

**KLASIFIKASI *E-MAIL SPAM* DENGAN MENGGUNAKAN
WEKA (*WAIKATO ENVIRONMENT FOR KNOWLEDGE
ANALYSIS*) *EXPLORER***

TUGAS AKHIR



MUHAMMAD FIRZA ARYOGI

1152001027

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER**

UNIVERSITAS BAKRIE

JAKARTA

2020

**KLASIFIKASI *E-MAIL SPAM* DENGAN MENGGUNAKAN
WEKA (*WAIKATO ENVIRONMENT FOR KNOWLEDGE
ANALYSIS*) *EXPLORER***

TUGAS AKHIR

**Diajukan sebagai salah satu syarat untuk memperoleh gelar Sarjana
Komputer**



MUHAMMAD FIRZA ARYOGI

1152001027

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER
UNIVERSITAS BAKRIE**

JAKARTA

2020

HALAMAN PERNYATAAN ORISINALITAS

Tugas Akhir adalah hasil karya saya sendiri,
dan semua sumber baik dikutip maupun
dirujuk telah saya nyatakan dengan benar.

Nama : Muhammad Firza Aryogi

NIM : 1152001027

Tanda Tangan : 

Tanggal : 10 Maret 2020

HALAMAN PENGESAHAN

Tugas Akhir ini diajukan oleh:

Nama : Muhammad Firza Aryogi
NIM : 1152001027
Program Studi : Informatika
Fakultas : Teknik dan Ilmu Komputer
Judul Skripsi : Klasifikasi E-mail Spam dengan Menggunakan WEKA Explorer

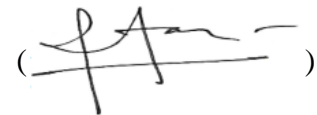
Telah berhasil dipertahankan di hadapan Dewan Penguji dan diterima sebagai bagian dari persyaratan untuk memperoleh gelar Sarjana Komputer pada Program Studi Informatika, Fakultas Teknik dan Ilmu Komputer, Universitas Bakrie.

DEWAN PENGUJI

Pembimbing : Prof. Dr. Hoga Saragih, ST., MT.



Penguji I : Reyful Rey Fatri, M.Sc.



Penguji II : Sigit Wijayanto, M.Sc.



Ditetapkan di : Jakarta

Tanggal : 10 Maret 2020

KATA PENGANTAR

Alhamdulillah, puji dan syukur penulis panjatkan kehadirat Allah SWT karena berkat segala rahmat karunia dan hidayah-Nya penulis diberikan kesehatan untuk menyelesaikan Tugas Akhir yang berjudul “Klasifikasi *E-mail Spam* dengan Menggunakan WEKA (*Waikato Environment for Knowledge Analysis Explorer*)” ini. Shalawat serta salam tidak lupa penulis panjatkan kepada Nabi besar Muhammad SAW.

Tugas Akhir ini dibuat guna memenuhi syarat lulus untuk mendapatkan gelar Sarjana. Selama proses penyusunan Tugas Akhir ini penulis mendapatkan banyak hambatan. Namun hal ini dapat teratasi berkat motivasi, dukungan, serta saran dari berbagai pihak. Untuk itu penulis ingin mengucapkan terimakasih kepada:

- Allah SWT atas limpahan berkah, nikmat sehat, keselamatan dan kelancaran setiap waktu hingga saat ini.
- Kedua orangtua penulis dan keluarga yang senantiasa selalu memberikan do’a, motivasi, serta semangat.
- Bapak Prof. Dr. Hoga Saragih, ST. MT. Ketua Program Studi Teknik Informatika Universitas Bakrie.
- Seluruh tim dosen Program Studi Teknik Informatika Universitas Bakrie yang telah memberikan pembelaran, saran dan arahan.
- Dina Audina, Rafi Afdan, Angga Aulia, Fathur Hadyan Syah, Imam Nugroho, Nandhika Dwi, Bhakti Nurhasan, Ramadhani Qodriansyah
- Ahmad Novel Gadran, Alifian Azmi, Wahyu Widodo, Aziz Sentosa, Alhamsya, Salmaa Badriatu, Prima Dona serta teman-teman Informatika 2015 lainnya yang tidak bisa disebutkan satu persatu yang telah memberikan dukungan, semangat dan kebersamaan dalam suka maupun duka dalam empat tahun perkuliahan.
- Kakak-kakak *senior* dan adik-adik angkatan 2016 dan 2017 Informatika Universitas Bakrie.

- Seluruh pihak Universitas Bakrie baik terlibat langsung maupun tidak yang telah memberikan pengalaman, motivasi, dukungan, dan fasilitas yang sangat membantu selama masa perkuliahan dan penyusunan Tugas Akhir ini.

Penulis sangat bersyukur dan mengucapkan banyak terimakasih untuk pihak0pihak tersebut atas bantuannya dalam bentuk apapun sehingga penulis dapat menyelesaikan Tugas Akhir ini. Semoga Tugas Akhir ini dapat bermanfaat baik bagi semua kalangan pendidikan, khususnya bidang Informatika.

Jakarta, 14 Maret 2020

Penulis

HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI

Sebagai *civitas* akademik Universitas Bakrie, saya yang bertanda tangan di bawah ini :

Nama : Muhammad Firza Aryogi

NIM : 1152002017

Program Studi : Informatika

Fakultas : Teknik dan Ilmu Komputer

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Bakrie **Hak Bebas Royalti Noneksklusif** (*Non-exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

**KLASIFIKASI E-MAIL SPAM DENGAN MENGGUNAKAN WEKA
(WAIKATO ENVIRONMENT FOR KNOWLEDGE ANALYSIS) EXPLORER**

Dengan Hak Bebas Royalti Noneksklusif ini Universitas Bakrie 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 untuk kepentingan akademis.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Jakarta

Pada Tanggal : 10 Maret 2020

Yang menyatakan



Muhammad Firza Aryogi

**KLASIFIKASI *E-MAIL SPAM* DENGAN MENGGUNAKAN WEKA
(*WAIKATO ENVIRONMENT FOR KNOWLEDGE ANALYSIS*) *EXPLORER***

Muhammad Firza Aryogi

ABSTRAK

E-mail adalah salah satu alat komunikasi antar individu maupun perusahaan (*professional*) yang sangat dibutuhkan pada saat sekarang ini. Hal ini disebabkan *e-mail* merupakan alat komunikasi yang dapat mengirim pesan berbentuk digital dengan efisien, cepat dan murah dimana dapat memberikan keuntungan yaitu memangkas waktu dan biaya. Tujuan *e-mail* yaitu sebagai media pertukaran pesan atau informasi digital antar manusia baik individu maupun kelompok. Konten pada *e-mail* berbentuk informasi antar individu dan informasi perusahaan (bisnis). Namun, terdapat informasi lain seperti mengenai promo suatu produk barang maupun jasa, informasi berupa iklan, dan konten yang tidak diinginkan bagi penerima. Hal tersebut merupakan *spam* yang dapat mengganggu kenyamanan para pengguna *e-mail*. Permasalahan tersebut dapat dicegah atau diminimalisir dengan mengklasifikasi data *e-mail* kedalam kategori *spam* dan *non-spam*. Metode klasifikasi yang dipakai untuk menyaring *e-mail* yaitu *naive bayes* dan J48 (*decision tree*) dengan melakukan *training* pada dataset dan melakukan *test* untuk menentukan prediksi kelas pada data *e-mail* tersebut.

Kata Kunci : *E-mail, Spam Detection, Spam Mail Filtration, Naive Bayes, Decision Tree, Klasifikasi*

**E-MAIL SPAM CLASSIFICATION USING WEKA (WAIKATO
ENVIRONMENT FOR KNOWLEDGE ANALYSIS) EXPLORER**

Muhammad Firza Aryogi

ABSTRACT

Today, e-mail is a crucial and necessary tool either for individuals (private) or for professionals (company). E-mail is an effective tool for communication as it saves a lot of time and cost. The purpose of e-mail as a media for exchanging digital message or information between human either individual or group. Contents of e-mails is information between individual or information between company (business). But, there is another information such as information regarding promotion of a goods and services product, advertisements, and an unsolicited messages which is unwanted by the recipient. That is a spam that can interfere e-mail user's convenience. This problem could be prevented or could be minimized by classifying the data of e-mails into spam or non-spam category. Classifier method that can be used for spam e-mail filtration are naive bayes and J48 (decision tree) by performing a training on the dataset, and do the test after it to determine the prediction of the class from the data.

Keywords: E-mail, Spam Detection, Spam Mail Filtration, Naive Bayes, Decision Tree, Classification.

DAFTAR ISI

ABSTRAK	vii
ABSTRACT	ii
DAFTAR ISI	iii
DAFTAR GAMBAR	v
DAFTAR TABEL	vii
DAFTAR RUMUS	viii
DAFTAR SINGKATAN	ix
BAB I PENDAHULUAN	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah	2
1.3. Batasan Masalah	2
1.4. Tujuan Penelitian	3
1.5. Sistematika Penulisan	3
BAB II TINJAUAN PUSTAKA	5
2.1 Penelitian Terkait	5
2.2 E-Mail	11
2.3 Spam Mail	11
2.4 Supervised Learning	14
BAB III METODOLOGI PENELITIAN	21
3.1 Kerangka Penelitian	21
3.2 Studi Literatur	22
3.3 Pembentukan Dataset	23
3.4 Implementasi Algoritma Naive Bayes & J48 (C4.5)	23
3.5 Merumuskan Masalah Penelitian	24
3.6 Time Table	25
BAB IV HASIL & PEMBAHASAN	26
4.1 Pengumpulan Data	26
4.2 Pre-proses Data	28
4.3 Klasifikasi Naive Bayes	30
4.4 Klasifikasi J48	37

4.5	Evaluasi	47
BAB V PENUTUP		52
5.1	Kesimpulan	52
5.2	Saran	52
DAFTAR PUSTAKA		53

DAFTAR GAMBAR

Gambar 2. 1 Supervised Learning Types	15
Gambar 2. 2 Konsep <i>Decision Tree</i>	16
Gambar 3. 1 Alur Pengerjaan Penelitian	21
Gambar 3. 2 Flowcart Implementasi <i>classifier</i>	22
Gambar 4. 1 <i>Raw Data</i>	26
Gambar 4. 2 Contoh <i>e-mail ham</i>	27
Gambar 4. 3 Contoh <i>e-mail spam</i>	27
Gambar 4. 4 <i>Command-line WEKA</i>	28
Gambar 4. 5 Setelah <i>TextDirectoryLoader</i>	29
Gambar 4. 6 File <i>arff</i>	29
Gambar 4. 7 <i>Input file arff</i>	31
Gambar 4. 8 <i>Input file arff</i>	32
Gambar 4. 9 Hasil <i>summary training set</i>	33
Gambar 4. 10 Hasil <i>summary training split 66%</i>	34
Gambar 4. 11 Hasil <i>summary training split 80%</i>	35
Gambar 4. 12 Data <i>test</i> baru.....	36
Gambar 4. 13 Opsi untuk <i>test set</i>	36
Gambar 4. 14 Hasil <i>summary test set</i>	37
Gambar 4. 15 <i>Input file arff (J48)</i>	38
Gambar 4. 16 <i>Input file arff (J48)</i>	39
Gambar 4. 17 <i>Evaluate</i> pada <i>training data (J48)</i>	40
Gambar 4. 18 Hasil <i>summary training set</i>	40
Gambar 4. 19 Hasil <i>summary training split 66%</i>	41
Gambar 4. 20 Hasil <i>summary training split 80%</i>	42
Gambar 4. 21 Jumlah daun dan ukuran pohon pada <i>training split 90%</i>	43
Gambar 4. 22 Data <i>test</i> baru (J48).....	44
Gambar 4. 23 Opsi untuk <i>test set (J48)</i>	44
Gambar 4. 24 Hasil <i>summary test set (J48)</i>	45

Gambar 4. 25 Hasil prediksi data <i>test</i> (J48).....	46
Gambar 4. 26 Jumlah pohon dan daun dari hasil prediksi (J48).....	46
Gambar 4. 27 Hasil Perbandingan <i>Training</i> Data.....	49
Gambar 4. 28 Hasil Perbandingan <i>Test</i> Data.....	51

DAFTAR TABEL

Tabel 2. 1 Rangkuman Penelitian	7
Tabel 2. 2 <i>Confusion Matrix</i>	19
Tabel 3. 1 <i>Time Table</i>	25
Tabel 4.1 Hasil Model Data <i>Text</i> Pada <i>Training</i>	47
Tabel 4.1 Hasil Model Data <i>Text</i> Pada <i>Testing</i>	49

DAFTAR RUMUS

Persamaan 2.1	Rumus <i>Entropy</i>
Persamaan 2.2	Rumus <i>Gain</i>
Persamaan 2.3	Rumus <i>Naive Bayes</i>
Persamaan 2.4	Rumus Penyesuaian <i>Naive Bayes</i>
Persamaan 2.5	Rumus <i>Posterior</i>
Persamaan 2.6	Rumus <i>Recall</i>
Persamaan 2.7	Rumus <i>Precision</i>
Persamaan 2.8	Rumus <i>Accuracy</i>
Persamaan 4.1	Rumus <i>Recall</i> pada <i>Training Naive Bayes</i>
Persamaan 4.2	Rumus <i>Precision</i> pada <i>Training Naive Bayes</i>
Persamaan 4.3	Rumus <i>Accuracy</i> pada <i>Training Naive Bayes</i>
Persamaan 4.4	Rumus <i>Recall</i> pada <i>Training J48</i>
Persamaan 4.5	Rumus <i>Precision</i> pada <i>Training J48</i>
Persamaan 4.6	Rumus <i>Accuracy</i> pada <i>Training J48</i>
Persamaan 4.7	Rumus <i>Recall</i> pada <i>Testing Naive Bayes</i>
Persamaan 4.8	Rumus <i>Precision</i> pada <i>Testing Naive Bayes</i>
Persamaan 4.9	Rumus <i>Accuracy</i> pada <i>Testing Naive Bayes</i>
Persamaan 4.10	Rumus <i>Recall</i> pada <i>Testing J48</i>
Persamaan 4.11	Rumus <i>Presicion</i> pada <i>Testing J48</i>
Persamaan 4.12	Rumus <i>Accuracy</i> pada <i>Testing J48</i>

DAFTAR SINGKATAN

SMD	<i>Spam Mail Detection</i>
KDD	<i>Knowledge Discovery from Databases</i>
ID3	<i>Iterative Dichotomiser 3</i>
IDS	<i>Intrusion Detection System</i>
SVM	<i>Support Vector Machine</i>
CHAID	<i>Chi-squared Automatic Interaction Detection</i>
CART	<i>Classification and Regression Tree</i>
ADABOOST	<i>Adaptive Boost</i>
TF-ISF	<i>Term Frequency – Inverse Sentence Frequency</i>
RDBMS	<i>Relational Database Management System</i>
CLI	<i>Command-Line Interface</i>
CSV	<i>Comma-Separated Values</i>