

SKRIPSI



**PERBANDINGAN HASIL PENERAPAN ALGORITMA KLASIFIKASI
DAN NATURAL LANGUAGE PROCESSING TERHADAP DATA
KEPUASAN PENGGUNA LAYANAN TRANSPORTASI UMUM MRT
JAKARTA**

MUHAMMAD NABIL NUFAIL PRIBADI

NIM. 1910511106

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
JUNI 2024**

SKRIPSI



**PERBANDINGAN HASIL PENERAPAN ALGORITMA KLASIFIKASI
DAN NATURAL LANGUAGE PROCESSING TERHADAP DATA
KEPUASAN PENGGUNA LAYANAN TRANSPORTASI UMUM MRT
JAKARTA**

MUHAMMAD NABIL NUFAIL PRIBADI

NIM. 1910511106

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
JUNI 2024**

PERNYATAAN ORISINALITAS

Skripsi ini adalah hasil karya sendiri dan semua sumber yang dikutip maupun yang ditunjuk telah saya nyatakan dengan benar:

Nama : Muhammad Nabil Nufail Pribadi
NIM. : 1910511106
Program Studi : S1 Informatika
Tanggal : 1 Juli 2024
Judul Skripsi/TA. : **Perbandingan Hasil Penerapan Algoritma Klasifikasi Dan Natural Language Processing Terhadap Data Kepuasan Pengguna Layanan Transportasi Umum MRT Jakarta**

Bilamana pada kemudian hari ditemukan ketidaksesuaian dengan pernyataan saya ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Jakarta 1 Juli 2024
Yang Menyatakan,



A handwritten signature of Muhammad Nabil Nufail Pribadi is overlaid on a digital stamp. The stamp features the Indonesian national emblem (Garuda Pancasila) in the center, surrounded by the text "REPUBLIK INDONESIA" and "REPUBLIQUE INDONÉSIE". Below the emblem, the text "EX-591-AK-283624271" is printed.

Muhammad Nabil Nufail Pribadi

PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS

Sebagai civitas akademika Universitas Pembangunan Nasional Veteran Jakarta, saya yang bertanda tangan di bawah ini:

Nama : Muhammad Nabil Nufail Pribadi
NIM. : 1910511106
Fakultas : Ilmu Komputer
Program Studi : S1 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 tulis ilmiah saya yang dipublikasikan dengan judul:

Perbandingan Hasil Penerapan Algoritma Klasifikasi Dan Natural Language Processing Terhadap Data Kepuasan Pengguna Layanan Transportasi Umum MRT Jakarta

Beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional Veteran Jakarta berhak menyimpan, mengalih media atau memformatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan Skripsi/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 : 1 Juli 2024

Yang Menyatakan,

Muhammad Nabil Nufail Pribadi

LEMBAR PENGESAHAN

Skripsi/Tugas Akhir ini diajukan oleh:

Nama : Muhammad Nabil Nufail Pribadi
NIM. : 1910511106
Program Studi : S1 Informatika
Judul Skripsi/TA. : Perbandingan Hasil Penerapan Algoritma Klasifikasi Dan *Natural Language Processing* Terhadap Data Kepuasan Pengguna Layanan Transportasi Umum MRT Jakarta

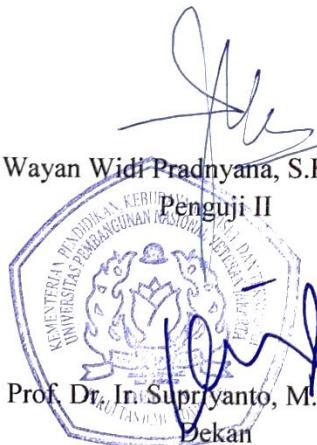
Telah berhasil dipertahankan di hadapan Tim Penguji dan diterima sebagai bagian persyaratan yang diperlukan untuk memperoleh gelar Sarjana pada Program Studi S1 Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional Veteran Jakarta.



Dr. Widya Cholil, S.Kom., M.IT.
Pengaji I



I Wayan Widi Pradnyana, S.Kom. M.TI.
Pengaji II



Prof. Dr. Ir. Supriyanto, M.Sc., IPM.
Bekan



Iin Ernawati, S.Kom., M.Si.
Pembimbing



Dr. Widya Cholil, S.Kom., M.IT.
Kepala Program Studi

Ditetapkan di : Jakarta
Tanggal Ujian : 27 Juni 2024

KATA PENGANTAR

Puji syukur penulis panjatkan kepada Allah SWT atas berkat dan rahmatnya, sehingga penulis mampu menyelesaikan penulisan proposal skripsi ini dengan lancar. Penulisan proposal skripsi dengan judul “Perbandingan Hasil Penerapan Algoritma Klasifikasi Dan *Natural Language Processing* Terhadap Data Kepuasan Pengguna Layanan Transportasi Umum MRT Jakarta” telah disusun sejak bulan Agustus 2023.

Dalam pelaksanaan penelitian ini, serta dalam penulisan proposal ini, penulis menyadari bahwa terdapat pihak lain yang turut membantu melancarkan penyelesaiannya. Untuk itu, penulis akan mengucapkan terima kasih kepada:

1. Orang tua penulis, yang telah memfasilitasi penulis dengan segala kebutuhan, sehingga penulis mampu menyelesaikan proposal ini dengan lancar.
2. Ibu Iin Ernawati, S. Kom., M.Si. selaku dosen pembimbing skripsi penulis yang telah membantu penulis selama penelitian berlangsung.
3. Bapak Prof. Dr. Ir. Supriyanto, M.Sc., IPM. selaku Dekan Fakultas Ilmu Komputer.
4. Ibu Dr. Widya Cholil, S. Kom., M. IT. selaku Kepala Program Studi Informatika.
5. Teman-teman penulis yang telah membantu dengan dukungan dan sedikit arahan yang diberikan kepada penulis.
6. Pihak lainnya yang turut membantu penulis dalam penyusunan skripsi ini yang tidak tercantum di atas, penulis sampaikan ucapan terima kasih.

Jakarta, 20 Mei 2024

Penulis

Abstrak

Seiring dengan bertambahnya jumlah penduduk dan terbatasnya lahan untuk tempat tinggal di DKI Jakarta, sudah tidak memungkinkan bagi setiap masyarakat untuk memiliki transportasi pribadi masing-masing. Hal ini dapat dibuktikan dengan kemacetan pada DKI Jakarta dengan penyebab terbesar berasal dari transportasi pribadi. MRT Jakarta merupakan salah satu solusi pemerintah DKI Jakarta dalam upaya mengatasi kemacetan. Namun tingkat kemacetan Jakarta kembali meningkat pada tahun 2023 dikarenakan bertambahnya kembali jumlah transportasi pribadi pada lalu lintas. Oleh karena itu, penulis memutuskan untuk melakukan klasifikasi terhadap kepuasan pengguna layanan transportasi umum MRT Jakarta untuk mengetahui alasan yang menyebabkan masyarakat bersedia atau enggan memilih untuk memanfaatkan sarana transportasi MRT Jakarta, yang diperoleh dari media sosial X, memanfaatkan *natural language processing* dan algoritma *Support Vector Machine*, algoritma *Random Forest Classifier*, serta algoritma *Logistic Regression* multinomial. Data tersebut berjumlah sebanyak 525 *post*, dengan kategori ‘positif’ sebanyak 222 data, kategori ‘negatif’ sebanyak 185 data, dan kategori ‘netral’ sebanyak 118 data. Dengan pembagian terhadap dataset berdasarkan perbandingan 80% data latih, dan 20% data uji, model klasifikasi dengan hasil paling akurat pada penelitian ini, yaitu model dengan algoritma *Random Forest Classifier*, menggunakan parameter terbaik yang diperoleh melalui teknik *hyperparameter tuning*, dengan nilai *class_weight='balanced'*, nilai *max_depth=350*, nilai *min_samples_split=5*, serta nilai *n_estimators=200*, menghasilkan nilai *accuracy* sebesar 81%, nilai *precision* sebesar 82.3%, nilai *recall* sebesar 81%, dan nilai *f1-score* sebesar 81.5%, serta berhasil secara akurat memprediksi 4 dari 4 sampel data baru berdasarkan target kelasnya.

Kata Kunci: *Natural Language Processing*, Klasifikasi, MRT Jakarta

Abstract

Along with the increasing population and limited amount of land for housing in DKI Jakarta, it is no longer possible for each person to privately own a car-based transportation. This can be proven by the occurrences of traffic jams in DKI Jakarta with the biggest cause coming from the usage of private transportation, namely cars. MRT Jakarta is one of the solutions proposed by the government of DKI Jakarta as one of their efforts to overcome traffic jams. However, Jakarta's congestion level increases again in 2023 due to the increased usage of private transportation. Therefore, the author decided to classify the opinion of users of MRT Jakarta to find out the reasons why people choose or are reluctant to utilize MRT Jakarta as their main method of transportation, obtained from the social media X, utilizing natural language processing and machine learning algorithm such as Support Vector Machine, Random Forest Classifier, and multinomial Logistic Regression. The data amounts to 525 posts, with 222 in the 'positive' category, 185 in the 'negative' category, and 118 in the 'neutral' category. By dividing the dataset based on a comparison of 80% training data and 20% test data, the classification model with the highest accuracy score in this research, namely the model with the Random Forest Classifier algorithm, using the best parameters obtained through the usage of hyperparameter tuning technique, with the value of class_weight='balanced', value max_depth=350, min_samples_split=5, and n_estimators=200, resulting in an accuracy value of 81%, a precision value of 82.3%, a recall value of 81%, and an f1-score value of 81.5%, and successfully predicted 4 of 4 new data samples based on their target classes.

Keywords: Natural Language Processing, Classification, MRT Jakarta

DAFTAR ISI

PERNYATAAN ORISINALITAS	i
PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS	ii
LEMBAR PENGESAHAN	iii
KATA PENGANTAR	iv
Abstrak	v
<i>Abstract</i>	vi
DAFTAR ISI.....	vii
DAFTAR GAMBAR	x
DAFTAR TABEL.....	xi
DAFTAR LAMPIRAN	xii
BAB I	1
PENDAHULUAN	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah	2
1.3. Ruang Lingkup	3
1.4. Luaran yang Diharapkan	3
1.5. Tujuan Penelitian.....	3
1.6. Manfaat Penelitian.....	4
BAB II.....	5
TINJAUAN PUSTAKA	5
2.1. <i>Machine Learning</i>	5
2.1.1. <i>Supervised Learning</i>	5
2.1.2. <i>Unsupervised Learning</i>	6
2.2. Klasifikasi.....	6
2.3. Analisis Sentimen.....	6
2.4. <i>Natural Language Processing</i>	7
2.5. <i>Text Mining</i>	7
2.6. Python.....	8
2.7. <i>Random Forest Classifier</i>	9
2.8. <i>Logistic Regression</i>	9
2.9. <i>Support Vector Machine</i>	10

2.10. <i>Pre-processing</i>	10
2.10.1. <i>Case Folding</i>	11
2.10.2. <i>Tokenizing</i>	11
2.10.3. <i>Stopword Removal</i>	11
2.10.4. <i>Stemming</i>	11
2.11. <i>Feature Engineering</i>	11
2.12. <i>Confusion Matrix</i>	12
2.13. Studi Relevan	13
BAB III	18
METODOLOGI PENELITIAN.....	18
3.1. Alur Penelitian.....	18
3.2. Tahap Penelitian	19
3.2.1. Perumusan Masalah	19
3.2.2. Studi Literatur	20
3.2.3. Pengumpulan Data.....	20
3.2.4. <i>Pre-processing</i>	21
3.2.4.1 <i>Exploratory Data Analysis</i>	21
3.2.4.2 <i>Case-Folding</i>	21
3.2.4.3 <i>Noise Reduction/Text Cleanup</i>	21
3.2.4.4 <i>Tokenization</i>	21
3.2.4.5 <i>Stopword Removal</i>	22
3.2.4.6 <i>Normalization/Spelling Correction</i>	22
3.2.4.7 <i>Stemming</i>	22
3.2.5. Pembagian Data	22
3.2.6. <i>Feature Engineering</i>	22
3.2.7. Perancangan Model Klasifikasi	23
3.2.8. Uji Coba Dengan Data Baru	23
3.2.9. Evaluasi.....	23
3.3. Waktu dan Tempat Penelitian	23
3.4. Perangkat Penelitian	24
3.4.1. Perangkat Keras	24
3.4.2. Perangkat Lunak	24
3.5. Jadwal Penelitian	24
BAB IV	26

HASIL DAN PEMBAHASAN.....	26
4.1. Data	26
4.2. <i>Pre-Processing</i>	26
4.2.1 <i>Exploratory Data Analysis</i>	27
4.2.2 <i>Case-Folding</i>	31
4.2.3 <i>Text Cleanup/Noise Reduction</i>	33
4.2.4 <i>Tokenization</i>	34
4.2.5 <i>Stopword Removal</i>	35
4.2.6 <i>Spelling Correction</i>	37
4.2.7 <i>Stemming</i>	39
4.3. Pembagian Data.....	42
4.4. <i>Feature Engineering</i> Menggunakan TF-IDF	43
4.5. Proses Klasifikasi	47
4.5.1 Klasifikasi dengan model <i>Support Vector Machine</i>	48
4.5.2 Klasifikasi dengan model <i>Random Forest Classifier</i>	51
4.5.3 Klasifikasi dengan model <i>Logistic Regression</i>	55
4.6. Evaluasi	59
BAB V	63
KESIMPULAN DAN SARAN.....	63
5.1. Kesimpulan.....	63
5.2. Saran	65
DAFTAR PUSTAKA	66
RIWAYAT HIDUP.....	69
LAMPIRAN	70
LAMPIRAN 1. DATASET PENELITIAN (100 DATA TERATAS).....	70
LAMPIRAN 2. HASIL PRAPROSSES DATA (100 DATA TERATAS)	75
LAMPIRAN 3. HASIL PEMERIKSAAN TURNITIN	111

DAFTAR GAMBAR

Gambar 2.1. Ilustrasi Proses <i>Text Mining</i>	8
Gambar 2.2. Confusion Matrix	12
Gambar 3.1. Alur Penelitian.....	18
Gambar 4.1. Visualisasi jumlah data berdasarkan kategori	28
Gambar 4.2. Visualisasi persebaran data berdasarkan panjang string	29
Gambar 4.3. Visualisasi <i>wordcloud</i> untuk keseluruhan dataset.....	29
Gambar 4.4. Visualisasi <i>wordcloud</i> untuk kategori positif.....	30
Gambar 4.5. Visualisasi <i>wordcloud</i> untuk kategori negatif.....	30
Gambar 4.6. Visualisasi <i>wordcloud</i> untuk kategori netral	31
Gambar 4.7. Visualisasi <i>wordcloud</i> kategori positif setelah <i>pre-processing</i>	41
Gambar 4.8. Visualisasi <i>wordcloud</i> kategori negatif setelah <i>pre-processing</i>	41
Gambar 4.9. Visualisasi <i>wordcloud</i> kategori netral setelah <i>pre-processing</i>	42
Gambar 4.10. <i>Confusion matrix</i> model klasifikasi <i>SVM</i> perbandingan 70:30.....	49
Gambar 4.11. <i>Confusion matrix</i> model klasifikasi <i>SVM</i> perbandingan 80:20.....	50
Gambar 4.12. <i>Classification report</i> model klasifikasi <i>SVM</i> perbandingan 70:30 .	50
Gambar 4.13. <i>Classification report</i> model klasifikasi <i>SVM</i> perbandingan 80:20 .	51
Gambar 4.14. <i>Confusion matrix</i> model klasifikasi <i>Random Forest Classifier</i> perbandingan 70:30	53
Gambar 4.15. <i>Confusion matrix</i> model klasifikasi <i>Random Forest Classifier</i> perbandingan 80:20	54
Gambar 4.16. <i>Classification report</i> model klasifikasi <i>Random Forest Classifier</i> perbandingan 70:30	54
Gambar 4.17. <i>Classification report</i> model klasifikasi <i>Random Forest Classifier</i> perbandingan 80:20	55
Gambar 4.18. <i>Confusion matrix</i> model klasifikasi <i>Logistic Regression</i> perbandingan 70:30	57
Gambar 4.19. <i>Confusion matrix</i> model klasifikasi <i>Logistic Regression</i> perbandingan 80:20	58
Gambar 4.20. <i>Classification report</i> model klasifikasi <i>Logistic Regression</i> perbandingan 70:30	58
Gambar 4.21. <i>Classification report</i> model klasifikasi <i>Logistic Regression</i> perbandingan 80:20	59

DAFTAR TABEL

Tabel 2.1. Perbandingan penelitian yang relevan	13
Tabel 3.1. Jadwal Penelitian.....	24
Tabel 4.1. Data opini pengguna MRT Jakarta	26
Tabel 4.2. Dimensi dari data	27
Tabel 4.3. Hasil pemeriksaan dari setiap kolom	27
Tabel 4.4. Jumlah data berdasarkan kategori	28
Tabel 4.5. Penerapan case-folding pada data	32
Tabel 4.6. Penerapan text cleanup pada data	33
Tabel 4.7. Penerapan proses tokenization pada data.....	34
Tabel 4.8. Pratinjau tambahan stopwords untuk data	36
Tabel 4.9. Penerapan proses stopword removal terhadap data	36
Tabel 4.10. Dictionary perbaikan penulisan kata.....	38
Tabel 4.11. Penerapan proses spelling correction terhadap data	38
Tabel 4.12. Penerapan proses stemming terhadap data.....	40
Tabel 4.13. Dimensi data setelah proses pembagian.....	42
Tabel 4.14. Sampel data untuk perhitungan teknik TF-IDF	43
Tabel 4.15. Perhitungan teknik TF-IDF terhadap tiga sampel data	44
Tabel 4.16. Parameter feature engineering dengan TfIdfVectorizer.....	46
Tabel 4.17. Value untuk hyperparameter tuning model SVM	48
Tabel 4.18. Parameter model klasifikasi algoritma SVM	48
Tabel 4.19. Hasil model klasifikasi Support Vector Machine	51
Tabel 4.20. Value hyperparameter tuning model Random Forest Classifier.....	52
Tabel 4.21. Parameter model klasifikasi algoritma Random Forest Classifier.....	52
Tabel 4.22. Hasil model klasifikasi Random Forest Classifier.....	55
Tabel 4.23. Value untuk hyperparameter tuning model Logistic Regression.....	56
Tabel 4.24. Parameter model klasifikasi algoritma Logistic Regression.....	56
Tabel 4.25. Hasil model klasifikasi Logistic Regression	59
Tabel 4.26. Perbandingan hasil model klasifikasi.....	60
Tabel 4.27. Sampel data baru	60
Tabel 4.28. Uji coba model klasifikasi dengan data baru	62

DAFTAR LAMPIRAN

LAMPIRAN 1. DATASET PENELITIAN (100 DATA TERATAS)	70
LAMPIRAN 2. HASIL PRAPROSES DATA (100 DATA TERATAS).....	75
LAMPIRAN 3. HASIL PEMERIKSAAN TURNITIN.....	111