuh mengerjakan skripsi.
6. Teman-teman dari S1 IF-09 Universitas Amikom Yogyakarta yang selalu mendukung saya dari semester 1 hingga sekarang.
7. Creator dan pengguna Google, Github, Stackoverflow, Mas-mas youtuber tutorial dari India yang telah banyak berjasa membantu saya dalam pemecahan masalah skripsi ini.

KATA PENGANTAR

Puji dan syukur saya panjatkan kepada Tuhan Yang Maha Esa yang telah memberikan rahmat, hidayah dan kekuatan sehingga saya dapat menyelesaikan skripsi yang berjudul Pengaruh Penggunaan Stemming Dan Lemmatization Terhadap Akurasi Analisis Sentimen.

Skripsi ini saya buat guna menyelesaikan studi jenjang Strata Satu (S1) pada program studi Informatika fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta. Selain itu juga merupakan suatu bukti bahwa mahasiswa telah menyelesaikan kuliah jenjang program strata 1 dan untuk memperoleh gelar Sarjana Komputer.

Dengan selesainya skripsi ini, Maka pada kesempatan ini saya mengucapkan terima kasih kepada :

1. Bapak Prof. Dr. M. Suyanto, MM. selaku Rektor Universitas AMIKOM Yogyakarta.
2. Bapak Sudarmawan, MT. selaku Ketua Program Studi S1 Informatika.
3. Ibu Mardhiya Hayaty, S.T, M.Kom. selaku Dosen Pembimbing, yang telah memberikan pengarahan yang sangat membantu dalam proses pembuatan skripsi ini.

Yogyakarta, 5 Juli 2020

Fadel Maulana Ichsan

DAFTAR ISI

| | |
|--|-------------|
| JUDUL | I |
| LEMBAR PERSETUJUAN | II |
| LEMBAR PENGESAHAN | III |
| PERNYATAAN..... | III |
| MOTTO | V |
| PERSEMBAHAN..... | VI |
| KATA PENGANTAR..... | VII |
| DAFTAR ISI..... | VIII |
| DAFTAR TABEL | XI |
| DAFTAR GAMBAR..... | XII |
| INTISARI | XIII |
| ABSTRACT | XIV |
| BAB I PENDAHULUAN..... | 1 |
| 1.1 LATAR BELAKANG..... | 1 |
| 1.2 RUMUSAN MASALAH | 3 |
| 1.3 HIPOTESIS | 3 |
| 1.4 BATASAN MASALAH | 3 |
| 1.5 MAKSUD DAN TUJUAN PENELITIAN | 3 |
| 1.6 MANFAAT PENELITIAN..... | 4 |

| | | |
|--|--|-----------|
| 1.7 | METODE PENELITIAN..... | 4 |
| 1.8 | SISTEMATIKA PENULISAN | 7 |
| BAB II LANDASAN TEORI | | 8 |
| 2.1 | KAJIAN PUSTAKA | 8 |
| 2.2 | DASAR TEORI..... | 11 |
| 2.2.1. | <i>DATA MINING</i> | 11 |
| 2.2.2. | <i>TEXT MINING</i> | 12 |
| 2.2.3. | ANALISIS SENTIMEN | 14 |
| 2.2.4. | <i>PREPROCESSING DATA</i> | 15 |
| 2.2.5. | <i>STEMMING</i> | 16 |
| 2.2.6. | <i>LEMMATIZATION</i> | 23 |
| 2.2.7. | <i>TERM FREQUENCY - INVERSE DOCUMENT FREQUENCY</i> | 25 |
| 2.2.8. | <i>SUPPORT VECTOR MACHINE</i> | 27 |
| 2.2.9. | <i>CONFUSION MATRIX</i> | 31 |
| 2.2.10. | VALIDASI..... | 32 |
| 2.2.11. | UJI SIGNIFIKASI..... | 34 |
| BAB III METODE PENELITIAN | | 36 |
| 3.1. | INSTRUMEN PENELITIAN..... | 36 |
| 3.1.1. | Perangkat Keras (Hardware)..... | 36 |
| 3.1.2. | Perangkat Lunak (Software)..... | 36 |
| 3.2. | TAHAPAN PENELITIAN | 37 |
| 3.2.1. | Persiapan Dataset..... | 38 |

| | |
|---|----|
| 3.2.2. <i>Preprocessing</i> | 38 |
| 3.2.3. <i>Stemming</i> | 39 |
| 3.2.4. <i>Lemmatization</i> | 40 |
| 3.2.5. Pembobotan Kata | 40 |
| 3.2.6. Implementasi Algoritma | 41 |
| 3.2.7. Evaluasi | 41 |
| 3.2.8. Validasi | 41 |
| 3.2.9. Uji Signifikasi | 42 |
| BAB IV ANALISIS DAN PEMBAHASAN | 43 |
| 4.1 PERSIAPAN DATASET | 43 |
| 4.2 PREPROCESSING | 44 |
| 4.2.1. <i>Stemming</i> | 51 |
| 4.2.2. <i>Lemmatization</i> | 53 |
| 4.3 PEMBOBOTAN KATA | 56 |
| 4.4 IMPLEMENTASI ALGORITMA | 60 |
| 4.5 EVALUASI | 64 |
| 4.6 VALIDASI | 68 |
| 4.7 UJI SIGNIFIKASI | 70 |
| BAB V PENUTUP | 75 |
| 5.1. KESIMPULAN | 75 |
| 5.2. SARAN | 75 |
| DAFTAR PUSTAKA | 76 |

DAFTAR TABEL

| | |
|---|----|
| Tabel 2. 1 Perbandingan penelitian..... | 10 |
| Tabel 2. 2 <i>Confusion Matrix</i> | 31 |
| Tabel 4. 1 <i>Data Cleaning</i> | 45 |
| Tabel 4. 2 Stopwords Removal..... | 47 |
| Tabel 4. 3 Case Folding | 48 |
| Tabel 4. 4 Tokenizing | 50 |
| Tabel 4. 5 Stemming | 52 |
| Tabel 4. 6 Lemmatization | 54 |
| Tabel 4. 7 Dokumen Text..... | 57 |
| Tabel 4. 8 Perhitungan TFIDF | 57 |
| Tabel 4. 9 Matriks Kernel Linear..... | 62 |
| Tabel 4. 10 Confusion Matrix Stemming..... | 66 |
| Tabel 4. 11 Confusion Matrix Lemmatization..... | 67 |
| Tabel 4. 12 Perbandingan k-Fold Cross Validation..... | 68 |
| Tabel 4. 13 Data kelas interval data Stemming..... | 73 |
| Tabel 4. 14 Data kelas interval data Lemmatization..... | 73 |
| Tabel 4. 15 Nilai $S_n(X)$ | 74 |

DAFTAR GAMBAR

| | |
|--|----|
| Gambar 2. 1 Alur proses <i>text mining</i> | 13 |
| Gambar 2. 2 Tipe Stemming | 17 |
| Gambar 2. 3 <i>Flowchart Lemmatization</i> | 24 |
| Gambar 2. 4 <i>Hyperplane</i> | 28 |
| Gambar 2. 5 Ilustrasi <i>k-fold validation</i> | 33 |
| Gambar 3. 1 Alur Penelitian..... | 37 |
| Gambar 4. 1 dataset..... | 43 |
| Gambar 4. 2 <i>Script Cleaning Tag HTML</i> | 44 |
| Gambar 4. 3 <i>Script Data Cleaning</i> | 45 |
| Gambar 4. 4 <i>Script Stopword Removal</i> | 46 |
| Gambar 4. 5 <i>Script Case Folding</i> | 48 |
| Gambar 4. 6 <i>Script Tokenizing</i> | 49 |
| Gambar 4. 7 <i>Script Stemming</i> | 51 |
| Gambar 4. 8 <i>Script Lemmatization</i> | 53 |
| Gambar 4. 9 <i>Split Dataset</i> | 55 |
| Gambar 4. 10 <i>Script TFIDF</i> | 56 |
| Gambar 4. 11 <i>Script Train SVM</i> | 60 |
| Gambar 4. 12 <i>Script Test SVM</i> | 63 |
| Gambar 4. 13 <i>Script Confusion Matrix</i> | 65 |
| Gambar 4. 14 <i>Script k-Fold Cross Validation</i> | 68 |

INTISARI

Analisis sentimen adalah studi komputasional dari opini-opini orang dengan menggunakan ilmu *text mining*. Analisis sentimen akan mengelompokkan teks yang ada dalam kalimat atau dokumen untuk mengetahui pendapat yang dikemukakan dalam kalimat atau dokumen tersebut, bisa berupa sentimen positif, maupun negatif. Sebelum dapat mengelompokkan sentimen menjadi positif maupun negatif, data teks akan melalui tahap *preprocessing* terlebih dahulu. *Preprocessing* berguna untuk membersihkan data dari *noise*, serta mempermudah *classifier* untuk melakukan klasifikasi. Salah satu tahapan *preprocessing* data adalah *stemming* dan *lemmatization*. *Stemming* dan *lemmatization* adalah salah satu tahap *preprocessing* untuk mengubah kata berimbuhan menjadi kata dasar. Walaupun memiliki fungsi yang sama namun terdapat perbedaan dalam implementasi *stemming* dan *lemmatization*. Pada *stemming* perubahan hanya dilakukan dengan memotong/ menghapus imbuhan kata. Sedangkan *lemmatization* memiliki alur yang lebih kompleks dengan melibatkan kamus bahasa untuk mencari kata dasar (*root*).

Pada penelitian ini dilakukan komparasi antara *stemming* dengan *lemmatization* untuk kasus analisis sentimen. Dataset akan diklasifikasi dengan algoritma *support vector machine*, kemudian dievaluasi dengan *confusion matrix* dan divalidasi menggunakan metode *k-fold cross validation* untuk melihat akurasi dari masing-masing *preprocessing*. Dan hasil akurasi keduanya akan dikomparasi menggunakan uji statistik untuk melihat apakah perbedaan diantara keduanya signifikan atau tidak.

Setelah dilakukan evaluasi dengan *confusion matrix*, *preprocessing* dengan *stemming* menghasilkan akurasi sebesar 85% sedangkan *lemmatization* menghasilkan akurasi 84%. Namun setelah dilakukan uji signifikansi, ternyata perbedaan dari keduanya tidak signifikan.

Kata-kunci: analisis sentimen, stemming, lemmatization, uji signifikan

ABSTRACT

Sentiment analysis is a computational study of people's opinions using text mining science. Sentiment analysis will classify the text in a sentence or document to find out the opinions expressed in the sentence or document, it can be in the form of positive or negative sentiments. Before being able to classify sentiments into positive or negative, text data will go through the preprocessing stage first. Preprocessing is useful for cleaning data from noise, as well as making it easier for classifiers to classify. One of the stages of data preprocessing is stemming and lemmatization. Stemming and lemmatization are preprocessing steps to convert affix words into basic words. Even though they have the same function, there are differences in the implementation of stemming and lemmatization. In stemming, changes are only done by cutting / removing affixes of words. Meanwhile, lemmatization has a more complex flow by involving a language dictionary to look for the root word.

In this study, a comparison between stemming and lemmatization was carried out for sentiment analysis cases. The dataset will be classified with a support vector machine algorithm, then evaluated with a confusion matrix and validated using the k-fold cross validation method to see the accuracy of each preprocessing. And the results of the accuracy of the two will be compared using statistical tests to see whether the difference between the two is significant or not.

After evaluating with confusion matrix preprocessing with stemming it produces an accuracy of 85% while lemmatization produces an accuracy of 84%. However, after the significance test was carried out, it turned out that the difference between the two was not significant.

Keywords: *sentiment analysis, stemming, lemmatization, significant test*