

**PERBANDINGAN KINERJA ALGORITMA SVM DAN C4.5
DALAM KLASIFIKASI SPAM EMAIL**

SKRIPSI NON-REGULER JURNAL

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi Informatika



disusun oleh

RAMADHANI KUSUMA HADI

19.11.3009

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2024

**PERBANDINGAN KINERJA ALGORITMA SVM DAN C4.5
DALAM KLASIFIKASI SPAM EMAIL**

SKRIPSI NON-REGULER JURNAL

untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi Informatika



disusun oleh

RAMADHANI KUSUMA HADI

19.11.3009

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2024

HALAMAN PERSETUJUAN

SKRIPSI

**PERBANDINGAN KINERJA ALGORITMA SVM DAN C4.5
DALAM KLASIFIKASI SPAM EMAIL.**

yang disusun dan diajukan oleh

RAMADHANI KUSUMA HADI

19.11.3009

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 27 Juni 2024

Dosen Pembimbing,


Mufia Sulistyono, M.kom

NIK. 190302248

HALAMAN PENGESAHAN
SKRIPSI
PERBANDINGAN KINERJA ALGORITMA SVM DAN C4.5
DALAM KLASIFIKASI SPAM EMAIL

yang disusun dan diajukan oleh

RAMADHANI KUSUMA HADI

19.11.3009

Telah dipertahankan di depan Dewan Penguji
pada tanggal 27 Juni 2024

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Yoga Prisyanto, S.Kom., M.Eng.
NIK. 190302412

Wiwid Widayani, M.Kom
NIK. 190302272

Mulia Sulistyono, M.Kom
NIK. 190302248



Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 27 Juni 2024

DEKAN FAKULTAS ILMU KOMPUTER



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

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

Nama mahasiswa : **Ramadhani Kusuma Hadi**
NIM : **19.11.3009**

Menyatakan bahwa Skripsi dengan judul berikut:

PERBANDINGAN KINERJA ALGORITMA SVM DAN C4.5 DALAM KLASIFIKASI SPAM EMAIL

Dosen Pembimbing : **Mulia Sulistiyono, M.Kom**

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, 27 juni 2024

Yang Menyatakan,



Ramadhani Kusuma Hadi

HALAMAN PERSEMBAHAN

Alhamdulillah, Segala Puji Bagi ALLAH SWT. Kepada-Nya kita memohon, memuji, dan meminta perlindungan, pengampunan serta petunjuk Kepada-Nya. Kita berlindung kepada ALLAH SWT dari kejahatan diri kita dan keburukan amal-amal kita. Dengan ini penulis mempersembahkan skripsi ini kepada semua pihak yang telah membantu dan memberikan motivasi serta dukungan yang besar sehingga skripsi ini dapat diselesaikan dengan sebaik-baiknya.

Melalui skripsi ini saya persembahkan terimakasih serta syukur kepada:

1. Rasa syukur kepada ALLAH SWT yang telah memberikan rahmat berupa karunia untuk terus mengucap syukur dan sabar dalam menyelesaikan tugas akhir skripsi.
2. Kepada orang tua saya yang saya cintai, bapak Susilo Hariyanto dan Wiwit Sulastri semoga selalu diberikan kesehatan atas doa dan motivasi serta dukungannya dalam membimbing penulis.
3. Saudara penulis Yoga Nourma Hadi Pratama serta keluarga kecilnya dengan Khainunnifa dengan malaikat kecilnya kaivan yang selalu mendukung penulis
4. Bapak Mulia Sulistiyono, M.Kom yang telah menjadi pembimbing penulis dalam menyelesaikan skripsi ini.
5. Bapak/Ibu dosen yang telah memberikan ilmu dan pendidikan selama berkuliah di Amikom.
6. Teman-teman pengurus inti HMIF periode 2021/2022 yang Bersama-sama berjuang dalam mencapai kesuksesan yang penuh suka dan duka. Anggota pengurus dan Kader HMIF 2021/2022 yang menjadi pemberi dukungan serta semangat saat menjadi Pengurus Himpunan.
7. Anggota grup "SKRIPSHIT KELAR MAKSIMAL SMSTER 8" yang telah mensupport dan berkontribusi dalam pengerjaan skripsi ini.
8. 'dla' sosok periang tapi ada kesedihan yang mendalam dibalikinya yang cukup kuat memberikan dorongan, motivasi, serta menyemangati sehingga dapat melewati semua rintangan yang ada, dan juga semangat terus kamu.

KATA PENGANTAR

Puji syukur kehadirat ALLAH SWT karena berkat rahmat serta hidayah-Nya penulis diberikan kesehatan dan kekuatan serta iman yang kuat sehingga dapat menyelesaikan skripsi yang berjudul **“PERBANDINGAN KINERJA ALGORITMA SVM DAN C4. 5 DALAM KLASIFIKASI SPAM EMAIL”**. Oleh karena itu penulis mengucapkan syukur karena dapat menyelesaikan skripsi ini. Skripsi ini diajukan kepada Program Studi S1 Informatika, Fakultas Ilmu Komputer, serta Universitas AMIKOM Yogyakarta

Penulis dalam menyelesaikan skripsi mendapatkan beberapa hambatan dalam berbagai hal, namun banyak pihak yang selalu membantu sehingga skripsi ini dapat diselesaikan dengan baik. Sengan begitu penulis hendak mengungkapkan terima kasih terhadap:

1. Bapak Prof. Dr. M Suyanto, M.M selaku Rektor Universitas AMIKOM Yogyakarta.
2. Bapak Hanif Al Fatta, S.Kom., M.Kom., Ph.D. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
3. Ibu Windha Mega Pradya D, M.Kom., selaku Ketua Program Studi S1 Informatika Universitas AMIKOM Yogyakarta.
4. Bapak Mulia Sulistiyono, M.Kom selaku dosen pembimbing yang selalu memberikan saran dan arahan dalam proses penulisan skripsi ini
5. Bapak Susilo Hariyanto dan Ibu Wiwit Sulastrri selaku orang tua serta Yoga Nourma Hadi Pratama beserta keluarga kecilnya dengan Khainunnisa Fasiha dengan malaikat kecilnya Kaivan dalam memberikan doa dan motivasi kepada penulis.

Yogyakarta, 27 Juni 2024

Penulis

DAFTAR ISI

HALAMAN JUDUL	i
HALAMAN PERSEMBAHAN	v
KATA PENGANTAR	vi
DAFTAR ISI.....	vii
DAFTAR TABEL.....	ix
DAFTAR GAMBAR.....	x
DAFTAR LAMPIRAN.....	Error! Bookmark not defined.
DAFTAR LAMBANG DAN SINGKATAN	xii
DAFTAR ISTILAH	xiii
INTISARI	xiv
<i>ABSTRACT</i>	xv
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	5
BAB II TINJAUAN PUSTAKA	6
2.1 Studi Literatur.....	6
2.2 Dasar Teori	18
2.2.1 Data Mining.....	18
2.2.2 Klasifikasi.....	19
2.2.3 Spam Email.....	21
2.2.7 Algoritma SVM	23
2.2.8 <i>Confusion Matrix</i>	25
BAB III METODE PENELITIAN	27
3.1 Alur Penelitian	27
3.2 Objek Penelitian.....	31
3.2.1 Pengumpulan Data	31

3.3	Alat dan Bahan.....	32
3.3.2	Bahan Penelitian.....	32
BAB IV HASIL DAN PEMBAHASAN		34
4.1	Pengumpulan Data.....	34
4.1.1	Membaca dan menampilkan <i>Dataset</i>	34
4.1.2	Perubahan nama kolom.....	35
4.2	Exploratory Data Analysis (EDA).....	35
4.2.1	Menampilkan dalam bentuk <i>Diagram Pie</i>	36
4.2.2	Penjumlahan karakter.....	36
4.3	<i>Pre-Processing</i>	40
4.3.1	<i>Case Folding</i>	40
4.3.2	<i>Stemming Data</i>	40
4.3.3	<i>Tokenizing</i>	41
4.3.4	Menghilangkan <i>Stopwords</i>	41
4.3.5	Memilih kolom yang relevan akan penelitian.....	42
4.3.6	Menampilkan <i>WordCloud</i>	42
4.4	Pembobotan Fitur (TF-IDF).....	43
4.5	<i>Splitting Data</i>	44
4.6	Klasifikasi.....	45
4.7	Evaluasi.....	46
4.7.1	Perbandingan klasifikasi Algoritma.....	46
4.7.2	Perhitungan <i>Confusion Matrix</i>	46
4.7.3	Hasil Evaluasi.....	50
BAB V PENUTUP		51
5.1	Kesimpulan.....	51
5.2	Saran	51
REFERENSI		52
LAMPIRAN.....		Error! Bookmark not defined.

DAFTAR TABEL

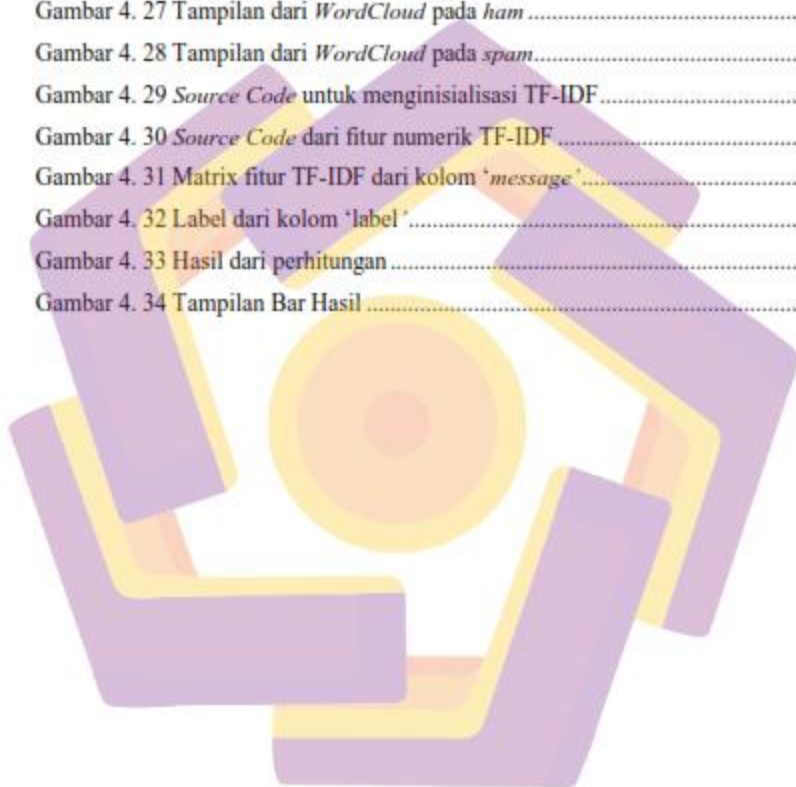
Tabel 2.1 Tabel Keaslian	9
Tabel 2.2 Stemming Data.....	18
Tabel 2.3 <i>Confusion Matrix</i>	25
Tabel 3. 1 Alur Penelitian	27
Tabel 3. 2 Deskripsi Variabel Penelitian	32
Tabel 4. 1 <i>Splitting Data</i>	45
Tabel 4. 2 Perbandingan kinerja Algoritma dengan <i>max_features</i> 3000	46
Tabel 4. 3 Perbandingan kinerja Algoritma dengan <i>max_features</i> 1500	46
Tabel 4. 4 Tampilan <i>Confusion Matrix</i> dari algoritma SVM	47
Tabel 4. 5 Tampilan <i>Confusion Matrix</i> dari algoritma C4. 5	48



DAFTAR GAMBAR

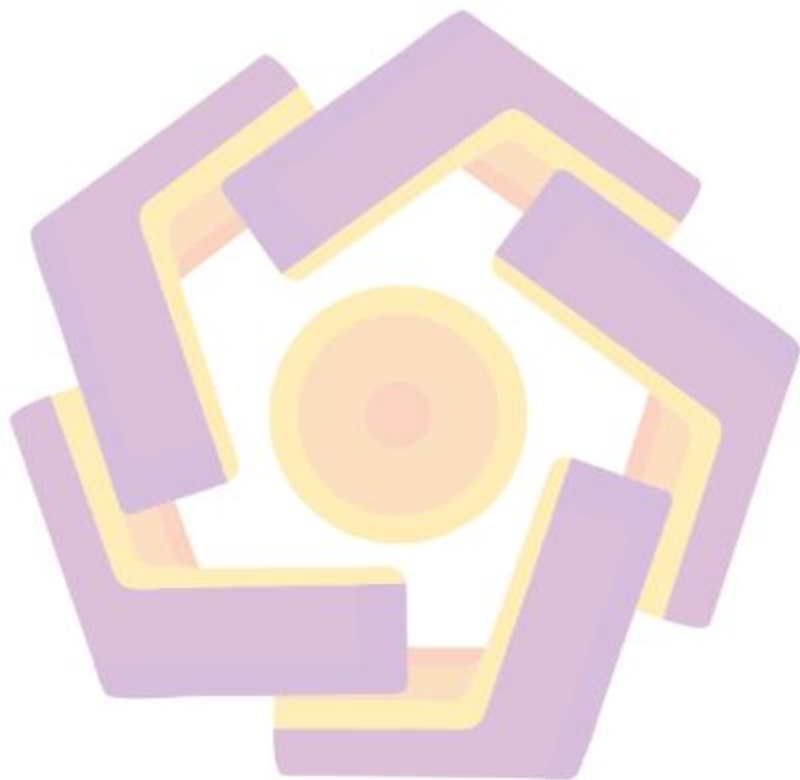
Gambar 3. 1 <i>Case Folding</i>	28
Gambar 3. 2 <i>Stemming</i>	29
Gambar 3. 3 <i>Tokenizing</i>	29
Gambar 3. 4 <i>Stopwords</i>	29
Gambar 3. 5 Pembobotan Fitur (TF-IDF)	30
Gambar 3. 6 Pelatihan model SVM	30
Gambar 3. 7 Pelatihan model C4.5	30
Gambar 3. 8 Evaluasi Model SVM dan C4. 5	31
Gambar 3. 9 Jumlah Dataset	31
Gambar 4. 1 <i>Source Code</i> pemanggilan <i>Dataset</i>	34
Gambar 4. 2 Tampilan <i>Dataset Spam Email</i>	34
Gambar 4. 3 Informasi <i>Dataset Spam Email</i>	35
Gambar 4. 4 <i>Source Code</i> untuk mengganti nama kolom	35
Gambar 4. 5 Tampilan <i>Dataset</i> setelah dirubah namanya	35
Gambar 4. 6 Data 'ham' dan 'spam' dalam <i>Dataset</i>	36
Gambar 4. 7 Tampilan Diagram Pie untuk data 'ham' dan 'spam'	36
Gambar 4. 8 Membuat kolom tabel baru bernama <i>num_characters</i>	37
Gambar 4. 9 Tampilan karakter yang terdapat pada data	37
Gambar 4. 10 Membuat kolom tabel baru bernama <i>num_words</i>	37
Gambar 4. 11 Membuat kolom tabel baru bernama <i>num_sentences</i>	38
Gambar 4. 12 Ekplorasi Data 'ham'	38
Gambar 4. 13 Ekplorasi Data 'spam'	38
Gambar 4. 14 Tampilan grafik dari <i>num_characters</i>	39
Gambar 4. 15 Tampilan grafik dari <i>num_words</i>	39
Gambar 4. 16 Contoh kalimat sebelum <i>Pre-processing</i>	40
Gambar 4. 17 <i>Source Code</i> dari <i>Case Folding</i>	40
Gambar 4. 18 Tampilan kalimat setelah <i>Case Folding</i>	40
Gambar 4. 19 <i>Source Code</i> dari <i>Stemming Data</i>	40
Gambar 4. 20 Tampilan kalimat setelah <i>Stemming Data</i>	41
Gambar 4. 21 <i>Source Code</i> dari <i>Tokenizing</i>	41

Gambar 4. 22 Tampilan kalimat setelah <i>Tokenizing</i>	41
Gambar 4. 23 <i>Source Code</i> dari menghilangkan <i>Stopwords</i>	41
Gambar 4. 24 Tampilan kalimat setelah menghilangkan <i>Stopwords</i>	42
Gambar 4. 25 <i>Source Code</i> kolom yang akan digunakan.....	42
Gambar 4. 26 Tampilan kolom yang digunakan.....	42
Gambar 4. 27 Tampilan dari <i>WordCloud</i> pada <i>ham</i>	43
Gambar 4. 28 Tampilan dari <i>WordCloud</i> pada <i>spam</i>	43
Gambar 4. 29 <i>Source Code</i> untuk menginisialisasi TF-IDF.....	44
Gambar 4. 30 <i>Source Code</i> dari fitur numerik TF-IDF	44
Gambar 4. 31 Matrix fitur TF-IDF dari kolom ' <i>message</i> '.....	44
Gambar 4. 32 Label dari kolom ' <i>label</i> '	44
Gambar 4. 33 Hasil dari perhitungan	50
Gambar 4. 34 Tampilan Bar Hasil	50



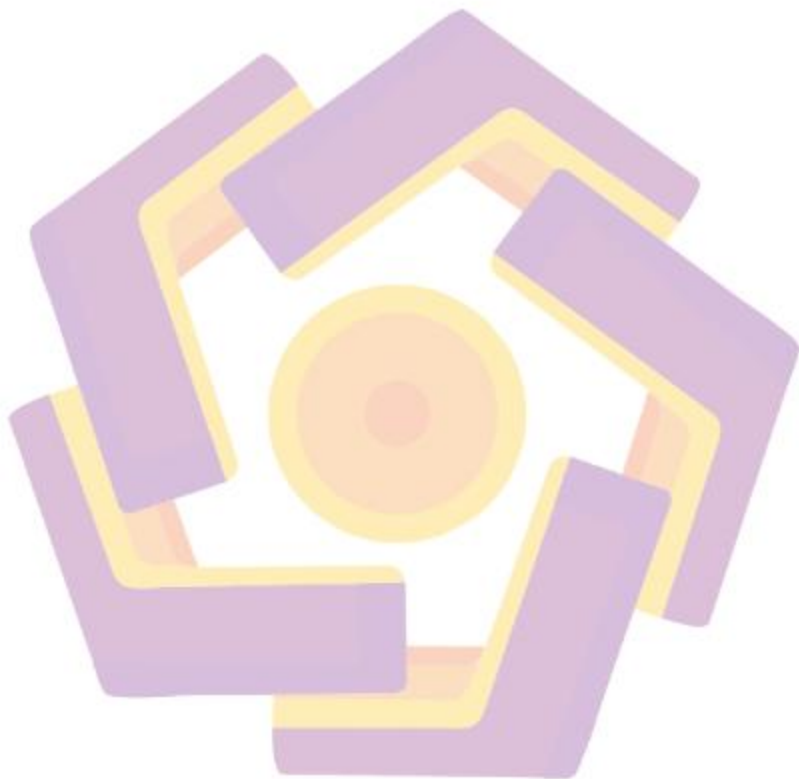
DAFTAR LAMBANG DAN SINGKATAN

SVM	Support Vector Machines
C4.5	Pohon Keputusan C4.5
TF-IDF	Term Frequency-Inverse Document Frequency



DAFTAR ISTILAH

Vektor	besaran yang mempunyai arah
Eigen Value	akar akar persamaan



INTISARI

Email merupakan surat elektronik yang memungkinkan seseorang untuk mengirimkan dan menerima pesan lebih dari satu orang. Email memiliki kelebihan diantara lain cepat, hemat dan dapat mengirimkan pesan dalam berbagai bentuk dokumen yang membuatnya populer di seluruh dunia. Banyaknya email yang dikirim setiap harinya sering disalahgunakan oleh beberapa pihak yang tidak bertanggungjawab untuk mengirim informasi iklan produk jasa dan berbagai informasi yang tidak berguna atau tidak diinginkan oleh pengguna, yang dimana itu disebut *SPAM*.

Untuk menyaring email *SPAM* digunakan banyak algoritma klasifikasi, algoritma tersebut diantaranya adalah *Support Vector Machine (SVM)* serta *Decision Tree (C4. 5)*. Pada penelitian ini akan dilakukan perbandingan algoritma SVM dan C4. 5. Untuk melihat kinerja dari permodelan algoritma maka dilakukan pembagian uji *max_features* 1500 dan 3000 data serta menggunakan perbandingan rasio 90% : 10%, 80% : 20%, 70% : 30% dan 60% : 40% yang bertujuan menilai data tersebut akan menghasilkan nilai akurasi tertinggi pada pengujian yang dilakukan.

Berdasarkan hasil pengujian yang menggunakan data sebanyak 10.000 baris data dan 2 kolom, menggunakan algoritma *Support Vector Machine (SVM)* serta *Decision Tree (C4. 5)* diperoleh hasil akurasi terbaik menggunakan *Max_features* 3000 dengan rasio 90% : 10% dimana SVM mendapatkan sebesar 98.20% dan C4. 5 sebesar 93.30%. Dari hasil uji coba tersebut membuktikan bahwa algoritma *Support Vector Machine (SVM)* merupakan metode terbaik dalam melakukan proses klasifikasi *spam email*.

Kata kunci: Klasifikasi, Spam-email, SVM, C4.5.

ABSTRACT

Email is an electronic mail that allows one to send and receive messages from more than one person. Email has advantages including fast, economical and can send messages in various forms of documents which make it popular around the world. The number of emails sent every day is often misused by some irresponsible parties to send information on advertisements for service products and various information that is useless or unwanted by users, which is called SPAM.

To filter SPAM emails, many classification algorithms are used, including Support Vector Machine (SVM) and Decision Tree (C4.5). In this study, a comparison of the SVM and C4.5. To see the performance of the algorithm modeling, the max_features test division is 1500 and 3000 data and uses a ratio of 90% : 10%, 80% : 20%, 70% : 30% and 60% : 40% which aims to assess the data that will produce the highest accuracy value in the test conducted.

Based on the test results using data as much as 10,000 rows of data and 2 columns, using the Support Vector Machine (SVM) and Decision Tree (C4.5) algorithms, the best accuracy results were obtained using Max_features 3000 with a ratio of 80%: 20% where SVM gets 98.20% and C4.5 amounted to 93.30%. From the test results, it proves that the Support Vector Machine (SVM) algorithm is the best method in performing the email spam classification process.

Keyword: *Classification, Spam-email, SVM, C4.5*