

**ANALISIS KOMPARATIF ALGORITMA MACHINE LEARNING
UNTUK DETEKSI EMAIL SPAM**

SKRIPSI

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



disusun oleh

HUSEIN ALI ASGHAR

18.83.0313

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2023

**ANALISIS KOMPARATIF ALGORITMA MACHINE LEARNING
UNTUK DETEKSI EMAIL SPAM**

SKRIPSI

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



disusun oleh

HUSEIN ALI ASGHAR

18.83.0313

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2023

HALAMAN PERSETUJUAN

SKRIPSI

**ANALISIS KOMPARATIF ALGORITMA MACHINE LEARNING
UNTUK DETEKSI EMAIL SPAM**

yang disusun dan diajukan oleh

Husein Ali Asghar

18.83.0313

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 21 Agustus 2023

Dosen Pembimbing,



Banu Santoso, S.T., M.Eng.

NIK. 190302327

HALAMAN PENGESAHAN

SKRIPSI

**ANALISIS KOMPARATIF ALGORITMA MACHINE LEARNING
UNTUK DETEKSI EMAIL SPAM**

yang disusun dan diajukan oleh

Husein Ali Asghar

18.83.0313

Telah dipertahankan di depan Dewan Penguji
pada tanggal 21 Agustus 2023

Susunan Dewan Penguji

Nama Penguji

Arifiyanto Hadinegoro, S.Kom., M.T.
NIK. 190302289

Wahid Miftahul Ashari, S.Kom., M.T.
NIK. 190302452

Banu Santoso, S.T., M.Eng.
NIK. 190302327

Tanda Tangan



Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 21 Agustus 2023

DEKAN FAKULTAS ILMU KOMPUTER



Hanif Al Fatta, S.Kom., M.Kom.
NIK. 190302096

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

Nama mahasiswa : Husein Ali Asghar
NIM : 18.83.0313

Menyatakan bahwa Skripsi dengan judul berikut:

Analisis Komparatif Algoritma Machine Learning Untuk Deteksi Email Spam

Dosen Pembimbing : Banu Santoso, S.T., M.Eng.

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, 21 Agustus 2023

Yang Menyatakan,



Husein Ali Asghar

HALAMAN PERSEMBAHAN

Puji syukur kehadirat Allah SWT atas berkat dan rahmat yang telah diberikan kepada saya sehingga dapat menyelesaikan skripsi ini. Oleh karena itu dengan rasa syukur dan bangga saya ucapkan terimakasih kepada:

1. Allah SWT, berkat rahmat kasih sayang dan pertolongan-Nya saya bisa sampai sejauh ini.
2. Kedua orang tua saya, Bapak dan Ibu, yang selalu mendukung dan mendidik saya sejak kecil hingga saat ini. Mereka selalu memberikan dukungan moril dan doa yang tulus sehingga saya dapat menyelesaikan masa studi S1 dengan baik. Tanpa dukungan dan doa mereka, saya tidak akan dapat mencapai titik ini.
3. Dosen pembimbing saya, Bapak Banu Santoso, S.T., M.Eng. yang selalu memacu dan membimbing saya dalam proses menyelesaikan skripsi.
4. Bapak/Ibu dosen pembimbing, pengajar, serta penguji di Universitas AMIKOM Yogyakarta yang selalu memberikan pengalaman ilmu dengan tulus dan mengarahkan saya sehingga saya dapat menyelesaikan masa studi S1 dengan baik. Saya berharap segala ilmu yang telah diberikan oleh Bapak/Ibu akan menjadi ladang amal yang memberikan pahala jariyah dan selalu diberikan keberkahan ilmu dan rezeki oleh Tuhan yang Maha Esa.
5. Bapak Ali Reza Mulachela, Bapak A.M Safwan, dan Bapak Alm. Didiek Sri Wiyono yang telah mendukung saya baik secara materiil maupun moril selama menempuh pendidikan S1.
6. Serta semua pihak yang tidak bisa penulis sebutkan satu persatu.

KATA PENGANTAR

Dengan rasa syukur yang mendalam kepada Tuhan yang Maha Esa karena dengan rahmat-Nya saya dapat menyelesaikan penyusunan skripsi ini. Judul skripsi yang diajukan adalah "Analisis Komparatif Algoritma Machine Learning Untuk Deteksi Email Spam ". Skripsi ini diajukan untuk memenuhi persyaratan kelulusan mata kuliah skripsi di Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.

Oleh karena itu, saya ingin mengucapkan terima kasih yang sebesar besarnya kepada :

1. Allah SWT, berkat rahmat kasih sayang dan pertolongan-Nya saya bisa sampai sejauh ini.
2. Bapak Prof. Dr. M. Suyanto, MM. selaku Rektor Universitas AMIKOM Yogyakarta.
3. Bapak Hanif Al Fatta, S.Kom., M.Kom. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
4. Bapak Dony Ariyus, M.Kom. selaku Kaprodi Teknik Komputer Universitas AMIKOM Yogyakarta.
5. Bapak Banu Santoso, S.T., M.Eng. selaku dosen pembimbing saya.
6. Dosen penguji dan segenap Dosen serta Karyawan Universitas AMIKOM Yogyakarta yang telah berbagi ilmu dan pengalamannya.
7. Kedua orang tua yang telah mendoakan, mendukung dan memberikan semangat kepada saya.
8. Semua pihak yang telah membantu dan tidak dapat disebutkan satu persatu.

Pada akhir kata pengantar ini, penulis menyadari bahwa tidak ada yang sempurna dan pasti terdapat kesalahan dalam penyusunan skripsi. Oleh karena itu, penulis meminta maaf atas segala kesalahan yang mungkin terjadi. Penulis berharap agar skripsi ini dapat bermanfaat bagi pembaca dan dapat dijadikan sebagai referensi untuk pengembangan ilmu yang lebih baik.

Yogyakarta, 21 Januari 2023

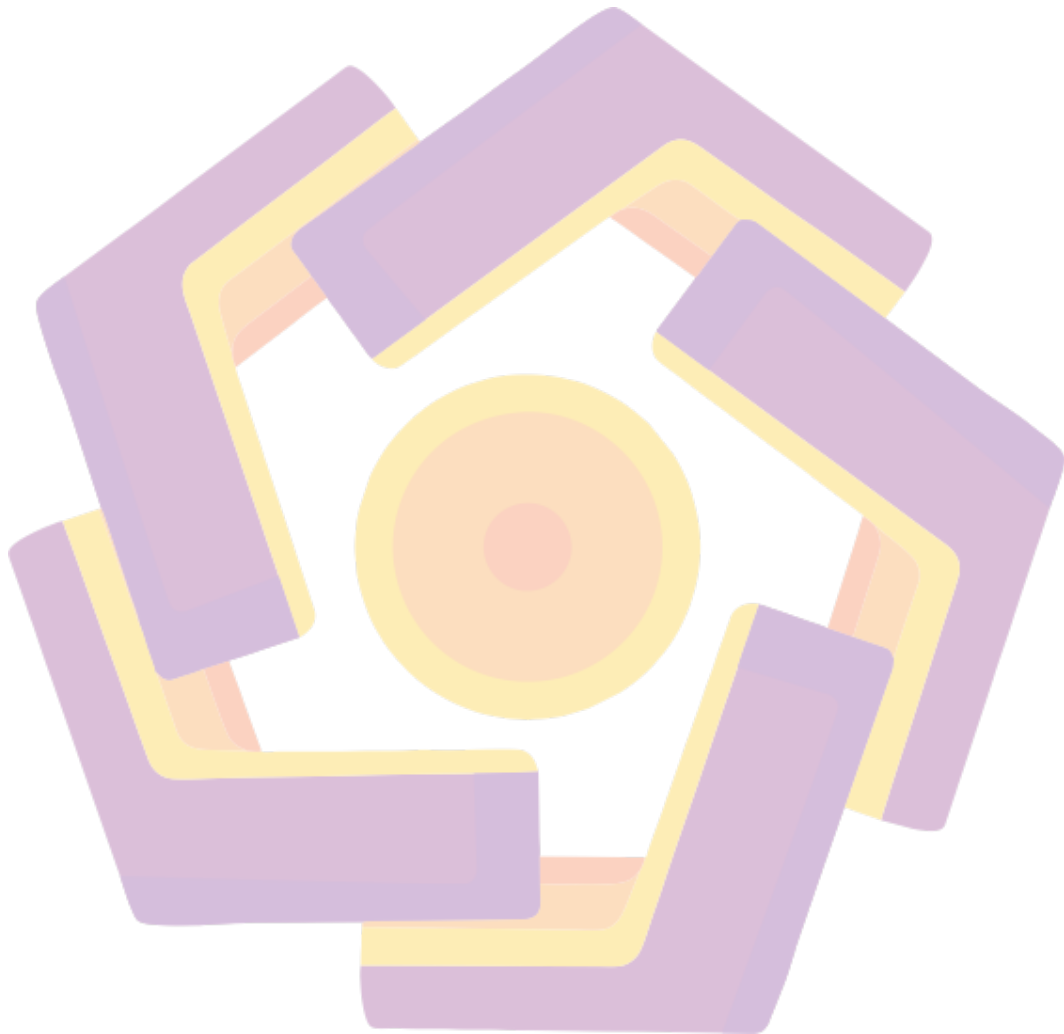
Penulis

DAFTAR ISI

HALAMAN JUDUL	i
HALAMAN PERSETUJUAN.....	Error! Bookmark not defined.
HALAMAN PENGESAHAN	Error! Bookmark not defined.
HALAMAN PERNYATAAN KEASLIAN SKRIPSI	iii
HALAMAN PERSEMBAHAN	v
KATA PENGANTAR	vi
DAFTAR ISI.....	vii
DAFTAR TABEL.....	x
DAFTAR GAMBAR.....	xi
DAFTAR LAMBANG DAN SINGKATAN	xiii
INTISARI	xiv
ABSTRACT.....	xv
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	3
1.3 Batasan Masalah	4
1.4 Tujuan Penelitian	4
1.5 Manfaat Penelitian	5
1.6 Sistematika Penulisan	5
BAB II TINJAUAN PUSTAKA	7
2.1 Studi Literatur	7
2.2 Dasar Teori	14
2.2.1 Surel dan Gambar <i>Spam</i>	15
2.2.2 Artificial Intelligence	16
2.2.3 Machine Learning	17
2.2.4 Klasifikasi	18
2.2.5 Algoritma Naïve Bayes.....	19
2.2.6 Algoritma Random Forest.....	21
2.2.7 Algoritma Support Vector Machine.....	23
2.2.8 Natural Language Processing.....	25

2.2.9 Optical Character Recognition	26
2.2.10 N-gram	30
2.2.11 Performance Evaluation Measure	30
2.2.12 Open-source Dataset	33
2.2.13 Google Collaboratory.....	33
2.2.14 Bahasa Pemrograman Python	34
2.2.15 CRISP-DM.....	34
BAB III METODE PENELITIAN	37
3.1 Objek Penelitian.....	37
3.2 Alur Penelitian	40
3.3 Alat dan Bahan.....	44
3.3.1 Data penelitian	45
3.3.2 Analisis Fungsional.....	45
3.3.3 Analisis Non-Fungsional	46
BAB IV HASIL DAN PEMBAHASAN	47
4.1 Collecting Data	47
4.2 Pre-processing Data	55
4.2.1 Data Cleansing dan Case Folding	55
4.2.2 Natural Language Processing.....	58
4.2.3 Ekstraksi Fitur dan Data Splitting.....	61
4.3 Penerapan Algoritma Multinomial Naïve Bayes	63
4.3.1 Implementasi Multinomial Naïve Bayes pada deteksi surel <i>spam</i>	64
4.3.2 Implementasi Multinomial Naïve Bayes pada deteksi gambar <i>spam</i> ..	67
4.4 Penerapan Algoritma Support Vector Machine.....	71
4.4.1 Implementasi Support Vector Machine untuk deteksi surel <i>spam</i>	71
4.4.2 Implementasi Support Vector Machine untuk deteksi gambar <i>spam</i> ..	74
4.5 Penerapan Algoritma Random Forest.....	77
4.5.1 Implementasi Random Forest untuk deteksi surel <i>spam</i>	77
4.5.2 Implementasi Random Forest untuk deteksi gambar <i>spam</i>	80

BAB V PENUTUP	84
5.1 Kesimpulan	84
5.2 Saran	85
REFERENSI	86



DAFTAR TABEL

Tabel 2. 1 Keaslian Penelitian	9
Tabel 3. 1 Contoh data surel spam dan non-spam	37
Tabel 3. 2 Contoh penerapan N-gram.....	43
Tabel 3. 3 Analisis Fungsional Sistem.....	45
Tabel 3. 4 Tabel Analisis Non-Fungsional	46
Tabel 4. 1 Data Splitting	62
Tabel 4. 2 Confusion Matrix MNB untuk surel <i>spam</i>	64
Tabel 4. 3 Hasil evaluasi pemodelan Multinomial Naive Bayes untuk surel <i>spam</i>	64
Tabel 4. 4 Confusion Matrix MNB pada gambar <i>spam</i>	68
Tabel 4. 5 Hasil evaluasi pemodelan MNB untuk gambar <i>spam</i>	68
Tabel 4. 6 Confusion Matrix SVM untuk deteksi surel <i>spam</i>	71
Tabel 4. 7 Hasil evaluasi pemodelan SVM untuk deteksi surel <i>spam</i>	72
Tabel 4. 8 Confusion Matrix SVM untuk deteksi gambar <i>spam</i>	74
Tabel 4. 9 Hasil evaluasi SVM untuk deteksi gambar <i>spam</i>	74
Tabel 4. 10 Confusion Matrix RF untuk deteksi surel <i>spam</i>	78
Tabel 4. 11 Hasil evaluasi pemodelan RF untuk deteksi surel <i>spam</i>	78
Tabel 4. 12 Tabel Confusion Matrix RF untuk deteksi gambar <i>spam</i>	80
Tabel 4. 13 Hasil evaluasi pemodelan RF untuk deteksi gambar <i>spam</i>	80

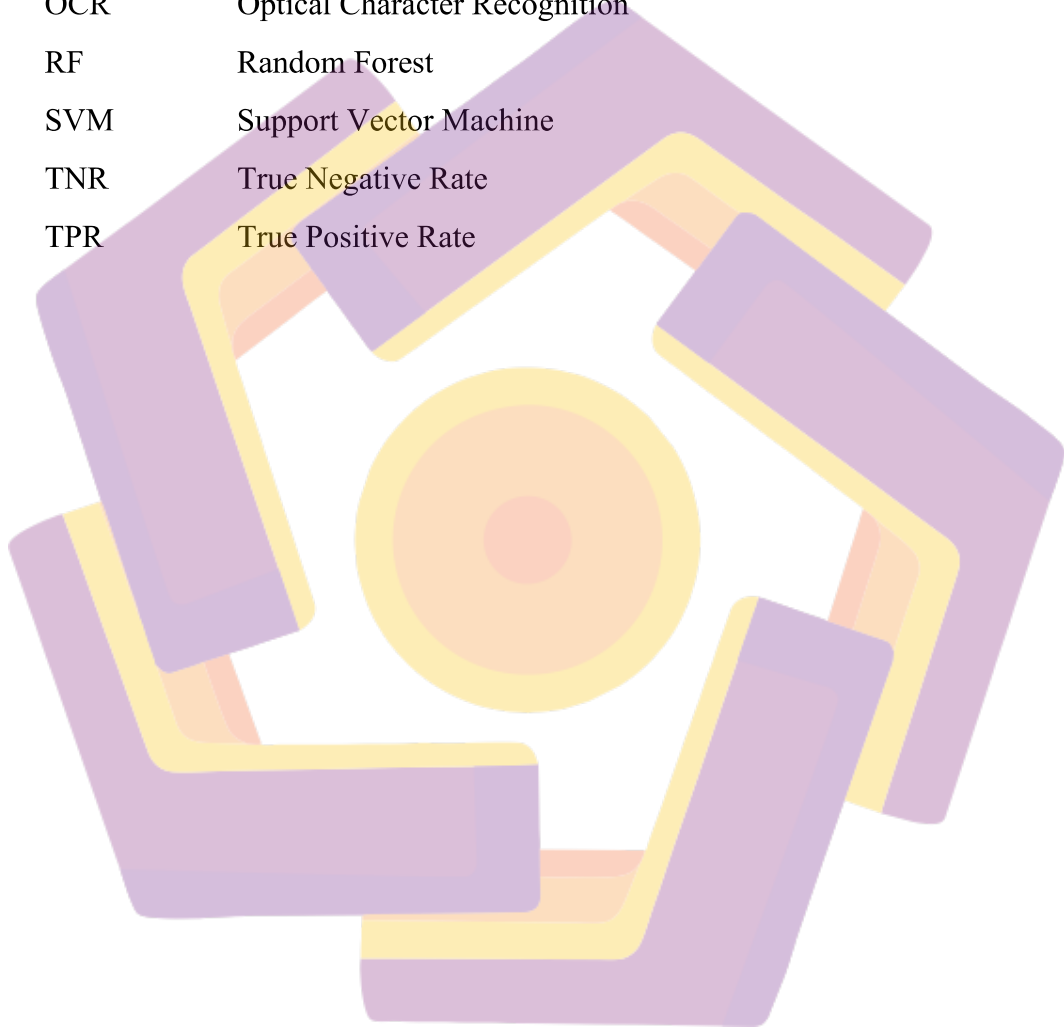
DAFTAR GAMBAR

Gambar 2. 1 Arsitektur Random Forest	22
Gambar 2. 2 Cara kerja OCR	27
Gambar 2. 3 Cara kerja Tesseract	29
Gambar 2. 4 Tabel Confusion Matrix	31
Gambar 2. 5 Tahapan dalam <i>CRISP-DM</i>	35
Gambar 3. 1 Contoh data gambar <i>spam</i>	39
Gambar 3. 2 Contoh gambar <i>natural</i> atau <i>non-spam</i>	39
Gambar 3. 3 Alur Penelitian	40
Gambar 4. 1 Informasi dataset surel <i>spam</i> yang didapat	47
Gambar 4. 2 Dataset surel <i>spam</i>	47
Gambar 4. 3 Informasi Dataset Gambar Spam dan Non-Spam	48
Gambar 4. 4 Instalasi OCR Pyteserract	48
Gambar 4. 5 Proses Importing Library yang dibutuhkan.....	49
Gambar 4. 6 Mounting Google Drive	49
Gambar 4. 7 Proses <i>unzipping dataset</i> gambar	49
Gambar 4. 8 Proses pencatatan nama file gambar dengan array	50
Gambar 4. 9 Beberapa nama file yang terdapat dalam dataset gambar	51
Gambar 4. 10 Proses Ekstraksi Teks pada Data Gambar Spam.....	52
Gambar 4. 11 Proses Ekstraksi Teks pada Data Gambar <i>Non-Spam</i>	52
Gambar 4. 12 Proses pembuatan file CSV ekstraksi gambar	53
Gambar 4. 13 Persentase data gambar spam dan non-spam	53
Gambar 4. 14 Perbandingan data gambar spam dan non-spam	54
Gambar 4. 15 Informasi dataset gambar yang berhasil didapat.....	54
Gambar 4. 16 Dataset gambar yang didapat	55
Gambar 4. 17 Jumlah data spam dan ham pada datset surel spam	56
Gambar 4. 18 Grafik perbandingan data spam dan ham pada dataset surel spam.....	56
Gambar 4. 19 Proses penyeimbangan dataset surel spam.....	56
Gambar 4. 20 Perubahan nama kolom dataset gambar spam	57
Gambar 4. 21 Merubah <i>value output</i> dataset surel spam	57
Gambar 4. 22 Tahapan data cleaning pada dataset	58
Gambar 4. 23 Kode program tokenisasi dataset surel <i>spam</i>	59
Gambar 4. 24 Kode program tokenisasi dataset gambar <i>spam</i>	59
Gambar 4. 25 Hasil dari proses tokenisasi pada dataset gambar <i>spam</i>	59
Gambar 4. 26 Hasil dari proses tokenisasi pada dataset surel <i>spam</i>	59
Gambar 4. 27 Instalasi NLTK Stopwords dan PorterStemmer.....	60
Gambar 4. 28 Kode program untuk <i>stopword removal</i> dan <i>word stemming</i>	60
Gambar 4. 29 Hasil penghapusan <i>stopwords</i> dan <i>stemming</i> dataset gambar spam	61
Gambar 4. 30 Hasil penghapusan <i>stopwords</i> dan <i>stemming</i> dataset surel <i>spam</i>	61
Gambar 4. 31 Proses Data Splitting untuk <i>dataset</i> gambar <i>spam</i>	61
Gambar 4. 32 Proses Data Splitting untuk <i>dataset</i> surel <i>spam</i>	62

Gambar 4. 33 Penerapan Ekstraksi Fitur N-gram	63
Gambar 4. 34 Kode Program Data Training Multinomial Naive Bayes.....	64
Gambar 4. 35 Diagram Garis Hasil Evaluasi Pemodelan MNB untuk surel <i>spam</i>	66
Gambar 4. 36 Diagram Garis Matthew Correlation Coefficient untuk MNB pada surel <i>spam</i>	67
Gambar 4. 37 Diagram Garis Hasil Evaluasi Pemodelan MNB untuk gambar <i>spam</i>	69
Gambar 4. 38 Diagram Garis Matthew Correlation Coefficient untuk MNB pada gambar <i>spam</i>	70
Gambar 4. 39 Kode Program Data Training Support Vector Machine	71
Gambar 4. 40 Diagram garis performa SVM dalam deteksi surel <i>spam</i>	73
Gambar 4. 41 Diagram Garis Matthew Correlation Coefficient untuk SVM pada surel <i>spam</i>	73
Gambar 4. 42 Diagram garis performa SVM dalam deteksi gambar <i>spam</i>	76
Gambar 4. 43 Diagram Garis Matthew Correlation Coefficient untuk SVM pada gambar <i>spam</i>	76
Gambar 4. 44 Kode Program Data Training Random Forest	77
Gambar 4. 45 Diagram garis performa Random Forest dalam deteksi surel <i>spam</i>	79
Gambar 4. 46 Diagram Garis Matthew Correlation Coefficient untuk Random Forest pada surel <i>spam</i>	79
Gambar 4. 47 Diagram garis performa Random Forest dalam deteksi gambar <i>spam</i>	82
Gambar 4. 48 Diagram Garis Matthew Correlation Coefficient untuk Random Forest pada gambar <i>spam</i>	82

DAFTAR LAMBANG DAN SINGKATAN

FPR	False Positive Rate
MCC	Matthew Correlation Coefficient
MNB	Multinomial Naïve Bayes
OCR	Optical Character Recognition
RF	Random Forest
SVM	Support Vector Machine
TNR	True Negative Rate
TPR	True Positive Rate



INTISARI

Penelitian ini bertujuan untuk melakukan perbandingan terhadap beberapa metode klasifikasi untuk klasifikasi surel spam dan gambar spam berdasarkan ekstraksi fitur N-gram. Penelitian ini menggunakan algoritma klasifikasi machine learning yaitu Multinomial Naïve Bayes, Random Forest dan Support Vector Machine. Dalam melakukan klasifikasi gambar spam, peneliti menggunakan metode Optical Character Recognition dengan menggunakan library Python yaitu pyteserract. Penelitian ini menggunakan open-source dataset yang berasal dari Kaggle untuk deteksi surel spam dan Spam Hunter untuk deteksi gambar spam.

Hasil dari penelitian ini menunjukkan bahwa algoritma Random Forest merupakan pemodelan machine learning terbaik dalam mendeteksi surel spam dengan tingkat akurasi 0,99, dan nilai Matthew Correlation Coefficient 0,999. Algoritma pemodelan machine learning dengan performa terbaik untuk mendeteksi gambar spam adalah algoritma Support Vector Machine dan algoritma Multinomial Naïve Bayes dengan tingkat akurasi yang sama yaitu 0,98 dan nilai Matthew Correlation Coefficient yang sama yaitu 0,976. Metode ekstraksi fitur N-gram terbaik pada penelitian ini adalah 1-gram.

Dari hasil penelitian ini, dapat disimpulkan bahwa algoritma Random Forest dengan ekstraksi fitur 1-gram memiliki performa terbaik dalam mendeteksi surel spam dan algoritma Multinomial Naïve Bayes serta Support Vector Machine dengan ekstraksi fitur 1-gram memiliki performa terbaik dalam mendeteksi gambar spam.

Kata kunci: *N-gram, Random Forest, Support Vector Machine, Multinomial Naïve Bayes, Deteksi Spam*

ABSTRACT

This study aims to compare several classification methods for the classification of spam emails and spam images based on N-gram feature extraction. This study uses a machine learning classification algorithm, namely Multinomial Naïve Bayes, Random Forest and Support Vector Machine. In classifying spam images, researchers use the Optical Character Recognition method using the Python library, namely pyteserract. This study uses open-source datasets from Kaggle for spam email detection and Spam Hunter for image spam detection.

The results of this study indicate that the Random Forest algorithm is the best machine learning model for detecting spam emails with an accuracy level of 0.99, and a Matthew Correlation Coefficient value of 0.999. The machine learning modeling algorithms with the best performance for detecting spam images are the Support Vector Machine algorithm and the Multinomial Naïve Bayes algorithm with the same accuracy level of 0.98 and the same Matthew Correlation Coefficient value of 0.976. The best N-gram feature extraction method in this study is 1-gram.

From the results of this study, it can be concluded that the Random Forest algorithm with 1-gram feature extraction has the best performance in detecting spam emails and the Multinomial Naïve Bayes algorithm and the Support Vector Machine with 1-gram feature extraction has the best performance in detecting spam images.

Keyword: *N-gram, Random Forest, Support Vector Machine, Multinomial Naïve Bayes, Spam Detection*