



PERBANDINGAN METODE *RANDOM FOREST* DAN *K-NEAREST NEIGHBORS* PADA ANALISIS SENTIMEN PENGGUNA TWITTER MENGENAI PROMO GOJEK

SKRIPSI

ASHIL HAFIDH DHIYA

1810511121

UNIVERSITAS PEMBANGUNAN NASIONAL “VETERAN”

JAKARTA

FAKULTAS ILMU KOMPUTER

PROGRAM STUDI INFORMATIKA

2023



**PERBANDINGAN METODE *RANDOM FOREST* DAN *K-NEAREST NEIGHBORS* PADA ANALISIS SENTIMEN
PENGGUNA TWITTER MENGENAI PROMO GOJEK**

SKRIPSI

**Diajukan Sebagai Salah Satu Syarat untuk Memperoleh Gelar
Sarjana Komputer**

ASHIL HAFIDH DHIYA

1810511121

UNIVERSITAS PEMBANGUNAN NASIONAL “VETERAN”

JAKARTA

FAKULTAS ILMU KOMPUTER

PROGRAM STUDI INFORMATIKA

2023

PERNYATAAN BEBAS PLAGIAT

PERNYATAAN BEBAS PLAGIAT

Skripsi ini adalah hasil karya sendiri dan sumber yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar.

Nama : Ashil Hafidh Dhiya

NIM : 1810511121

Tanggal : 30 November 2022

Judul Skripsi : Perbandingan Metode Random Forest dan K-Nearest Neighbors pada Analisis Sentimen Pengguna Twitter mengenai Promo Gojek

Bilamana dikemudian hari ditemukan ketidaksamaan dengan pernyataan ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Jakarta, 30 November 2022



Ashil Hafidh Dhiya

Scanned with CamScanner

PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS

Sebagai civitas akademik Universitas Pembangunan Nasional Veteran Jakarta, saya yang bertanda tangan dibawah ini :

Nama : Ashil Hafidh Dhiya

NIM : 1810511121

Fakultas : Ilmu Komputer

Program Studi : 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 ilmiah saya yang berjudul:

Perbandingan Metode *Random Forest* dan *K-Nearest Neighbors* pada Analisis Sentimen Pengguna Twitter mengenai Promo Gojek

Beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional Veteran Jakarta berhak menyimpan, mengalih-media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan 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 : 30 November
2022

Yang menyatakan,



(Ashil Hafidh Dhiya)

LEMBAR PENGESAHAN

LEMBAR PENGESAHAN

Tugas Akhir ini diajukan oleh:

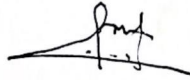
Nama : Ashil Hafidh Dhiya

NIM : 1810511121

Program Studi : SI Informatika

Judul Tugas Akhir : Perbandingan Metode Random Forest dan K-Nearest Neighbors pada Analisis Sentimen Pengguna Twitter mengenai Promo Gojek

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



Dr. Ermatita, M. Kom.

Penguji I



Helena Nurramdhani Irmanda,
S.Pd., M.Kom.

Penguji II



Ati Zaidiah, S.Kom, M.TI.

Pembimbing



Dr. Ermatita, M.Kom

Dekan



Dr. Widya Cholil, S.Kom, M.I.T.

Ketua Program Studi

Ditetapkan di : Jakarta

Tanggal Ujian : 11 Januari 2023



Scanned with CamScanner

PERBANDINGAN METODE *RANDOM FOREST* DAN *K-NEAREST NEIGHBORS* TERHADAP ANALISIS SENTIMEN PADA TWITTER MENGENAI PROMO GOJEK

Ashil Hafidh Dhiya

ABSTRAK

Layanan transportasi *online* merupakan salah satu topik yang sedang hangat dibicarakan. Banyak orang mengandalkan lalu lintas *online*, atau pengemudi dan pelanggan. Transportasi *online* menjadi topik yang ramai diperbincangkan karena kebutuhan transportasi umum di beberapa daerah sulit terpenuhi. Respon masyarakat terhadap pelayanan yang diberikan oleh penyedia jasa transportasi *online* bervariasi, ada yang memberikan hasil positif dan negatif. Analisis sentimen merupakan salah satu cara untuk memahami pendapat seseorang atau sekelompok orang. Data berupa tweet akan dikumpulkan melalui Twitter sebanyak 453 tweet dan dibagi dua menjadi 80% data latih dan 20% data uji. Setelah itu data akan di *preprocessing* menggunakan *case folding*, *cleansing data*, normalisasi, *stemming*, *stopword removal* dan tokenisasi. Setelah data di *preprocess*, maka data akan diberi *term weighting* menggunakan TF-IDF. Setelah data diberi bobot, data akan diolah menggunakan *K-Nearest Neighbors* dan *Random Forest*. Penelitian ini diharapkan bisa mendapatkan informasi akan sentimen opini publik terhadap promo Gojek serta mengetahui performa metode *K-Nearest Neighbors* dan *Random Forest*. Pada pengujian pertama yaitu menggunakan metode *K-Nearest Neighbors* diperoleh nilai akurasi sebesar 75%, nilai presisi sebesar 80%, nilai recall sebesar 64%, dan nilai *f-1 score* sebesar 71%. Pengujian kedua yaitu dengan menerapkan metode *Random Forest* dan diperoleh nilai akurasi sebesar 77%, nilai presisi sebesar 75%, nilai recall sebesar 77%, dan nilai *f-1 score* sebesar 75%. Berdasarkan perolehan hasil evaluasi tersebut, metode *Random Forest* lebih baik dibandingkan dengan metode *K-Nearest Neighbors* dengan nilai akurasi sebesar 77%.

Kata Kunci: Promo Gojek, Analisis Sentimen, *K-Nearest Neighbors*, *Random Forest*

COMPARISON OF RANDOM FOREST AND K-NEAREST NEIGHBORS METHODS TOWARDS ANALYSIS OF SENTIMENT ON TWITTER ABOUT GOJEK PROMOTION

Ashil Hafidh Dhiya

ABSTRACT

Online transportation services are one of the hotly discussed topics. Many people rely on online traffic, or drivers and customers. Online transportation is a topic that is often discussed because public transportation needs in some areas are difficult to meet. Public responses to the services provided by online transportation service providers vary, some giving positive and negative results. Sentiment analysis is one way to understand the opinion of a person or group of people. Data in the form of tweets will be collected via Twitter as many as 453 tweets and divided into 80% training data and 20% test data. After that, the data will be preprocessed using case folding, data cleansing, normalization, stemming, stopword removal and tokenization. After the data is preprocessed, the data will be given term weighting using TF-IDF. After the data is weighted, the data will be processed using K-Nearest Neighbors and Random Forest. This research is expected to be able to obtain information on the sentiment of public opinion towards the Gojek promo and determine the performance of the K-Nearest Neighbors and Random Forest methods. In the first test, using the K-Nearest Neighbors method, the accuracy value is 75%, the precision value is 80%, the recall value is 64%, and the f-1 score is 71%. The second test is by applying the Random Forest method and obtained an accuracy value of 77%, a precision value of 75%, a recall value of 77%, and an f-1 score of 75%. Based on the evaluation results, the Random Forest method is better than the K-Nearest Neighbors method with an accuracy value of 77%.

Keyword: *Gojek Promotion, Sentiment Analysis, K-Nearest Neighbors, Random Forest.*

KATA PENGANTAR

Puji dan syukur penulis panjatkan atas berkat, rahmat dan karunia-Nya kepada Allah SWT, penulis dapat menyelesaikan Tugas Akhir/Skripsi yang berjudul “Perbandingan Metode Random Forest dan K-Nearest Neighbors pada Analisis Sentimen Pengguna Twitter mengenai Promo Gojek” dengan baik.

Dalam penyelesaian tugas akhir ini, penulis tidak lupa mengucapkan terima kasih kepada pihak-pihak yang telah memberikan dukungan dan masukan kepada penulis, serta membantu penulis dalam merealisasikan tugas akhir ini, pihak-pihak tersebut adalah :

1. Kedua orang tua penulis beserta keluarga yang selalu memberikan doa dan dukungan penuh sehingga dapat menyelesaikan tugas akhir ini.
2. Ibu Ati Zaidiah S.Kom., MTI. Selaku dosen pembimbing yang telah memberikan saran dan masukan yang sangat bermanfaat.
3. Bapak/Ibu dosen Informatika Universitas Pembangunan Nasional Veteran Jakarta, terima kasih atas ilmu-ilmu yang selama ini sudah diajarkan.
4. Teman-teman seperti Nadhifa Zhafira, Vanesa, Hanif Razka, Kavindra Razik Afif, Nabila Adhari, Daniel Manurung, Krisna Jonathan, serta teman-teman terdekat yang tiada henti membantu dan memberikan dukungan dalam setiap proses penulisan tugas akhir.
5. Seluruh pihak yang terlibat dalam kelancaran pembuatan tugas akhir ini dan yang belum disebutkan di atas, penulis ucapkan terima kasih.

Akhir kata, semoga tugas akhir/skripsi ini dapat bermanfaat bagi para pembacanya.

Jakarta, 30 November 2022

Penulis

DAFTAR ISI

PERNYATAAN BEBAS PLAGIAT.....	ii
PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS	iii
LEMBAR PENGESAHAN	iv
ABSTRAK	v
ABSTRACT.....	vi
KATA PENGANTAR	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	xii
DAFTAR GAMBAR	xiii
DAFTAR LAMPIRAN.....	xiv
BAB I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	3
1.3 Tujuan Penelitian.....	3
1.4 Manfaat Penelitian.....	3
1.5 Ruang Lingkup	4
1.6 Luaran Yang Diharapkan	4
1.7 Sistematika Penulisan.....	4
BAB II TINJAUAN PUSTAKA.....	6
2.1 Promo	6
2.2 Gojek	6
2.3 Twitter	6
2.4 <i>Data Mining</i>	7

2.5	<i>Preprocess Data</i>	7
2.5.1	<i>Case Folding</i>	8
2.5.2	<i>Cleansing Data</i>	8
2.5.3	Normalisasi	8
2.5.4	<i>Stemming</i>	8
2.5.5	<i>Stopword Removal</i>	8
2.5.6	Tokenisasi	8
2.6	Analisis Sentimen.....	9
2.7	<i>Term Frequency-Invers Document Frequency (TF-IDF)</i>	9
2.8	Metode Klasifikasi	10
2.8.1	<i>K-Nearest Neighbors</i>	10
2.8.2	<i>Random Forest</i>	11
2.9	Evaluasi	12
2.9.1	Akurasi	13
2.9.2	Presisi	13
2.9.3	<i>Recall</i>	13
2.9.4	F-1 Score	13
2.10	Penelitian Terdahulu	13
BAB III METODE PENELITIAN.....		16
3.1	Tahapan Penelitian	16
3.1.1	Mengidentifikasi Masalah.....	17
3.1.2	Studi Literatur	18
3.1.3	<i>Crawling Data</i>	18
3.1.4	<i>Labeling Data</i>	18
3.1.5	<i>Preprocess Data</i>	18
3.1.6	Pembobotan TF-IDF	20

3.1.7	<i>Split Data</i>	20
3.1.8	Klasifikasi	20
3.1.9	Evaluasi	21
3.2	Perangkat Penelitian	21
3.2.1	Perangkat Keras (<i>Hardware</i>)	21
3.2.2	Perangkat Lunak (<i>Software</i>).....	21
BAB IV HASIL DAN PEMBAHASAN		23
4.1	Data	23
4.2	Pelabelan Data	23
4.3	<i>Preprocess Data</i>	25
4.3.1	<i>Case Folding</i>	25
4.3.2	<i>Cleaning Data</i>	26
4.3.3	Normalisasi	27
4.3.4	<i>Stemming</i>	28
4.3.5	<i>Stopwords Removal</i>	28
4.3.6	Tokenisasi	29
4.4	Visualisasi	30
4.5	Pembobotan TF-IDF.....	32
4.6	Klasifikasi.....	35
4.6.1	Klasifikasi dengan <i>K-Nearest Neighbors</i>	35
4.6.2	Klasifikasi dengan <i>Random Forest</i>	37
4.7	Evaluasi	38
4.7.1	Evaluasi Klasifikasi <i>K-Nearest Neighbors</i>	38
4.7.2	Evaluasi Klasifikasi <i>Random Forest</i>	39
4.7.3	Perbandingan Performa Klasifikasi <i>K-Nearest Neighbors</i> dengan <i>Random Forest</i>	41

BAB V PENUTUP.....	44
5.1 Kesimpulan.....	44
5.2 Saran.....	45
DAFTAR PUSTAKA	47
RIWAYAT HIDUP.....	51
LAMPIRAN.....	52

DAFTAR TABEL

Tabel 2. 1 Confusion Matrix	12
Tabel 4. 1 Hasil Pelabelan Data	23
Tabel 4. 2 Kuantitas Pelabelan Data	25
Tabel 4. 3 Tahapan <i>Case Folding</i>	25
Tabel 4. 4 Tahapan <i>Cleansing Data</i>	26
Tabel 4. 5 Tahapan Normalisasi.....	27
Tabel 4. 6 Tahapan <i>Stemming</i>	28
Tabel 4. 7 Tahapan <i>Stopwords Removal</i>	29
Tabel 4. 8 Tahapan Tokenisasi.....	29
Tabel 4. 9 <i>Data Training</i>	32
Tabel 4. 10 Pembobotan TF-IDF	33
Tabel 4. 11 Pembagian Data Klasifikasi	35
Tabel 4. 12 Tabel Langkah Perhitungan <i>Cosine Similarity</i>	36
Tabel 4. 13 Hasil Klasifikasi dengan K-NN	37
Tabel 4. 14 Hasil Klasifikasi dengan <i>Random Forest</i>	37
Tabel 4. 15 Confusion Matrix Klasifikasi dengan K-NN	38
Tabel 4. 16 Confusion Matrix Klasifikasi dengan <i>Random Forest</i>	40
Tabel 4. 17 Perbandingan Performa Klasifikasi	41

DAFTAR GAMBAR

Gambar 3. 1 Tahapan Penelitian	17
Gambar 4. 1 Visualisasi Data.....	30
Gambar 4. 2 Visualisasi Data Label Positif	31
Gambar 4. 3 Visualisasi Data Label Negatif.....	32

DAFTAR LAMPIRAN

Lampiran 1 Kamus Normalisasi	52
Lampiran 2 Stopwords Bahasa Indonesia	71
Lampiran 3 Hasil <i>Random Forest</i>	77
Lampiran 4 Hasil Lengkap KNN	90