

**IMPLEMENTASI METODE FITUR EKSTRAKSI DAN
ALGORITMA SUPERVISED LEARNING DALAM
MENGIDENTIFIKASI JENIS KELAMIN
BERDASARKAN NAMA**

SKRIPSI SARJANA

Oleh

Muhammad Rizkiansyah
183112706450091



**PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI KOMUNIKASI DAN
INFORMATIKA
UNIVERSITAS NASIONAL
2022**

**IMPLEMENTASI METODE FITUR EKSTRAKSI DAN
ALGORITMA SUPERVISED LEARNING DALAM
MENGIDENTIFIKASI JENIS KELAMIN
BERDASARKAN NAMA**

SKRIPSI SARJANA

Karya ilmiah sebagai salah satu syarat untuk memperoleh gelar
Sarjana Teknik Teknologi Informatika dari Fakultas Teknologi Komunikasi dan
Informatika

Oleh

Muhammad Rizkiansyah
183112706450091



**PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI KOMUNIKASI DAN
INFORMATIKA
UNIVERSITAS NASIONAL
2022**

**HALAMAN PENGESAHAN
TUGAS SARJANA**

**IMPLEMENTASI METODE FITUR EKSTRAKSI DAN
ALGORITMA SUPERVISED LEARNING DALAM
MENGIDENTIFIKASI JENIS KELAMIN
BERDASARKAN NAMA**



Pembimbing I

Pembimbing II

(Ratih Titi Komala Sari, ST, MM, MMSI)
NID. 0103150850

(Albaar Rubhasy, S.Si, MTI)
NID. 050020069

KATA PENGANTAR

Puji dan syukur penulis panjatkan kepada Allah SWT Tuhan Yang Maha Esa yang telah memberikan rahmat dan karunia sehingga penulis dapat menyelesaikan skripsi dengan judul **“Implementasi Metode Fitur Ekstraksi dan Algoritma Supervised Learning dalam Mengidentifikasi Jenis Kelamin berdasarkan Nama”** sebagai salah satu syarat kelulusan Program Studi Sarjana Informatika Fakultas Teknologi Komunikasi dan Informatika.

Penelitian dan penulisan skripsi ini tidak terlepas dari bantuan berbagai pihak, oleh karena itu penulis menyampaikan banyak terima kasih terutama kepada dosen pembimbing Tugas Akhir, ibu Ratih Titi Komala Sari, ST, MM, MMSI dan bapak Albaar Rubhasy, S.Si, MTI yang telah meluangkan banyak waktu, tenaga, pikiran, bimbingan, arahan, motivasi serta memaklumi segala kekurangan penulis selama penelitian tugas akhir dan penyusunan skripsi. Penulis juga mengucapkan banyak terima kasih kepada:

1. Bapak dan Mamah selaku orangtua penulis yang telah banyak memberi dukungan dan doa yang tak terhitung.
2. Seluruh dosen pengajar di program studi Informatika FTKI maupun dosen di program studi lain, yang memberikan banyak pengetahuan dan ilmu.
3. Teman-teman seangkatan dan sehimpuan yang telah membantu dan memberikan dukungan.
4. Kaka dan sahabat-sahabat yang telah memberikan banyak dukungan secara mental maupun fisik.

Akhir kata, semoga Allah SWT Tuhan Yang Maha Esa membalas kebaikan dan bantuan yang telah diberikan dengan hal yang lebih baik. Penulis mengharapkan kritik dan saran yang bersifat membangun dan semoga skripsi ini dapat memberikan manfaat dibidang teknologi informatika.

Tangerang Selatan, 17 Desember 2022


Muhammad Rizkiansyah

**HALAMAN PERNYATAAN PERSETUJUAN
PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN
AKADEMIS**

Sebagai sivitas akademik Program Studi Teknik Informatika, Fakultas Teknologi Komunikasi dan Informatika, saya yang bertanda tangan di bawah ini:

Nama : Muhammad Rizkiansyah

NIM : 183112706450091

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Fakultas Teknologi Komunikasi dan Informatika, Hak Bebas Royalti Noneksklusif (*Non-exclusive Royalti Free Right*) atas karya ilmiah saya yang berjudul:

**IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA
SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS
KELAMIN BERDASARKAN NAMA**

Beserta perangkat yang ada (jika diperlukan). Dengan Hak ini Fakultas Teknologi Komunikasi dan Informatika berhak menyimpan, mengalihmedia/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 : Tangerang Selatan

Pada tanggal : 17 Desember 2022

Yang menyatakan



(Muhammad Rizkiansyah)

HALAMAN PENGESAHAN

TUGAS AKHIR

IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA
SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS
KELAMIN BERDASARKAN NAMA



Muhammad Rizkiansyah
183112706450091

Dosen Pembimbing 1

A handwritten signature in black ink, appearing to be 'Ratih Titi Komalasari'.

(Ratih Titi Komalasari, S. T, MM., MMSI)

Dosen Pembimbing 2

A handwritten signature in black ink, appearing to be 'Albaar Rubhasy'.

(Albaar Rubhasy, S.Si, MTI)

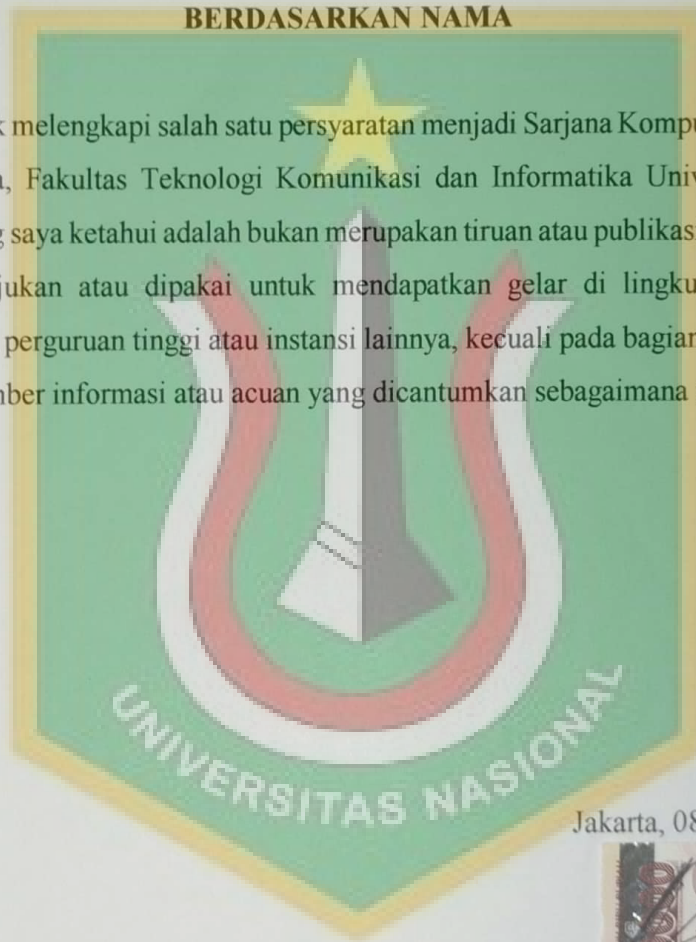
PERNYATAAN KEASLIAN TUGAS AKHIR

Saya menyatakan dengan sesungguhnya bahwa Tugas Akhir dengan judul :

IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS KELAMIN

BERDASARKAN NAMA

Yang dibuat untuk melengkapi salah satu persyaratan menjadi Sarjana Komputer pada Program Studi Informatika, Fakultas Teknologi Komunikasi dan Informatika Universitas Nasional, sebagaimana yang saya ketahui adalah bukan merupakan tiruan atau publikasi dari Tugas Akhir yang pernah diajukan atau dipakai untuk mendapatkan gelar di lingkungan Universitas Nasional maupun perguruan tinggi atau instansi lainnya, kecuali pada bagian – bagian tertentu yang menjadi sumber informasi atau acuan yang dicantumkan sebagaimana mestinya.



Jakarta, 08 Maret 2023



Muhammad Rizkiansyah

183122706450091

LEMBAR PERSETUJUAN TUGAS AKHIR

Tugas Akhir dengan judul :

IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS KELAMIN

BERDASARKAN NAMA

Dibuat untuk melengkapi salah satu persyaratan menjadi Sarjana Komputer pada Program Studi Informatika, Fakultas Teknologi Komunikasi dan Informatika Universitas Nasional. Tugas Akhir ini diujikan pada Sidang Akhir Semester Ganjil 2022-2023 pada tanggal 22 Februari Tahun 2023



Dosen Pembimbing 1

Ratih Titi Komalasari, S. T.

MM., MMSI

NID. 0103150850

Ketua Program Studi

A circular stamp with the text 'UNIVERSITAS NASIONAL' around the perimeter. Inside the stamp is a handwritten signature.

Ratih Titi Komalasari, S. T.

MM., MMSI

NID. 0103150850

LEMBAR PERSETUJUAN JUDUL YANG TIDAK ATAU YANG DIREVISI

Nama : Muhammad Rizkiansyah
NPM : 183112706450091
Fakultas/Akademi : Fakultas Teknologi Komunikasi dan Informatika
Program Studi : Informatika
Tanggal Sidang : 22 Februari 2023

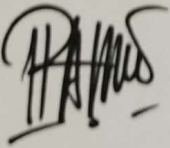

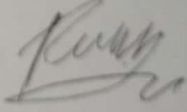
JUDUL DALAM BAHASA INDONESIA :

IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS KELAMIN BERDASARKAN NAMA

JUDUL DALAM BAHASA INGGRIS :

IMPLEMENTATION OF FEATURE EXTRACTION METHOD AND SUPERVISED LEARNING ALGORITHM IN IDENTIFYING GENDER BASED ON NAME

TANDA TANGAN DAN TANGGAL

| Pembimbing 1 | Ka. Prodi | Mahasiswa |
|---|---|---|
| TGL : | TGL : | TGL : 13 - 03 - 2023 |
|  |  |  |

LEMBAR PERSETUJUAN JUDUL YANG TIDAK ATAU YANG DIREVISI

Nama : Muhammad Rizkiansyah
NPM : 183112706450091
Fakultas/Akademi : Fakultas Teknologi Komunikasi dan Informatika
Program Studi : Informatika
Tanggal Sidang : 22 Februari 2023

JUDUL DALAM BAHASA INDONESIA :

IMPLEMENTASI METODE FITUR EKSTRAKSI DAN ALGORITMA SUPERVISED LEARNING DALAM MENGIDENTIFIKASI JENIS KELAMIN BERDASARKAN NAMA

JUDUL DALAM BAHASA INGGRIS :

IMPLEMENTATION OF FEATURE EXTRACTION METHOD AND SUPERVISED LEARNING ALGORITHM IN IDENTIFYING GENDER BASED ON NAME

TANDA TANGAN DAN TANGGAL

| Pembimbing 2 | Ka. Prodi | Mahasiswa |
|---|--|---|
| TGL : 10-03-2023 | TGL : | TGL : 10-03-2023 |
|  |   |  |

ABSTRAK

Pemegang data memiliki dokumen terkait data personal akan tetapi tidak adanya variabel gender dalam dokumen tersebut. Dalam hal ini, memungkinkan pemilik data mengklasifikasikan variabel gender berdasarkan nama yang ada secara manual. Mengklasifikasikan dokumen secara manual dinilai tidak seefisien seperti dulu, karena jumlah data yang semakin meningkat. Algoritma supervised learning yang ada pada machine learning dapat berperan, sebagai cara alternatif dalam pengklasifikasian dokumen. Penelitian ini dilakukan untuk mencari model yang terbaik dalam mengidentifikasi jenis kelamin berdasarkan nama, yang diimplementasikan ke dalam aplikasi prediksi berbasis web, dengan menggunakan fitur ekstraksi countvectorizer dan pemanfaatan n-gram serta membandingkan kedua algoritma supervised learning yaitu logistic regression dan multinomial naive bayes. Hasil dari penelitian ini didapatkan model terbaik yaitu model logistic regression pada rentan 2-12 gram dengan split data 80:20 yang memiliki tingkat akurasi 94.76%, dan berdasarkan uji validasi model menggunakan confusion matrix logistic regression memperoleh 0.95 f1-score pada semua labeling, beserta dari hasil uji prediksi yang mendapatkan kesalahan prediksi tersedikit yaitu 28. Maka dari itu model logistic regression dengan split data 80:20 yang akan diterapkan ke dalam aplikasi prediksi berbasis web.

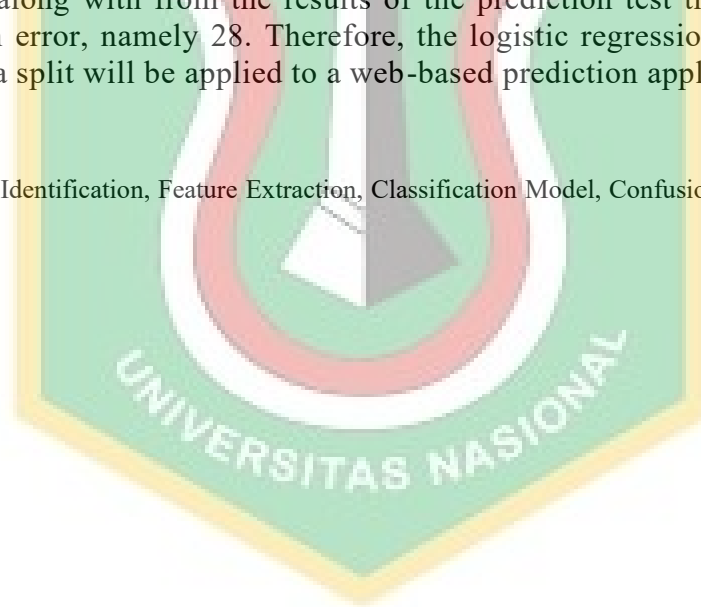
Kata Kunci: Identifikasi, Fitur Ekstraksi, Model Klasifikasi, Confusion Matrix, Aplikasi Prediksi.



ABSTRACT

The data owner has documents containing personal data, but there is no mention of gender in those records. In this instance, it enables the data owner to manually categorize gender variables based on names that already exist. The growing volume of data makes manual document classification less effective than it formerly was. As a substitute method of classifying documents, the supervised learning technique used in machine learning may be useful. By comparing the two supervised learning algorithms, logistic regression and multinomial naive Bayes, and using the countvectorizer extraction feature and n-grams, this research sought to identify the most effective model for classifying gender based on name. This model was then implemented into a web-based prediction application. The results of this study obtained the best model, namely the logistic regression model at a susceptibility of 2-12 grams with a data split of 80:20 which had an accuracy rate of 94.76%, and based on the model validation test using a confusion matrix logistic regression obtained 0.95 f1-score on all labeling, along with from the results of the prediction test that gets the least prediction error, namely 28. Therefore, the logistic regression model with an 80:20 data split will be applied to a web-based prediction application.

Keywords: Identification, Feature Extraction, Classification Model, Confusion Matrix, Prediction Application.



DAFTAR ISI

| | |
|--|------------|
| HALAMAN PERNYATAAN ORISINALITAS | i |
| HALAMAN PENGESAHAN..... | ii |
| TUGAS SARJANA | ii |
| KATA PENGANTAR..... | iii |
| HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS | iv |
| ABSTRAK..... | v |
| ABSTRACT..... | vi |
| DAFTAR ISI..... | vii |
| DAFTAR TABEL | x |
| DAFTAR GAMBAR..... | xi |
| BAB I..... | 1 |
| PENDAHULUAN..... | 1 |
| 1.1. Latar Belakang | 1 |
| 1.2. Perumusan Masalah..... | 2 |
| 1.3. Tujuan Penelitian | 2 |
| 1.4. Manfaat Penelitian | 3 |
| 1.5. Batasan Masalah..... | 3 |
| BAB II | 4 |
| TINJAUAN PUSTAKA..... | 4 |
| 2.1. Landasan Teori..... | 4 |
| 2.1.1. Data terpilih | 4 |
| 2.1.2. Identifikasi jenis kelamin berdasarkan nama | 4 |
| 2.1.3. Data Preprocessing..... | 4 |
| 2.1.4. Fitur Ekstraksi | 5 |

| | | |
|--------------------------|---|----|
| 2.1.5. | Supervised learning..... | 7 |
| 2.1.6. | Logistic Regression | 8 |
| 2.1.7. | Multinomial Naive Bayes..... | 8 |
| 2.1.8. | Evaluasi | 8 |
| 2.2. | Studi Literatur | 10 |
| 2.2.1. | Matriks Penelitian | 13 |
| BAB III | | 15 |
| METODE PENELITIAN | | 15 |
| 3.1. | Teknik Pengumpulan Data..... | 15 |
| 3.2. | Sumber Data | 15 |
| 3.3. | Perancangan Sistem | 16 |
| 3.3.1. | Preprocessing data | 16 |
| 3.3.2. | Split Data | 17 |
| 3.3.3. | Fitur Ekstraksi | 17 |
| 3.3.4. | Modeling | 17 |
| 3.3.5. | Evaluasi Model | 18 |
| 3.4. | Pengujian Sistem | 18 |
| 3.5. | Implementasi Sistem | 18 |
| BAB IV | | 19 |
| HASIL DAN DISKUSI | | 19 |
| 4.1. | Pembahasan Data dan Preprocessing..... | 19 |
| 4.1.1. | Data | 19 |
| 4.1.2. | Preprocessing..... | 20 |
| 4.2. | Pembahasan Split Data | 22 |
| 4.3. | Pembahasan Fitur dan Modeling..... | 23 |
| 4.3.1. | Fitur | 23 |
| 4.3.2. | Modeling | 24 |
| 4.4. | Evaluasi Model | 25 |
| 4.5. | Pengujian Sistem | 28 |
| 4.6. | Implementasi Model..... | 29 |
| 4.6.1. | Tampilan implementasi model pada aplikasi web | 29 |
| 4.6.2. | Uji prediksi model pada aplikasi berbasis web | 31 |
| BAB V | | 32 |

| | |
|-----------------------------------|-----------|
| KESIMPULAN DAN SARAN | 32 |
| 5.1. Kesimpulan | 32 |
| 5.2. Saran..... | 32 |
| DAFTAR PUSTAKA | 34 |



DAFTAR TABEL

| | |
|--|----|
| Tabel 2.1 Tahap vektorisasi bag-of-word..... | 6 |
| Tabel 2.2 Tabel yang dihasilkan confusion matrix | 9 |
| Tabel 2.3 Tabel perhitungan confusion matrix..... | 9 |
| Tabel 2.4 Matrix penelitian | 13 |
| Tabel 3.3 Parameter yang digunakan | 17 |
| Tabel 3.4 Representasi hasil | 18 |
| Tabel 4.1 Tampilan dataset 1..... | 19 |
| Tabel 4.2 Tampilan dataset 2..... | 19 |
| Tabel 4.3 Hasil transformasi data pada dataset 1 | 20 |
| Tabel 4.4 Hasil transformasi data pada dataset 2 | 20 |
| Tabel 4.5 Hasil integration data | 21 |
| Tabel 4.6 Sebelum preprocessing data..... | 21 |
| Tabel 4.7 Sesudah preprocessing data..... | 22 |
| Tabel 4.8 Hasil split data..... | 22 |
| Tabel 4.9 Hasil term frekuensi beserta pemanfaatan n-gram | 23 |
| Tabel 4.10 Hasil nilai rata-rata 2 sampai 12 gram (LR)..... | 24 |
| Tabel 4.11 Hasil nilai rata-rata 2 sampai 12 gram (MNB)..... | 24 |
| Tabel 4.12 Data nilai hasil confusion matrix setiap model | 26 |
| Tabel 4.13 Nilai validasi dari setiap model pada pembagian data 70:30 | 26 |
| Tabel 4.14 Nilai validasi dari setiap model pada pembagian data 80:20 | 27 |
| Tabel 4.15 Hasil uji prediksi | 28 |
| Tabel 4.16 Persentase uji prediksi | 28 |

DAFTAR GAMBAR

| | |
|--|----|
| Gambar 3.2 Rancangan sistem penelitian..... | 16 |
| Gambar 4.1 Heatmap confusion matrix lr (a) 70:30 (b) 80:20..... | 25 |
| Gambar 4.2 Heatmap confusion matrix mnb (a) 70:20 (b) 80:20 | 25 |
| Gambar 4.3 Tampilan home | 29 |
| Gambar 4.4 Tampilan halaman prediksi dokumen..... | 30 |
| Gambar 4.5 Tampilan halaman prediksi satuan | 30 |
| Gambar 4.6 Uji prediksi pada web | 31 |

