

***MULTINOMIAL NAÏVE BAYES* UNTUK DETEKSI *PHISHING*
URL DENGAN PERBANDINGAN REPRESENTASI FITUR
BAG OF WORDS DAN *TF-IDF***

SKRIPSI

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi S1 Teknik Komputer



disusun oleh

GEOREL JEFERSON FRANSISKUS BONAI

22.83.0833

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2026

***MULTINOMIAL NAÏVE BAYES UNTUK DETEKSI PHISHING
URL DENGAN PERBANDINGAN REPRESENTASI FITUR
BAG OF WORDS DAN TF-IDF***

SKRIPSI

untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi S1 Teknik Komputer



disusun oleh

GEOREL JEFERSON FRANSISKUS BONAI

22.83.0833

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2026

HALAMAN PERSETUJUAN

SKRIPSI

***MULTINOMIAL NAÏVE BAYES* UNTUK DETEKSI *PHISHING*
URL DENGAN PERBANDINGAN REPRESENTASI FITUR
BAG OF WORDS DAN *TF-IDF***

yang disusun dan diajukan oleh

Georel Jeferson Fransiskus Bonai

22.83.0833

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 23 Januari 2026

Dosen Pembimbing,



Tonny Hidayat, S.Kom., M.Kom., Ph.D

NIK. 190302182

HALAMAN PENGESAHAN

SKRIPSI

**MULTINOMIAL NAÏVE BAYES UNTUK DETEKSI PHISHING
URL DENGAN PERBANDINGAN REPRESENTASI FITUR
BAG OF WORDS DAN TF-IDF**

yang disusun dan diajukan oleh

Georel Jeferson Fransiskus Bonai

22.83.0833

Telah dipertahankan di depan Dewan Penguji
pada tanggal 23 Januari 2026

Susunan Dewan Penguji

Nama Penguji

Melwin Syafrizal, S.Kom., M.Eng., Ph.D
NIK. 190302105

Senie Destya, M.Kom
NIK. 190302312

Tonny Hidayat, S.Kom., M.Kom., Ph.D
NIK. 190302182

Tanda Tangan



Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 23 Januari 2026

DEKAN FAKULTAS ILMU KOMPUTER



Prof. Dr. Kusrini, M.Kom.
NIK. 190302106

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

Nama mahasiswa : Georel Jeferson Fransiskus Bonai
NIM : 22.83.0833

Menyatakan bahwa Skripsi dengan judul berikut:

MULTINOMIAL NAÏVE BAYES UNTUK DETEKSI PHISHING URL DENGAN PERBANDINGAN REPRESENTASI FITUR BAG OF WORDS DAN TF-IDF

Dosen Pembimbing : Tonny Hidayat, S.Kom., M.Kom., Ph.D

1. Karya tulis ini adalah benar-benar ASLI dan BELUM PERNAH diajukan untuk mendapatkan gelar akademik, baik di Universitas AMIKOM Yogyakarta maupun di Perguruan Tinggi lainnya.
2. Karya tulis ini merupakan gagasan, rumusan dan penelitian SAYA sendiri, tanpa bantuan pihak lain kecuali arahan dari Dosen Pembimbing.
3. Dalam karya tulis ini tidak terdapat karya atau pendapat orang lain, kecuali secara tertulis dengan jelas dicantumkan sebagai acuan dalam naskah dengan disebutkan nama pengarang dan disebutkan dalam Daftar Pustaka pada karya tulis ini.
4. Perangkat lunak yang digunakan dalam penelitian ini sepenuhnya menjadi tanggung jawab SAYA, bukan tanggung jawab Universitas AMIKOM Yogyakarta.
5. Pernyataan ini SAYA buat dengan sesungguhnya, apabila di kemudian hari terdapat penyimpangan dan ketidakbenaran dalam pernyataan ini, maka SAYA bersedia menerima SANKSI AKADEMIK dengan pencabutan gelar yang sudah diperoleh, serta sanksi lainnya sesuai dengan norma yang berlaku di Perguruan Tinggi.

Yogyakarta, 23 Januari 2026

Yang Menyatakan,



Georel Jeferson Fransiskus Bonai

HALAMAN PERSEMBAHAN

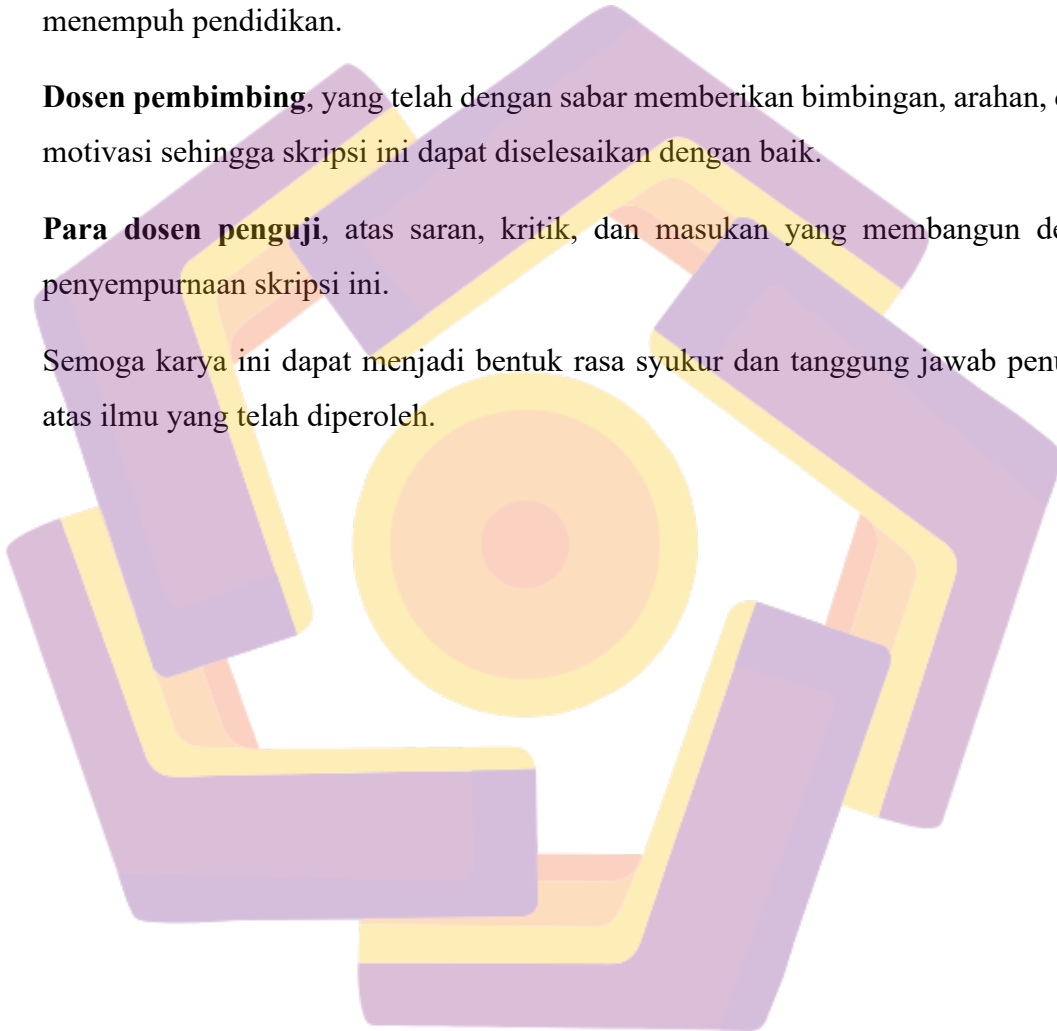
Skripsi ini penulis persembahkan kepada:

Orang tua tercinta dan wali serta keluarga, yang senantiasa memberikan doa, kasih sayang, dukungan moral, serta pengorbanan yang tak ternilai selama penulis menempuh pendidikan.

Dosen pembimbing, yang telah dengan sabar memberikan bimbingan, arahan, dan motivasi sehingga skripsi ini dapat diselesaikan dengan baik.

Para dosen penguji, atas saran, kritik, dan masukan yang membangun demi penyempurnaan skripsi ini.

Semoga karya ini dapat menjadi bentuk rasa syukur dan tanggung jawab penulis atas ilmu yang telah diperoleh.



KATA PENGANTAR

Puji dan syukur penulis panjatkan ke hadirat Tuhan Yesus Kristus, karena atas rahmat dan karunia-Nya penulis dapat menyelesaikan skripsi ini dengan judul ”*Multinomial Naïve Bayes* untuk Deteksi *Phishing* URL dengan Perbandingan Representasi Fitur *Bag of Words* dan TF-IDF” Skripsi ini disusun sebagai salah satu syarat untuk memperoleh gelar Sarjana pada Program Studi Teknik Komputer, Fakultas Ilmu Komputer, Universitas Amikom Yogyakarta.

Penulis menyadari bahwa dalam proses penyusunan skripsi ini tidak terlepas dari bantuan, bimbingan, dan dukungan dari berbagai pihak. Oleh karena itu, pada kesempatan ini penulis ingin menyampaikan ucapan terima kasih yang sebesar-besarnya kepada:

1. Bapak Tonny Hidayat, S.Kom., M.Kom., Ph.D selaku Dosen Pembimbing, yang telah meluangkan waktu, tenaga, dan pikiran untuk memberikan bimbingan, arahan, serta masukan yang sangat berarti selama proses penyusunan skripsi ini hingga dapat diselesaikan dengan baik.
2. Bapak dan Ibu Dosen Penguji, yang telah memberikan saran, kritik, dan masukan yang membangun dalam proses pendadaran, sehingga skripsi ini dapat disempurnakan.
3. Orang Tua dan Keluarga tercinta, yang telah memberikan dukungan moral, doa, serta bantuan finansial selama penulis menempuh pendidikan hingga tahap penyelesaian skripsi ini.
4. Yoga Pangestu teman saya, yang telah banyak membantu penulis, baik melalui diskusi, pemberian contoh kerangka presentasi, contoh skripsi, naskah, maupun berbagai informasi yang berkaitan dengan isi dan penyusunan skripsi ini.

Yogyakarta, 23 Januari 2026

Penulis

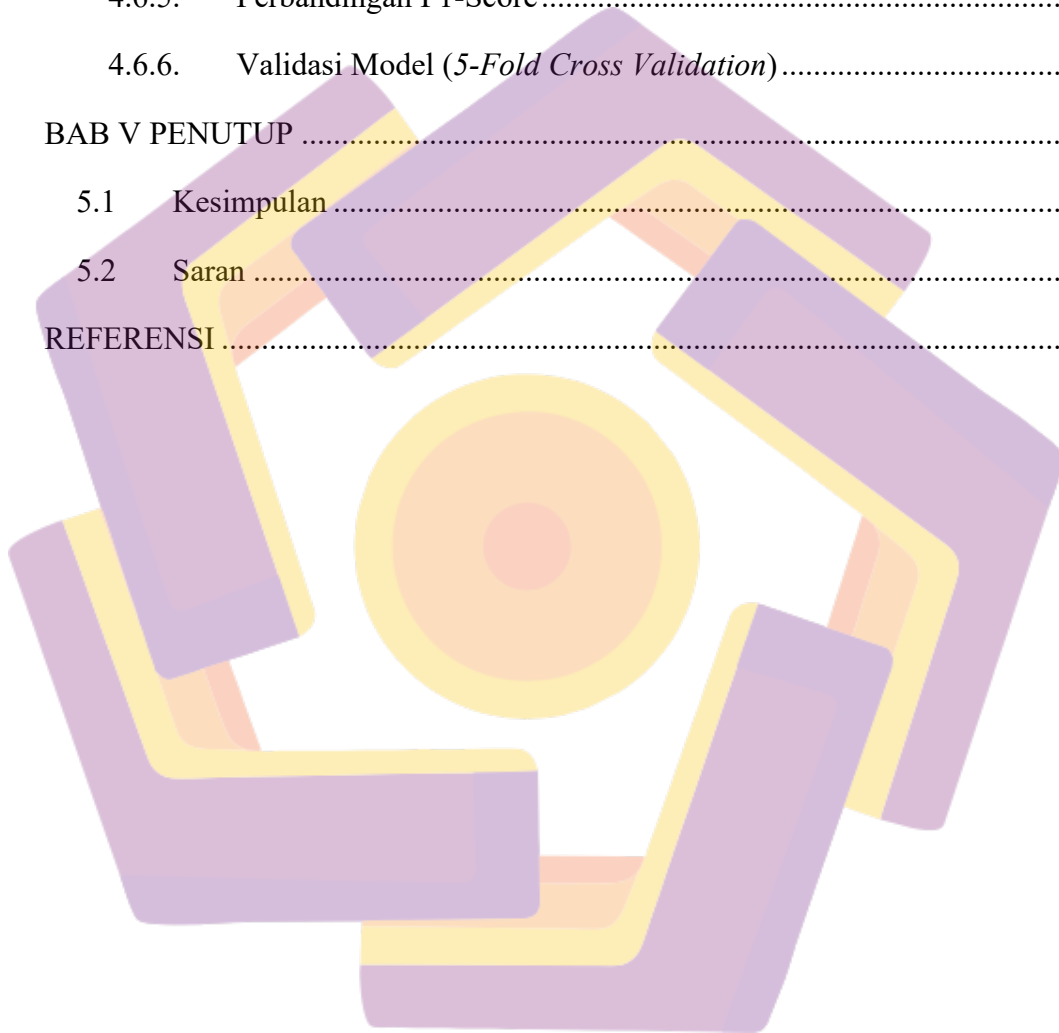
DAFTAR ISI

HALAMAN JUDUL	i
HALAMAN PERSETUJUAN.....	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERNYATAAN KEASLIAN SKRIPSI	iv
HALAMAN PERSEMBAHAN	v
KATA PENGANTAR	vi
DAFTAR ISI.....	vii
DAFTAR TABEL.....	xi
DAFTAR GAMBAR.....	xii
DAFTAR LAMPIRAN.....	xiv
DAFTAR LAMBANG DAN SINGKATAN	xv
DAFTAR ISTILAH.....	xvi
INTISARI	xvii
<i>ABSTRACT</i>	xviii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah.....	3
1.3 Batasan Masalah	3
1.4 Tujuan Penelitian	4
1.5 Manfaat Penelitian	4
1.6 Sistematika Penulisan	4
BAB II TINJAUAN PUSTAKA	6
2.1 Studi Literatur	6

2.1.1.	Keamanan Siber dan Ancaman <i>Phishing</i>	6
2.1.2.	Pendeteksian <i>Phishing</i> Menggunakan <i>Machine Learning</i>	6
2.1.3.	Penerapan <i>Naïve Bayes</i> dalam Klasifikasi <i>Phishing</i>	8
2.1.4.	Pengaruh Representasi Teks terhadap Kinerja Klasifikasi	9
2.1.5.	Perbandingan Metode Tokenisasi	10
2.2	Dasar Teori.....	17
2.2.1.	<i>Machine Learning</i>	17
2.2.2.	Klasifikasi Teks	17
2.2.3.	<i>Phishing</i>	21
2.2.4.	<i>Bag-of-Words</i> (BoW).....	24
2.2.5.	<i>Term Frequency–Inverse Document Frequency</i> (TF-IDF).....	26
2.2.6.	SMOTE	27
2.2.7.	<i>Confusion Matrix</i>	27
2.2.8.	<i>Naïve Bayes Classifier</i>	30
BAB III METODE PENELITIAN		31
3.1	Jenis dan Pendekatan Penelitian	31
3.2	Objek Penelitian.....	32
3.3	Alur Penelitian	33
3.4	Analisis Kebutuhan.....	34
3.4.1.	Kebutuhan Fungsional	34
3.4.2.	Kebutuhan Non-Fungsional	35
3.5.	Rencana Pengujian.....	37
3.5.1.	Skenario Pengujian	37
3.5.2.	Skema Pengujian Model	38
3.5.3.	Metode Evaluasi.....	39

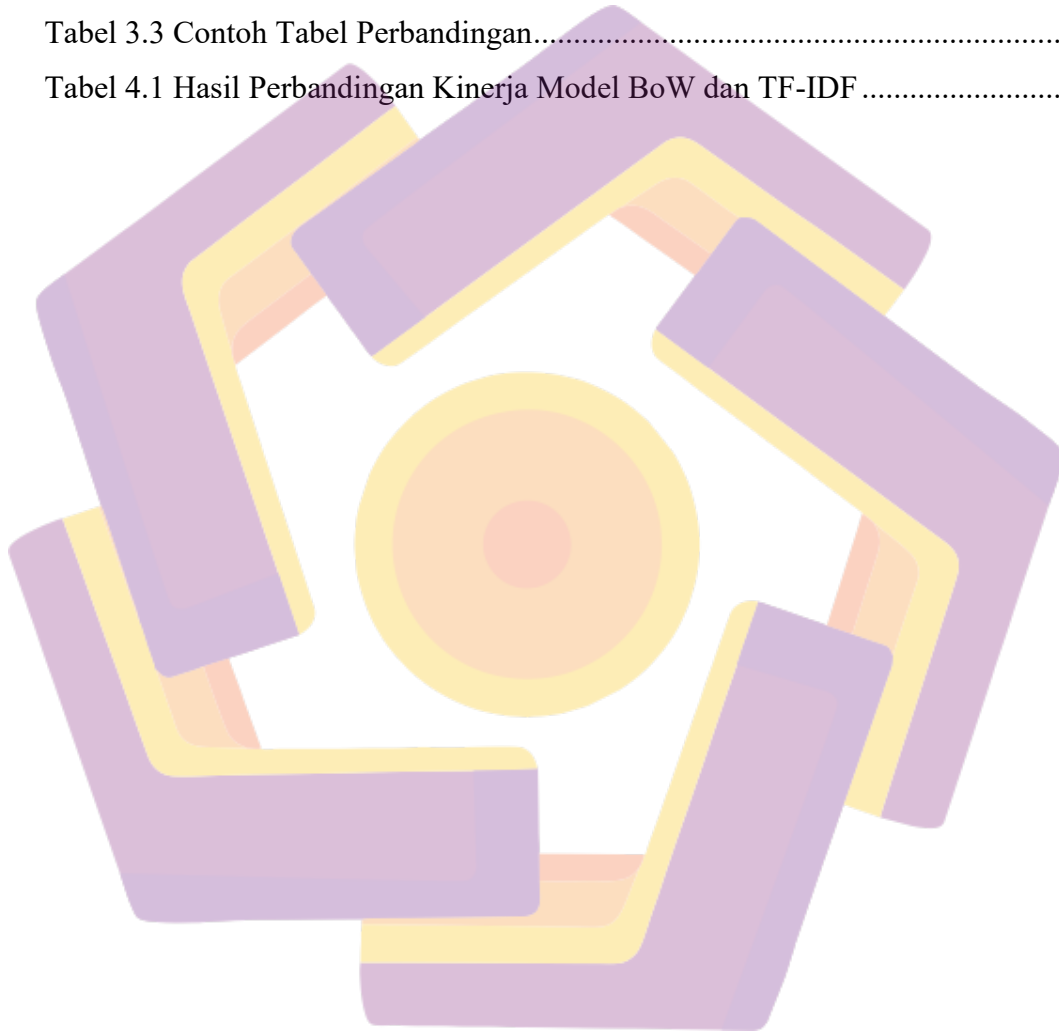
3.5.4.	Metode Pengujian Kinerja	39
3.5.5.	Validasi Pengujian	39
3.5.6.	Rencana Perbandingan Metode.....	41
BAB IV HASIL DAN PEMBAHASAN		43
4.1.	Persiapan Penelitian	43
4.1.1.	Persiapan Lingkungan <i>Google Colaboratory</i>	43
4.1.2.	<i>Library</i> dan <i>Tools</i> yang Digunakan	44
4.2.	Tahapan Eksperimen.....	46
4.2.1.	Deskripsi Dataset	46
4.2.2.	<i>Text Preprocessing</i>	47
4.2.3.	Pembagian Data (<i>Split Data</i>)	48
4.2.4.	Representasi Teks (BoW dan TF-IDF)	49
4.2.5.	Penyeimbangan Data dengan SMOTE	50
4.3.	Pelatihan dan Pengujian Model	51
4.3.1.	Pelatihan Model <i>Multinomial Naïve Bayes</i>	51
4.3.2.	Pengujian Model Menggunakan Data Uji.....	51
4.4.	Evaluasi dan Validasi Model	52
4.4.1.	Evaluasi Kinerja Model	52
4.4.2.	<i>Confusion Matrix</i>	53
4.4.3.	Analisis Waktu Komputasi	53
4.4.4.	Validasi Model (<i>5-Fold Cross Validation</i>).....	54
4.5.	Perbandingan dan Pemilihan Model	54
4.5.1.	Perbandingan Kinerja BoW dan TF-IDF	55
4.5.2.	Pemilihan Model Terbaik	55
4.6.	Hasil Pengujian Model.....	56

4.6.1.	Hasil Pengujian BoW.....	56
4.6.2.	Hasil Pengujian TF-IDF.....	58
4.6.3.	Analisis Waktu Komputasi	61
4.6.4.	Perbandingan Akurasi BoW vs TF-IDF.....	62
4.6.5.	Perbandingan F1-Score.....	63
4.6.6.	Validasi Model (<i>5-Fold Cross Validation</i>).....	64
BAB V PENUTUP		66
5.1	Kesimpulan	66
5.2	Saran	66
REFERENSI		68



DAFTAR TABEL

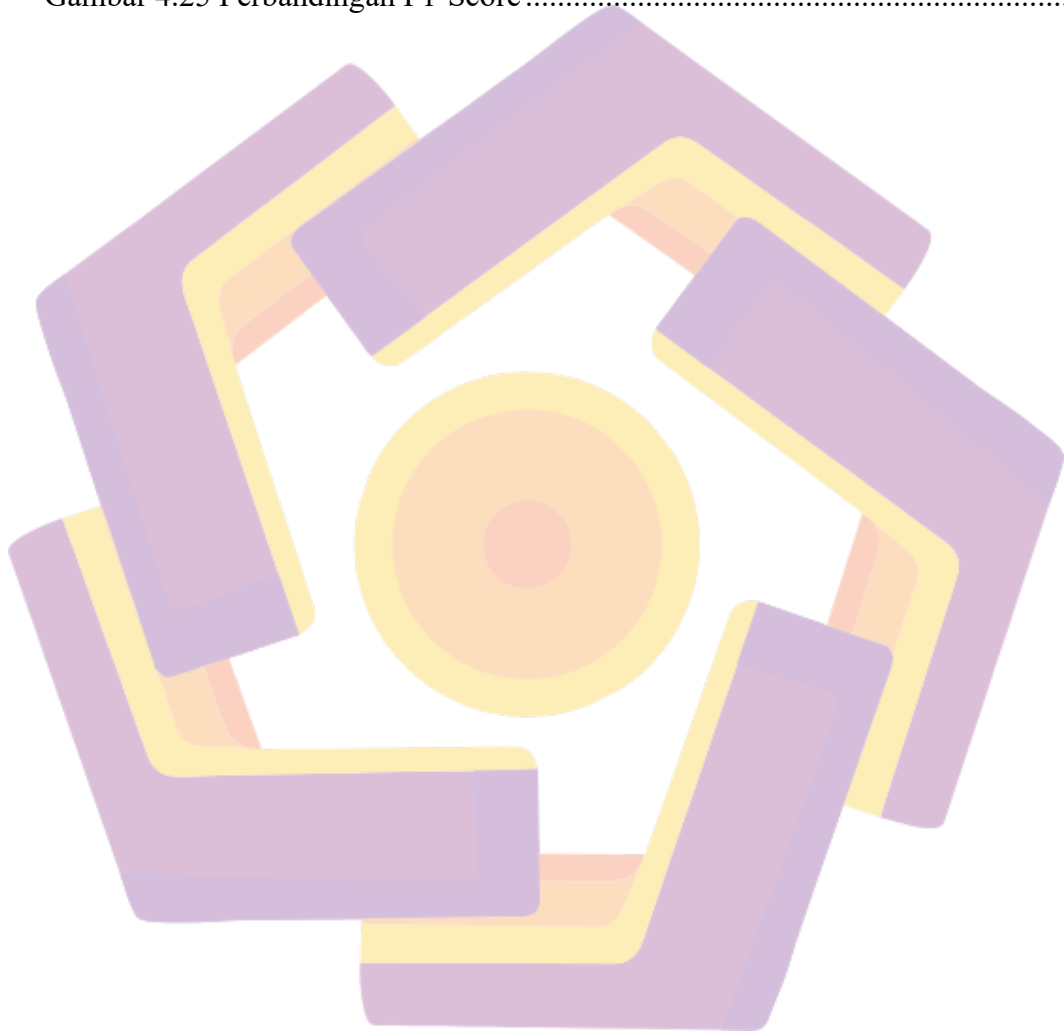
Tabel 2.1 Keaslian Penelitian	13
Tabel 2.2 Frekuensi Kata dalam Dokumen.....	26
Tabel 2.3 Confusion Matrix	28
Tabel 3.1 Tabel Spesifikasi Perangkat Keras.....	35
Tabel 3.2 Tabel Spesifikasi Perangkat Lunak	36
Tabel 3.3 Contoh Tabel Perbandingan.....	41
Tabel 4.1 Hasil Perbandingan Kinerja Model BoW dan TF-IDF	64



DAFTAR GAMBAR

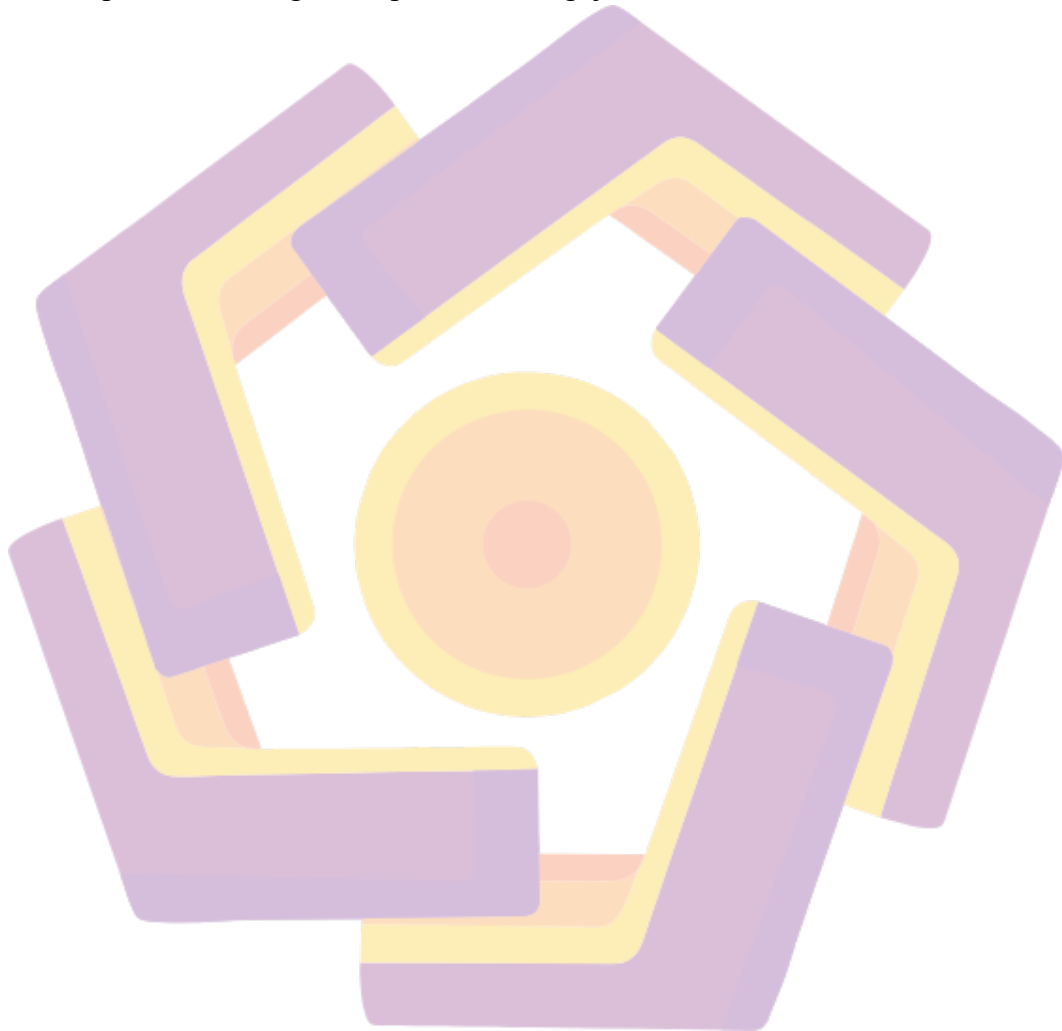
Gambar 2.1 Alur proses klasifikasi.....	18
Gambar 2.2 Alur Preprocessing	19
Gambar 2.3 Alur Pembagian Dataset Menjadi Data Latih dan Data Uji.....	21
Gambar 2.4 Contoh link phishing.....	23
Gambar 2.5 Contoh lain dari URL phishing.....	24
Gambar 2.6 Analogi Bag of Words	25
Gambar 3.1 Contoh Struktur Dataset URL Phishing.....	33
Gambar 3.2 Alur Penelitian	33
Gambar 3.1 Proses pengujian Model.....	38
Gambar 3.2 Ilustrasi Proses 5-Fold Cross Validation.....	40
Gambar 4.1 Tampilan Antarmuka Google Colaboratory	44
Gambar 4.2 Library yang Digunakan	45
Gambar 4.3 Distribusi Data URL Phishing	46
Gambar 4.4 Kode untuk Memuat Dataset.....	47
Gambar 4.5 Text Preprocessing	48
Gambar 4.6 Split Data.....	49
Gambar 4.7 Tokenisasi dan Representasi Teks	50
Gambar 4.8 Representasi Teks dan SMOTE	50
Gambar 4.9 Multinomial Naïve Bayes	51
Gambar 4.10 Pengujian Model menggunakan Data Uji	51
Gambar 4.11 Evaluasi Kinerja Model.....	52
Gambar 4.12 Confusion matrix hasil klasifikasi URL phishing.....	53
Gambar 4.13 Waktu Komputasi Vektorisasi (BoW dan TF-IDF).....	53
Gambar 4.14 Waktu Training + Prediction.....	54
Gambar 4.15 Proses validasi model menggunakan 5-fold cross validation	54
Gambar 4.16 Kode tabel perbandingan BoW dan TF-IDF	55
Gambar 4.17 Grafik Akurasi.....	55
Gambar 4.18 Proses penyimpanan model terbaik.....	56
Gambar 4.19 Classification Report BoW	57

Gambar 4.20 Confusion Matrix BoW.....	58
Gambar 4.21 Classification Report TF-IDF	59
Gambar 4.22 Confusion Matrix TF-IDF.....	60
Gambar 4.23 Waktu Komputasi BoW dan TF-IDF dengan SMOTE.....	61
Gambar 4.24 Diagram perbandingan Akurasi BoW dan TF-IDF.....	62
Gambar 4.25 Perbandingan F1-Score	63

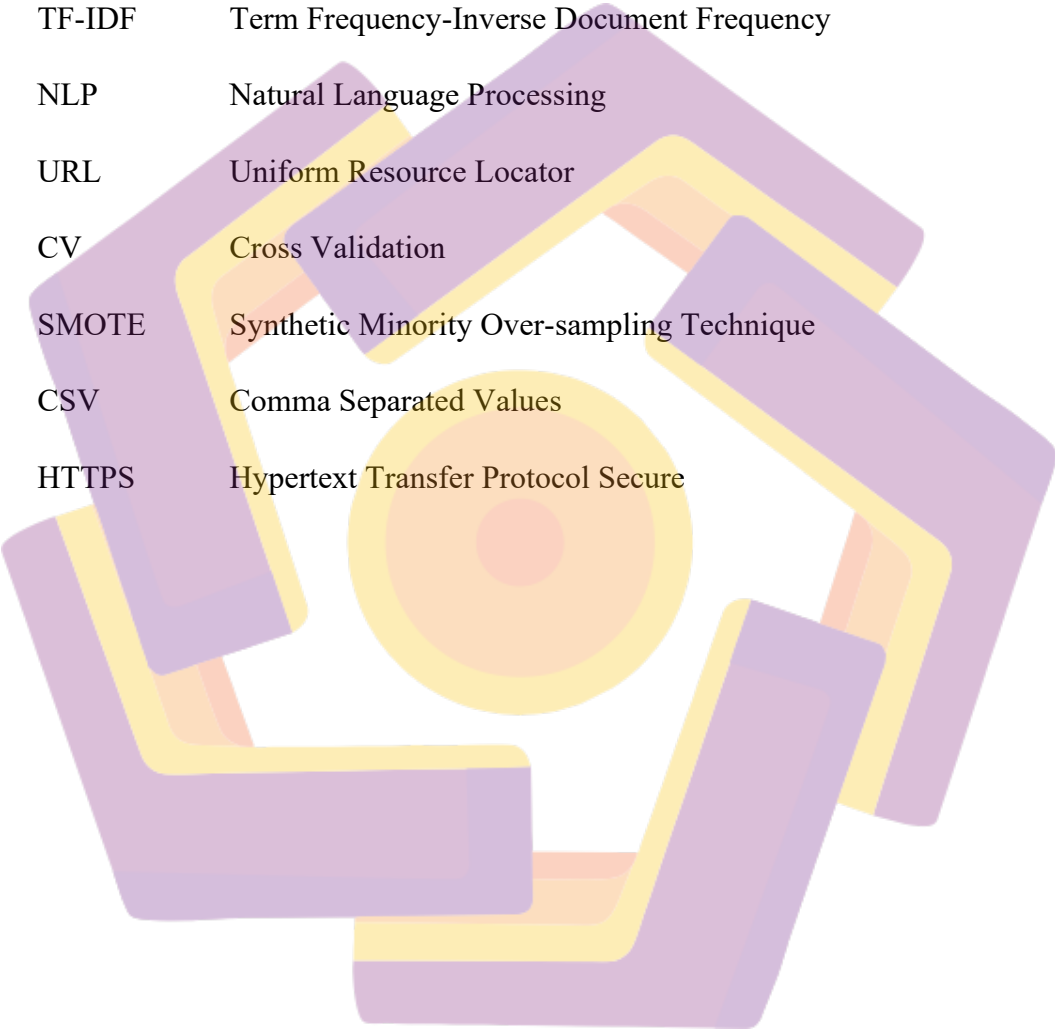


DAFTAR LAMPIRAN

Lampiran A. Source Code Deteksi Phishing URL Menggunakan Multinomial Naïve Bayes	71
Lampiran B. Potongan Output Hasil Pengujian	72

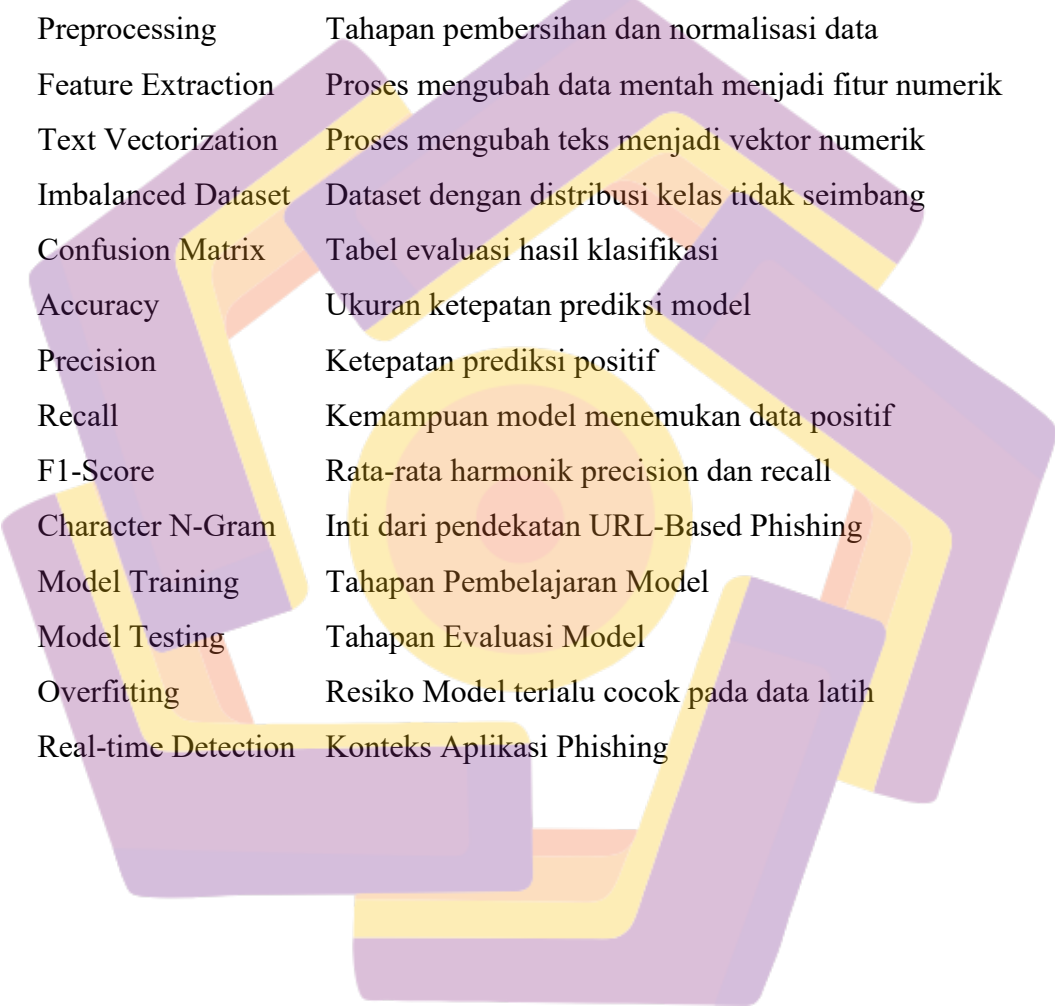


DAFTAR LAMBANG DAN SINGKATAN



SVM	Support Vector Machine
BoW	Bag of Words
TF-IDF	Term Frequency-Inverse Document Frequency
NLP	Natural Language Processing
URL	Uniform Resource Locator
CV	Cross Validation
SMOTE	Synthetic Minority Over-sampling Technique
CSV	Comma Separated Values
HTTPS	Hypertext Transfer Protocol Secure

DAFTAR ISTILAH



Phishing	Teknik penipuan untuk memperoleh informasi sensitif melalui URL palsu
Klasifikasi	Proses pengelompokan data ke dalam kelas tertentu
Dataset	Kumpulan data yang digunakan dalam penelitian
Preprocessing	Tahapan pembersihan dan normalisasi data
Feature Extraction	Proses mengubah data mentah menjadi fitur numerik
Text Vectorization	Proses mengubah teks menjadi vektor numerik
Imbalanced Dataset	Dataset dengan distribusi kelas tidak seimbang
Confusion Matrix	Tabel evaluasi hasil klasifikasi
Accuracy	Ukuran ketepatan prediksi model
Precision	Ketepatan prediksi positif
Recall	Kemampuan model menemukan data positif
F1-Score	Rata-rata harmonik precision dan recall
Character N-Gram	Inti dari pendekatan URL-Based Phishing
Model Training	Tahapan Pembelajaran Model
Model Testing	Tahapan Evaluasi Model
Overfitting	Resiko Model terlalu cocok pada data latih
Real-time Detection	Konteks Aplikasi Phishing

INTISARI

Phishing merupakan salah satu bentuk serangan siber yang memanfaatkan URL palsu untuk menipu pengguna agar memberikan informasi sensitif. Seiring meningkatnya jumlah serangan phishing, diperlukan metode deteksi otomatis yang efektif dan efisien. Penelitian ini bertujuan untuk membangun dan mengevaluasi model deteksi phishing URL menggunakan algoritma Multinomial Naïve Bayes dengan dua metode representasi teks, yaitu Bag of Words (BoW) dan Term Frequency–Inverse Document Frequency (TF-IDF).

Dataset yang digunakan terdiri dari beberapa kelas URL, yaitu benign, phishing, malware, dan defacement, yang memiliki distribusi data tidak seimbang. Untuk mengatasi permasalahan tersebut, diterapkan teknik Synthetic Minority Over-sampling Technique (SMOTE) pada data latih. Proses penelitian meliputi tahapan preprocessing teks, representasi fitur menggunakan character n-gram, pelatihan model, serta evaluasi menggunakan metrik accuracy, precision, recall, dan F1-score. Selain itu, dilakukan 5-Fold Cross Validation untuk mengukur stabilitas dan kemampuan generalisasi model.

Hasil pengujian menunjukkan bahwa kedua metode representasi teks mampu menghasilkan performa klasifikasi yang baik. Representasi BoW menghasilkan accuracy sebesar (78,20%) dan F1-score sebesar (78,59%), sedangkan TF-IDF memperoleh accuracy sebesar (77,90%) dan F1-score sebesar (79,16%). Dari sisi stabilitas, hasil validasi menunjukkan nilai CV-Mean yang konsisten pada kedua metode. Selain itu, BoW memiliki waktu komputasi vektorisasi yang lebih rendah dibandingkan TF-IDF, sehingga lebih efisien dalam pemrosesan data skala besar. Berdasarkan evaluasi kinerja, stabilitas model, dan efisiensi waktu komputasi, model Multinomial Naïve Bayes dengan representasi BoW dan penerapan SMOTE dipilih sebagai model terbaik dalam penelitian ini.

Kata kunci: phishing URL, Multinomial Naïve Bayes, Bag of Words, TF-IDF, SMOTE

ABSTRACT

Phishing is one of the most common cyber threats that exploits fraudulent URLs to deceive users into disclosing sensitive information. As phishing attacks continue to increase, effective and efficient automated detection methods are required. This study aims to develop and evaluate a phishing URL detection model using the Multinomial Naïve Bayes algorithm with two text representation approaches, namely Bag of Words (BoW) and Term Frequency–Inverse Document Frequency (TF-IDF).

The dataset used in this study consists of several URL classes, including benign, phishing, malware, and defacement, which exhibit an imbalanced class distribution. To address this issue, the Synthetic Minority Over-sampling Technique (SMOTE) was applied to the training data. The research process includes text preprocessing, feature representation using character n-grams, model training, and performance evaluation using accuracy, precision, recall, and F1-score metrics. In addition, 5-Fold Cross Validation was conducted to assess model stability and generalization capability.

The experimental results show that both text representation methods achieve good classification performance. The BoW representation obtained an accuracy of (78.20%) and an F1-score of (78.59%), while TF-IDF achieved an accuracy of (77.90%) and an F1-score of (79.16%). Cross-validation results indicate consistent performance across folds for both approaches. In terms of computational efficiency, BoW requires lower vectorization time compared to TF-IDF, making it more suitable for large-scale phishing URL detection. Based on overall classification performance, model stability, and computational efficiency, the Multinomial Naïve Bayes model with BoW representation and SMOTE is selected as the best-performing model in this study.

Keyword: *phishing URL, Multinomial Naïve Bayes, Bag of Words, TF-IDF, SMOTE*