



**PERBANDINGAN ALGORITMA K-NEAREST NEIGHBOR,
RANDOM FOREST, DAN LOGISTIC REGRESSION
TERHADAP KATEGORI REVIEW PADA MARKETPLACE
XYZ**

SKRIPSI

ARIA NANDA HERDIWAN

1910511083

**PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
2023**



**PERBANDINGAN ALGORITMA K-NEAREST NEIGHBOR,
RANDOM FOREST, DAN LOGISTIC REGRESSION
TERHADAP KATEGORI REVIEW PADA MARKETPLACE
XYZ**

SKRIPSI

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

ARIA NANDA HERDIWAN

1910511083

**PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
2023**

PERNYATAAN ORISINALITAS

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

Nama : Aria Nanda Herdiawan

NIM : 1910511083

Tanggal : 19 Januari 2024

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

Jakarta, 19 Januari 2024



(Aria Nanda Herdiawan)

LEMBAR PERSETUJUAN PUBLIKASI

Sebagai civitas akademik Universitas Pembangunan Nasional Veteran Jakarta, saya yang bertandatangan di bawah ini:

Nama : Aria Nanda Herdiawan
NIM : 1910511083
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 ilmiah saya yang berjudul:

PERBANDINGAN ALGORITMA K-NEAREST NEIGHBOR, RANDOM FOREST, DAN LOGISTIC REGRESSION TERHADAP KATEGORI REVIEW PADA MARKETPLACE XYZ

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 kata (basis data), merawat dan mempublikasikan Skripsi saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta. Demikian pernyataan ini saya buat dengan sebenarnya.

Jakarta, 19 Januari 2024

Yang menyatakan,



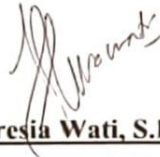
(Aria Nanda Herdiawan)

LEMBAR PENGESAHAN

Tugas akhir ini diajukan oleh:

Nama : Aria Nanda Herdiawan
NIM : 1910511083
Program Studi : S1 Informatika
Judul Tugas Akhir : Perbandingan Algoritma K-Nearest Neighbor,
Random Forest, dan Logistic Regression Terhadap
Kategori Review pada Marketplace XYZ

Telah berhasil dipertahankan dihadapan 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.



(Theresia Wati, S.Kom., MTL)

Penguji I



(Anita Muliawati, S.Kom., MTL)

Penguji II



(Jayanta, S. Kom., M. Si)

Pembimbing



(Dr. Krisnanik S.Kom., MM.)

Plt Dekan



(Dr. Widya Cholil, M.I.T.)

Kepala Program Studi

Ditetapkan di : Jakarta

Tanggal Ujian : 20 Desember 2023

Perbandingan Algoritma K-Nearest Neighbor, Random Forest, dan Logistic Regression Terhadap Kategori Review pada Marketplace XYZ

Aria Nanda Herdiawan

ABSTRAK

Dengan perkembangan teknologi di dunia yang sangat cepat, dapat mempermudah dalam melakukan segala hal. Salah satunya adalah berbelanja secara *online*. Salah satu platform yang dapat digunakan untuk berbelanja *online* adalah *marketplace*. Banyak *review* yang diberikan pembeli ketika selesai melakukan transaksi pada *marketplace*. Sehingga dibutuhkan klasifikasi untuk membedakan kategori *review* pembeli, selain itu dapat diketahui perbedaan kinerja dari algoritma klasifikasi *K-nearest neighbor*, *random forest*, dan *logistic regression*. Dengan menggunakan data *review* sebanyak 1000 data yang dibagi menjadi kelas produk 618 data, pengiriman 214 data, dan pelayanan 168 data. Data dilakukan *preprocessing* dengan tahap *case folding*, *cleansing*, normalisasi, tokenisasi, *filtering*, dan *stemming*. Kemudian dihitung bobot kata pada data dengan *tf-idf*. Data akan dibagi menjadi *data train* dan *data test*. Setelah itu dilakukan klasifikasi dengan *K-nearest neighbor*, *random forest*, dan *logistic regression*. Yang menghasilkan akurasi tertinggi dari *K-nearest neighbor* 76%, *random forest* 81%, dan *logistic regression* 80%.

Kata Kunci: *Review, Marketplace, Klasifikasi Teks, K-Nearest Neighbor, Random Forest, Logistic Regression*

ABSTRACT

With the rapid development of technology in the world, it can make doing everything easier. One of them is shopping online. One platform that can be used for online shopping is a marketplace. Many reviews are given by buyers when completing transactions in marketplace. So classification is needed to distinguish categories of buyer reviews, besides that we can know the difference performance of the K-nearest neighbor, random forest, and logistic regression classification algorithms. By using review data of 1000 data which is divided into product classes 618 data, delivery 214 data, and service 168 data. The data was preprocessed using case folding, cleansing, normalization, tokenization, filtering, and stemming stages. Then the word weights in the data are calculated using tf-idf. The data will be divided into train data and test data. After that, classification was carried out with K-nearest neighbor, random forest, and logistic regression. Which produces the highest accuracy of K-nearest neighbor 76%, random forest 81%, and logistic regression 80%.

Key Word: Review, Marketplace, Text Classification, K-Nearest Neighbor, Random Forest, Logistic Regression

KATA PENGHANTAR

Puji dan syukur selalu tercurahkan kepada Allah SWT atas nikmat dan karunia yang diberikan-Nya, sehingga penulis dapat menyelesaikan skripsi dengan judul “Perbandingan Algoritma K-Nearest Neighbor, Random Forest, dan Logistic Regression Terhadap Kategori Review pada Marketplace XYZ”.

Dalam penulisan skripsi ini tidak mungkin selesai tanpa adanya dukungan, bantuan, saran, serta bimbingan dari berbagai pihak. Oleh karena itu penulis ingin memberikan terimakasih pada berbagai pihak yang telah membantu selama proses penulisan skripsi ini, yaitu kepada:

1. Bapak Andy Sosiawan dan Ibu Heti Herliati selaku orang tua penulis yang selalu memberikan semangat serta doa selama pengerjaan skripsi ini.
2. Bapak Jayanta, S. Kom., M. Si. selaku dosen pembimbing skripsi, yang telah membantu dan memberikan saran selama proses penulisan skripsi ini.
3. Bapak Hamonangan Kinantan Prabu, M.T. selaku dosen pembimbing akademik penulis.
4. Ibu Dr. Widya Cholil, M.I.T selaku Kaprodi Informatika Universitas Pembangunan Nasional Veteran Jakarta.
5. Luthfiah Putri Ayuningtyas dan Ayuni Fitria selaku teman penulis yang telah membantu dalam proses pelabelan data.
6. Serta keluarga, kerabat, dan rekan penulis yang telah membantu dan memberikan dukungan selama penulisan skripsi ini.

Tangerang Selatan, 19 Januari 2023

Penulis,



Aria Nanda Herdiawan

DAFTAR ISI

| | |
|-----------------------------------|------|
| PERNYATAAN ORISINALITAS | ii |
| LEMBAR PERSETUJUAN PUBLIKASI..... | iii |
| LEMBAR PENGESAHAN | iv |
| ABSTRAK | v |
| ABSTRACT | vi |
| KATA PENGHANTAR | vii |
| DAFTAR ISI..... | viii |
| DAFTAR TABEL..... | xi |
| DAFTAR GAMBAR | xii |
| DAFTAR SIMBOL..... | xiii |
| BAB 1 PENDAHULUAN | 1 |
| 1.1 Latar Belakang | 1 |
| 1.2 Rumusan Masalah | 3 |
| 1.3 Tujuan Penelitian | 3 |
| 1.4 Manfaat Penelitian | 4 |
| 1.5 Ruang Lingkup..... | 4 |
| 1.6 Luaran Penelitian | 4 |
| 1.7 Sistematika Penulisan | 4 |
| BAB 2 TINJAUAN PUSTAKA | 6 |
| 2.1 Marketplace..... | 6 |
| 2.2 Review | 6 |
| 2.3 Text Mining..... | 6 |
| 2.4 Klasifikasi Teks..... | 7 |
| 2.5 <i>Pre Processing</i> | 7 |
| 2.5.1 <i>Case Folding</i> | 7 |
| 2.5.2 <i>Cleansing</i> | 8 |
| 2.5.3 Normalisasi | 8 |
| 2.5.4 Tokenisasi | 8 |
| 2.5.5 <i>Filtering</i> | 8 |

| | |
|---|----|
| 2.5.6 <i>Stemming</i> | 9 |
| 2.6 Tf-Idf | 9 |
| 2.7 <i>K-Nearest Neighbor</i> | 10 |
| 2.8 <i>Random Forest</i> | 10 |
| 2.9 <i>Logistic Regression</i> | 11 |
| 2.10 <i>Confusion Matrix</i> | 12 |
| 2.11 Penelitian Terkait | 13 |
| BAB 3 METODE PENELITIAN | 16 |
| 3.1 Identifikasi Masalah | 17 |
| 3.2 Studi Literatur | 17 |
| 3.3 Akuisisi Data..... | 17 |
| 3.4 Pelabelan Data..... | 18 |
| 3.5 <i>Preprocessing</i> | 18 |
| 3.5.1 <i>Case Folding</i> | 19 |
| 3.5.2 <i>Cleansing</i> | 19 |
| 3.5.3 Normalisasi | 20 |
| 3.5.4 Tokenisasi | 20 |
| 3.5.5 <i>Filtering</i> | 21 |
| 3.5.6 <i>Stemming</i> | 21 |
| 3.6 Pembobotan..... | 22 |
| 3.7 Pembagian Data | 22 |
| 3.8 Klasifikasi | 23 |
| 3.9 Evaluasi..... | 23 |
| 3.10 Perbandingan Kinerja <i>K-Nearest Neighbor</i> , <i>Random Forest</i> , dan <i>Logistic Regression</i> | 24 |
| 3.11 Alat Bantu Penelitian | 24 |
| 3.11.1 Perangkat Keras | 24 |
| 3.11.2 <i>Perangkat Lunak</i> | 24 |
| BAB 4 HASIL DAN PEMBAHASAN | 25 |
| 4.1 Data | 25 |

| | |
|---|-----|
| 4.2 Pelabelan Data..... | 27 |
| 4.3 <i>Preprocessing</i> | 28 |
| 4.3.1 <i>Case Folding</i> | 28 |
| 4.3.2 <i>Cleansing</i> | 29 |
| 4.3.3 Normalisasi | 31 |
| 4.3.4 Tokenisasi | 32 |
| 4.3.5 <i>Filtering</i> | 33 |
| 4.3.6 <i>Stemming</i> | 35 |
| 4.4 Pembobotan Kata | 36 |
| 4.5 Pembagian Data | 38 |
| 4.6 <i>K-Nearest Neighbor</i> | 39 |
| 4.7 <i>Random Forest</i> | 41 |
| 4.8 <i>Logistic Regression</i> | 42 |
| 4.9 Evaluasi..... | 44 |
| 4.9.1 Evaluasi <i>K-Nearest Neighbor</i> | 44 |
| 4.9.2 Evaluasi <i>Random Forest</i> | 48 |
| 4.9.3 Evaluasi <i>Logistic Regression</i> | 50 |
| 4.10 Perbandingan Kinerja <i>K-Nearest Neighbor</i> , <i>Random Forest</i> , dan <i>Logistic Regression</i> | 53 |
| 4.11 Implementasi pada <i>Marketplace</i> | 55 |
| BAB 5 PENUTUP | 61 |
| 5.1 Kesimpulan | 61 |
| 5.2 Saran..... | 62 |
| DAFTAR PUSTAKA | 63 |
| LAMPIRAN..... | 67 |
| Lampiran 1: Dataset | 67 |
| Lampiran 2: Turnitin | 139 |




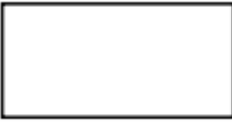
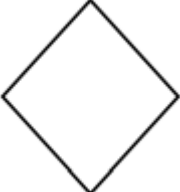
DAFTAR TABEL

| | |
|--|----|
| Tabel 1.1 Rata-rata pengunjung web marketplace bulanan | 2 |
| Tabel 2.1 Confusion Matrix | 13 |
| Tabel 3.1 Confusion Matrix 3 kelas | 23 |
| Tabel 4.1 Hasil case folding | 28 |
| Tabel 4.2 Hasil cleansing | 30 |
| Tabel 4.3 Hasil normalisasi | 31 |
| Tabel 4.4 Hasil tokenisasi | 32 |
| Tabel 4.5 Hasil filtering | 34 |
| Tabel 4.6 Hasil stemming | 35 |
| Tabel 4.7 Tf-Idf | 36 |
| Tabel 4.8 Pembagian data | 38 |
| Tabel 4.9 Tf-Idf setiap kata pada D3 | 39 |
| Tabel 4.10 klasifikasi k-nearest neighbor | 41 |
| Tabel 4.11 klasifikasi random forest | 42 |
| Tabel 4.12 Klasifikasi Logistic Regression | 44 |
| Tabel 4.13 Kinerja K-nearest neighbor data train 80% : data test 20% | 47 |
| Tabel 4.14 Kinerja K-nearest neighbor data train 50% : data test 50% | 47 |
| Tabel 4.15 Kinerja K-nearest neighbor data train 20% : data test 80% | 47 |
| Tabel 4.16 Kinerja random forest | 50 |
| Tabel 4.17 Kinerja logistic regression | 53 |

DAFTAR GAMBAR

| | |
|--|----|
| Gambar 2.1 Random Forest (Zivkovic, 2022) | 11 |
| Gambar 3.1 Kerangka berpikir..... | 16 |
| Gambar 3.2 Preprocessing | 19 |
| Gambar 3.3 Proses case folding | 19 |
| Gambar 3.4 Proses cleansing | 19 |
| Gambar 3.5 Proses Normalisasi | 20 |
| Gambar 3.6 Proses tokenisasi | 21 |
| Gambar 3.7 Proses filtering..... | 21 |
| Gambar 3.8 Proses stemming..... | 22 |
| Gambar 4.1 Data hasil Scraping pertama..... | 25 |
| Gambar 4.2 Data hasil Scraping kedua | 26 |
| Gambar 4.3 Data hasil Scraping ketiga..... | 26 |
| Gambar 4.4 Data terkumpul | 27 |
| Gambar 4.5 Label data | 28 |
| Gambar 4.6 Nilai intercept dari masing-masing kelas | 43 |
| Gambar 4. 7 Nilai koefisien dari masing-masing kelas | 43 |
| Gambar 4.8 Confusion Matrix K-NN | 45 |
| Gambar 4.9 confusion matrix random forest | 48 |
| Gambar 4.10 confusion matrix logistic regression | 51 |
| Gambar 4.11 Perbandingan kinerja algoritma | 54 |
| Gambar 4.12 Tampilan review pada marketplace pada umumnya | 56 |
| Gambar 4.13 Tampilan fitur review dengan tambahan filter kategori..... | 57 |
| Gambar 4.14 Tampilan dari pilihan fitur filter kategori | 58 |
| Gambar 4.15 Tampilan fitur filter ketika diklik kategori produk | 59 |
| Gambar 4.16 Tampilan review yang hanya berisi kategori produk | 60 |

DAFTAR SIMBOL

| GAMBAR | NAMA | KETERANGAN |
|---|-------------------|--|
|  | <i>Terminator</i> | Melambangkan awal dimulai dan akhir dari proses |
|  | <i>Flow</i> | Melambangkan penghubung antar simbol untuk alur proses |
|  | <i>Data</i> | Melambangkan data <i>input</i> atau <i>output</i> dalam proses |
|  | <i>Process</i> | Melambangkan proses yang dijalankan |
|  | <i>Decision</i> | Melambangkan percabangan proses yang memiliki putusan yang berbeda |