

**PERBANDINGAN ALGORITMA PADA *MACHINE LEARNING* UNTUK  
KLASIFIKASI *EMAIL SPAM***

**SKRIPSI**



disusun oleh

**Hery Iswanto**

**17.11.1713**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA**

**2021**

**PERBANDINGAN ALGORITMA PADA *MACHINE LEARNING* UNTUK  
KLASIFIKASI *EMAIL SPAM***

**SKRIPSI**

Untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Hery Iswanto**

**17.11.1713**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA**

**2021**

**PERSETUJUAN**

**SKRIPSI**

**PERBANDINGAN ALGORITMA PADA *MACHINE LEARNING* UNTUK  
KLASIFIKASI *EMAIL SPAM***

yang dipersiapkan dan disusun oleh

**Hery Iswanto**

**17.11.1713**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 8 November 2021

**Dosen Pembimbing**

**Erni Sentwatl, S.Kom., M.Cs.**

**NIK. 190302231**

**PENGESAHAN**

**SKRIPSI**

**PERBANDINGAN ALGORITMA PADA *MACHINE LEARNING* UNTUK  
KLASIFIKASI *EMAIL SPAM***

yang dipersiapkan dan disusun oleh

**Hery Iswanto**

**17.11.1713**

telah dipertahankan di depan Dewan Penguji  
pada tanggal 18 November 2021

**Susunan Dewan Penguji**

**Nama Penguji**

**Yuli Astuti, M.Kom**

**NIK. 190302146**

**Anna Balta, M.Kom**

**NIK. 190302290**

**Erni Seniwati, S.Kom, M.Cs**

**NIK. 190302231**

**Tanda Tangan**

Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 18 November 2021

**DEKAN FAKULTAS ILMU KOMPUTER**

**Hanif Al Fatta, M.Kom**

**NIK. 190302096**

## PERNYATAAN

Saya yang bertandatangan dibawah ini menyatakan bahwa, skripsi ini merupakan karya saya sendiri (ASLI), dan isi dalam skripsi ini tidak terdapat karya yang pernah diajukan oleh orang lain untuk memperoleh gelar akademis di suatu institusi pendidikan tinggi manapun, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis dan/atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Segala sesuatu yang terkait dengan naskah dan karya yang telah dibuat adalah menjadi tanggungjawab saya pribadi.

Yogyakarta, 18 November 2021



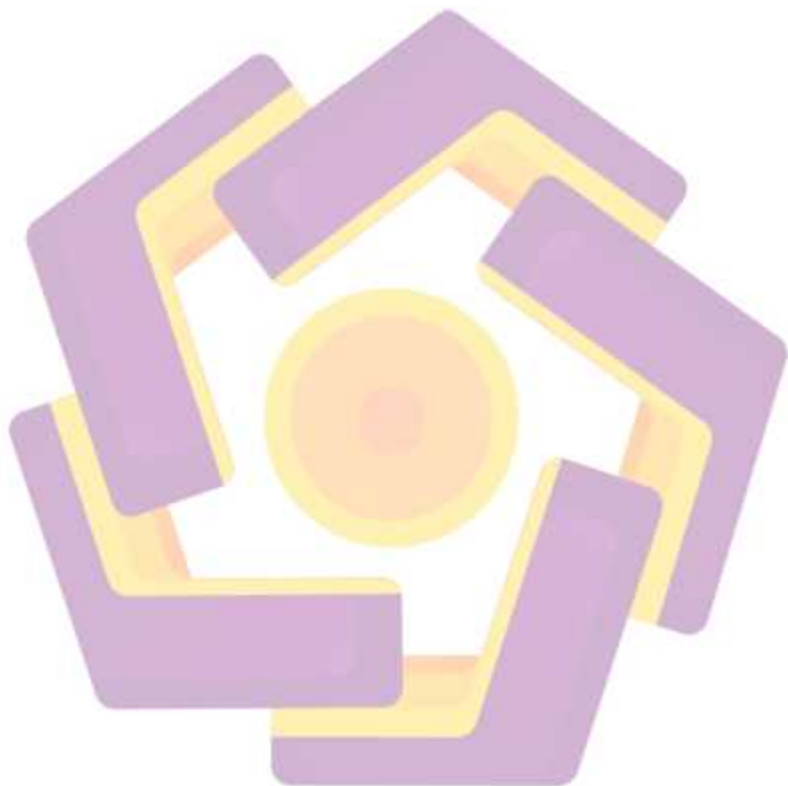
Hery Iswanto

NIM. 17.11.1713

**MOTTO**

“Didalam hidup ada saat untuk berhati-hati, atau berhenti berlari.”

**Talfun – Barasuara**



## PERSEMBAHAN

Puji syukur kepada Allah SWT yang telah memberikan saya sehat rohani maupun jasmani, nikmat yang tak terkira sehingga dapat menyelesaikan tugas akhir ini. Dengan ini saya persembahkan untuk :

1. Allah SWT yang telah memberikan kemudahan dan kelancaran dalam menyelesaikan skripsi ini.
2. Kedua orang tua yang selalu memberikan dukungan doa dan support kepada saya sehingga dapat berada dititik saat ini.
3. Untuk ibu dosen pembimbing saya Ibu Erni Seniwati, S.Kom, M.Cs terimakasih atas waktu dan ilmunya dalam membimbing saya menyelesaikan skripsi ini.
4. Untuk teman – teman kelas IF-12 yang sudah banyak membantu saat berkuliah di Universitas Amikom Yogyakarta.

Dan seluruh pihak yang tidak dapat saya sebutkan satu per satu, terimakasih atas segala bantuan dan doanya sehingga saya dapat menyelesaikan skripsi ini dengan baik.

## KATA PENGANTAR

*Assalamualaikum Warahmatullahi Wabarakatuh*

Allhamdulillah Puji dan syukur senantiasa peneliti panjatkan kepada Allah SWT, karena berkat nikmat, rahmat, dan pertolongan-Nya peneliti dapat menyelesaikan laporan skripsi ini dengan baik. Laporan skripsi yang dibuat untuk memenuhi syarat memperoleh gelar kesarjanaan Strata-1 (S1) jurusan Informatika Universitas AMIKOM Yogyakarta diharapkan bisa menjadi salah satu referensi pembuatan skripsi di Universitas AMIKOM Yogyakarta serta dapat memberikan penambahan ide yang dapat dikembangkan dimasa depan.

Atas segala bantuan serta amal baik semua pihak diatas, semoga mendapat ridlo Allah SWT. Penulis sangat menyadari bahwa penulisan tugas akhir ini masih kurang sempurna mengingat kurangnya kemampuan dan pengetahuan penulis. Oleh karena itu, saran dan ktirik yang membangun dari pembaca sangat penulis harapkan demi kesempurnaan dan kebaikan tugas akhir ini. Penulis berharap bahwa penulisan tugas akhir ini dapat bermanfaat bagi penulis, pembaca, maupun penelitian di masa depan.

*Wasalamualaikum Warahmatullahi Wabarakatuh*

18 November 2021

Hery Iswanto



## DAFTAR ISI

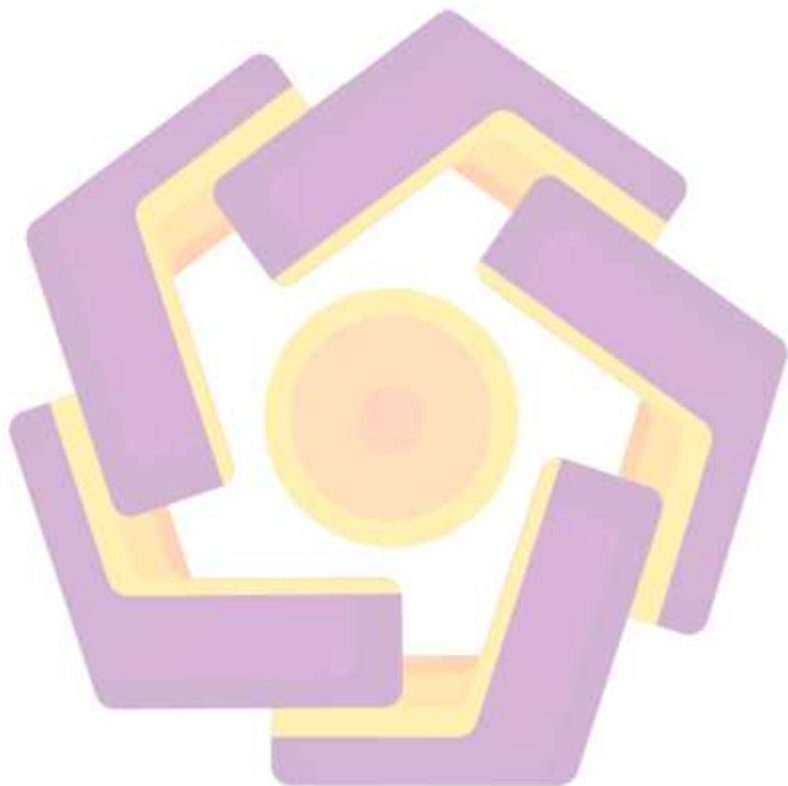
HALAMAN JUDUL .....	ii
LEMBAR PERSETUJUAN .....	iii
LEMBAR PENGESAHAN .....	iv
LEMBAR PENYATAAN .....	v
MOTTO .....	vi
PERSEMBAHAN .....	vii
KATA PENGANTAR .....	viii
DAFTAR ISI .....	ix
DAFTAR TABEL .....	xi
DAFTAR GAMBAR .....	xiii
DAFTAR PERSAMAAN .....	xiv
INTISARI .....	xv
ABSTRACT .....	xvi
<b>BAB 1 PENDAHULUAN .....</b>	<b>1</b>
1.1 Latar Belakang .....	1
1.2 Rumusan Masalah .....	4
1.3 Batasan Maslaah .....	4
1.4 Maksud dan Tujuan Penelitian .....	4
1.5 Metode Penelitian .....	5
1.6 Sistematika Penulisan .....	6
<b>BAB 2 LANDASAN TEORI .....</b>	<b>7</b>
2.1 Kajian Pustaka .....	7
2.2 Email .....	10
2.3 SPAM .....	10
2.4 Machine Learning .....	11
2.5 Text Mining .....	11
2.6 Text Pre-Processing .....	12
2.7 Pembobotan TF-IDF .....	13
2.8 Klasifikasi .....	14
2.9 Naïve Bayes .....	14
2.10 Random Forest .....	17
2.11 K-Nearest Neighbor (KNN) .....	19
2.12 Pengujian Dengan Confusion Matrix .....	20
2.13 Bahasa Pemrograman Phyton .....	22

<b>BAB 3 METODOLOGI PENELITIAN .....</b>	<b>23</b>
3.1 Alat dan Bahan Penelitian.....	23
3.1.1 Alat Penelitian.....	23
3.1.2 Bahan Penelitian.....	24
3.2 Alur Penelitian .....	24
3.2.1 Studi Literatur .....	25
3.2.2 Pengumpulan Data .....	25
3.2.3 Preprocessing Data.....	26
3.2.4 Pembobotan dengan TF-IDF.....	27
3.2.5 Uji Coba dan Perbandingan Metode.....	36
3.2.6 Hasil dan Kesimpulan .....	44
3.2.7 Perancangan Proses Penelitian.....	44
<b>BAB 4 HASIL DAN PEMBAHASAN.....</b>	<b>46</b>
4.1 Implementasi.....	46
4.1.1 Input Dataset .....	46
4.1.2 Preprocessing Data.....	47
4.1.2.1 Cleaning Data.....	47
4.1.2.2 Tokenizing.....	49
4.1.2.3 Stopword Removal / Filtering .....	50
4.1.2.4 Stemming .....	51
4.1.3 Pembobotan TF-IDF .....	52
4.1.4 Data Split.....	53
4.2 Hasil Pengujian.....	54
4.2.1 Hasil Pengujian Naïve Bayes.....	54
4.2.2 Hasil Pengujian <i>Random Forest</i> .....	58
4.2.3 Hasil Pengujian K-Nearest Neighbor (KNN) .....	63
4.3 Analisis Hasil Pengujian.....	67
4.3.1 Accuracy .....	67
4.3.2 Precision, Recall, dan F-Measure Email Ham .....	69
4.3.3 Precision, Recall dan F-measure Email SPAM.....	70
4.3.4 Perbandingan Menggunakan Grafik ROC .....	72
<b>BAB 5 PENUTUP.....</b>	<b>74</b>
5.1 Kesimpulan .....	74
5.2 Saran.....	74

## DAFTAR TABEL

Tabel 2.1 Persamaan dan Perbedaan dengan Penelitian Sebelumnya .....	8
Tabel 2.2 Confusion Matrix .....	21
Tabel 3.1 Spesifikasi Perangkat Keras.....	23
Tabel 3.2 Contoh Data <i>Email</i> .....	28
Tabel 3.3 Proses Pembobotan Setiap Kata Pada Dokumen.....	29
Tabel 3.4 Hasil Proses Dari Metode <i>IDF</i> .....	31
Tabel 3.5 Proses Perhitungan Pembobot Pada Setiap Dokumen .....	33
Tabel 3.6 Hasil Pembobotan Kata Kunci Pada Setiap Dokumen.....	36
Tabel 3.7 Confusion Matrix .....	37
Tabel 3.8 Pembobotan Setiap Kata Pada Dokumen .....	38
Tabel 3.9 Pembobotan Kata Yang Diklasifikasikan Sebagai <i>SPAM</i> .....	39
Tabel 3.10 Pembobotan Kata Yang Diklasifikasikan Sebagai <i>Ham</i> .....	39
Tabel 3.11 Contoh Data Latih dan Data Uji.....	42
Tabel 3.12 Urutan Hasil Perhitungan .....	43
Tabel 4.1 Hyperparameter Value Train Test Split Algoritma Naïve Bayes .....	55
Tabel 4.2 Hyperparameter Algoritma Naïve Bayes.....	55
Tabel 4.3 Confusion Matrix Algoritma Naïve Bayes .....	56
Tabel 4.4 Kinerja Algoritma <i>Naïve Bayes</i> .....	57
Tabel 4.5 Hyperparameter Value Train Test Split Algoritma Random Forest.....	59
Tabel 4.6 Hyperparameter Algoritma Random Forest .....	59
Tabel 4.7 Confusion Matrix Algoritma Random Forest .....	60
Tabel 4.8 Kinerja Algoritma <i>Random Forest</i> .....	61
Tabel 4.9 Hyperparameter Value Train Test Split Algoritma KNN .....	63
Tabel 4.10 <i>Hyperparameter</i> Algoritma <i>KNN</i> .....	64
Tabel 4.11 <i>Confusion Matrix</i> Algoritma <i>KNN</i> .....	64
Tabel 4.12 Kinerja Algoritma <i>KNN</i> .....	65

<b>Tabel 4.13 Hasil Kinerja Akurasi Algoritma <i>Naïve Bayes</i>, <i>Random Forest</i>, dan <i>KNN</i>.....</b>	<b>68</b>
<b>Tabel 4.14 Precision, Recall, dan F-measure Email Ham .....</b>	<b>69</b>
<b>Tabel 4.15 Precision, Recall, dan F-measure Email SPAM .....</b>	<b>70</b>



## DAFTAR GAMBAR

Gambar 2.1 Diagram Tahap Proses Menggunakan Metode <i>Random Forest</i> .....	19
Gambar 3.1 Isi Dataset <i>email</i> yang Belum Melalui <i>Preprocessing</i> .....	26
Gambar 3.2 Alur Proses <i>Pre-processing</i> .....	27
Gambar 3.3 Flowchart Diagram.....	45
Gambar 4.1 Source Code dan Tampilan Dataset.....	47
Gambar 4.2 Source Code dan Hasil Menghilangkan Data yang Tidak Diperlukan .....	48
Gambar 4.3 Source code dan Hasil Dari Proses Pembersihan Data .....	49
Gambar 4.4 Source code dan Hasil Proses Tokenizing Dataset .....	50
Gambar 4.5 Source code dan Hasil Proses Stopword Removal / Filtering pada Dataset .....	51
Gambar 4.6 Source code dan Hasil Proses Stemming pada Dataset.....	52
Gambar 4.7 Source code dan Hasil Proses Pembobotan Dengan <i>TF-IDF</i> pada Dataset .....	53
Gambar 4.8 Source code dan Hasil Proses <i>Data Split</i> pada Dataset .....	53
Gambar 4.9 Perbandingan Akurasi Algoritma <i>Naïve Bayes</i> , <i>Random Forest</i> , dan <i>KNN</i> .....	68
Gambar 4.10 Perbandingan Nilai Accuracy, Precision, dan F-measure Email Ham .....	70
Gambar 4.11 Perbandingan Nilai Accuracy, Precision, dan F-measure Email SPAM .....	71
Gambar 4.12 Hasil Perhitungan dan Grafik Menggunakan AUROC.....	71

## DAFTAR PERSAMAAN

Persamaan 2.1.....	13
Persamaan 2.2.....	13
Persamaan 2.3.....	13
Persamaan 2.4.....	15
Persamaan 2.5.....	16
Persamaan 2.6.....	16
Persamaan 2.7.....	16
Persamaan 2.8.....	16
Persamaan 2.9.....	17
Persamaan 2.10.....	20
Persamaan 2.11.....	21
Persamaan 2.12.....	21
Persamaan 2.13.....	22
Persamaan 2.14.....	22

## INTISARI

*Email* merupakan bukti perkembangan teknologi di dunia dalam menggantikan peran surat konvensional sebagai sarana komunikasi. Fitur dan fasilitas yang diberikan oleh *email* memberikan kenyamanan dan kemudahan untuk pengguna mengirimkan berbagai informasi dalam waktu singkat meskipun terpaut jarak yang jauh. Kemampuan *email* yang dapat mengirim berbagai informasi ke penerima dengan mudah dan tanpa biaya ini sering disalahgunakan oleh berbagai pihak untuk mengirim informasi berisikan iklan produk atau jasa, dan berbagai informasi lainnya yang tidak diinginkan oleh pengguna *email*. Informasi inilah yang sering disebut dengan *email spam*.

Untuk mencegah *email spam*, terdapat banyak metode dalam melakukan penyaringan *email* untuk mengklasifikasikan jenis *email* kedalam *spam* atau *ham*. Pada penelitian ini akan melakukan perbandingan metode *Naïve Bayes*, *Random Forest*, dan *K-Nearest Neighbor* yang bertujuan untuk mencari metode klasifikasi yang baik untuk klasifikasi *email*. Sebelum melakukan proses klasifikasi, penelitian ini menambahkan proses *pre-processing* dan pembobotan dengan *TF-IDF*.

Berdasarkan hasil pengujian yang menggunakan 1998 sampel data, nilai akurasi tertinggi didapatkan pada metode *Naïve Bayes* dan *Random Forest* yaitu sebesar 83,5% diikuti dengan metode *K-Nearest Neighbor* sebesar 82,75%. Dengan penambahan uji coba menggunakan grafik *AUROC*, metode yang memiliki nilai tertinggi didapatkan pada metode *Naïve Bayes* yaitu sebesar 0,859, metode *K-Nearest Neighbor* sebesar 0,845 dan metode *Random Forest* sebesar 0,840. Dari uji coba tersebut membuktikan bahwa metode *Naïve Bayes* merupakan metode terbaik dalam melakukan proses klasifikasi *email spam*.

**Kata Kunci** : *Email spam*, Klasifikasi, *Naïve Bayes*, *Random Forest*, *K-Nearest Neighbor*

## ABSTRACT

*Email is an evidence to technological developments in the world in replace the role of conventional letters as a form of communication. The features and facilities provided by e-mail provide comfort and convenience for users to send various information in a short time even though they are far apart. The ability of email to convey numerous types of information to recipients easily and cheaply is frequently abused by various parties to send information such as product or service advertisements, and as well as other information that email users do not want. This information is often referred to as spam email.*

*To prevent spam emails, there are many methods in filtering emails to classify email types into spam or ham. The study will compare the Methods of Naïve Bayes, Random Forest, and K-Nearest Neighbor which aims to find a good classification method for email classification. Before conducting the classification process, the study added a pre-processing and weighting process with TF-IDF.*

*Based on the results of tests using 1998 sample data, the highest accuracy value was obtained in the Naïve Bayes and Random Forest methods of 83.5% followed by the K-Nearest Neighbor method of 82.75%. With the addition of trials using AUROC charts, the method that has the highest value is obtained in the Naïve Bayes method of 0.859, the K-Nearest Neighbor method of 0.845 and the Random Forest method of 0.840. From the trial proved that the Naïve Bayes method is the best method in the process of classifying spam emails*

**Keyword :** *Spam Email, Classification, Naïve Bayes, Random Forest, K-Nearest Neighbor*