

**PERBANDINGAN LOGISTIC REGRESSION DAN RANDOM FOREST
MENGUNAKAN CORRELATION-BASED FEATURE SELECTION
UNTUK DETEKSI WEBSITE PHISING**

SKRIPSI



disusun oleh

Farida

17.11.1708

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2021**

**PERBANDINGAN LOGISTIC REGRESSION DAN RANDOM FOREST
MENGUNAKAN CORRELATION-BASED FEATURE SELECTION
UNTUK DETEKSI WEBSITE PHISING**

SKRIPSI

untuk memenuhi sebagian persyaratan
mencapai gelar Sarjana
pada Program Studi Informatika



disusun oleh

Farida

17.11.1708

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2021**

PERSETUJUAN

SKRIPSI

PERBANDINGAN LOGISTIC REGRESSION DAN RANDOM FOREST MENGUNAKAN CORRELATION-BASED FEATURE SELECTION UNTUK DETEKSI WEBSITE PHISING

yang dipersiapkan dan disusun oleh

Farida

17.11.1708

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 29 Agustus 2021

Dosen Pembimbing,

All Mustopa, M.Kom.
NIK. 190302192

PENGESAHAN

SKRIPSI

PERBANDINGAN LOGISTIC REGRESSION DAN RANDOM FOREST MENGUNAKAN CORRELATION-BASED FEATURE SELECTION UNTUK DETEKSI WEBSITE PHISING

yang dipersiapkan dan disusun oleh

Farida

17.11.1708

telah dipertahankan di depan Dewan Penguji
pada tanggal 18 November 2021

Susunan Dewan Penguji

Nama Penguji

All Mustopa, M.Kom
NIK. 190302192

Tanda Tangan

Rizqi Sukma Kharisma, M.Kom
NIK. 190302215

Akhmad Dahlan, M.Kom
NIK. 190302174

Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 08 Desember 2021.

DEKAN FAKULTAS ILMU KOMPUTER

Hanif Al Fatta, S.Kom., 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.

Banjarnegara, 04 Desember 2021

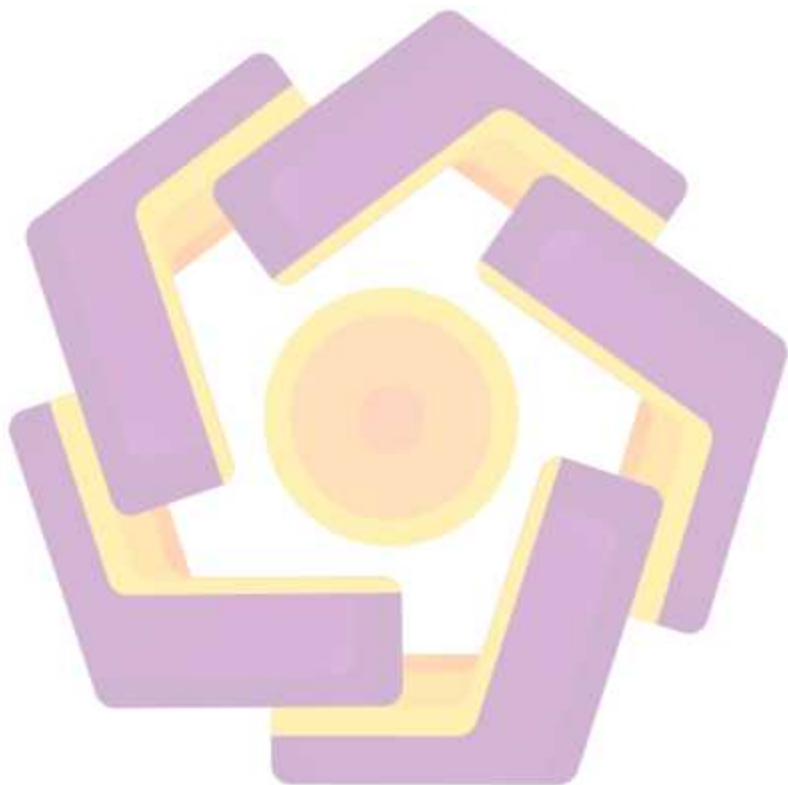


Farida
NIM. 17.11.1708

MOTTO

“Waktu bagaikan pedang. Jika kamu tidak memanfaatkannya dengan baik, maka ia akan memanfaatkanmu.”

(HR. Muslim)



PERSEMBAHAN

Alhamdulillah rabbil'aalamiin

Puji syukur kepada Allah SWT atas segala rahmat dan hidayah yang telah memberikan kelancaran, kemudahan, kesehatan, dan kesabaran dalam menyusun skripsi ini. Pada kesempatan ini penulis ingin mempersembahkan karya penelitian ini sekaligus mengucapkan terima kasih yang sebesar-besarnya kepada:

1. Kedua orang tua, Bapak Darno Nasrudin dan Ibu Parsiah, atas segala doa, dukungan dan kasih sayang yang tidak ada hentinya diberikan. Terima kasih atas segala pengorbanan, nasihat dan bantuan baik secara moril maupun materi hingga saat ini dan seterusnya.
2. Kedua adik saya, Rehana dan Dawam yang selalu menemani, memberikan doa dan semangat hingga penelitian skripsi ini dapat terselesaikan.
3. Untuk seluruh keluarga saya yang selalu memberikan doa dan dukungan.
4. Bapak Ali Mustopa, M.Kom selaku dosen pembimbing, terimakasih banyak atas bimbingannya sehingga skripsi ini dapat terselesaikan.
5. Untuk mas Rafi dan Anggun sebagai support system yang selalu ada dan membantu dalam proses skripsi ini sehingga mendapatkan hasil yang terbaik.

Seluruh pihak yang tidak dapat disebutkan satu per satu, terimakasih atas segala bantuannya dan do'anya sehingga terselesaikan skripsi ini.

KATA PENGANTAR

Assalamu'alaikum Warahmatullaahi Wabarakaatuh

Dengan rahmat Allah SWT Yang Maha Pengasih dan Penyayang, puji syukur penulis panjatkan kehadirat Allah SWT yang telah melimpahkan rahmat, karunia, dan hidayah-Nya sehingga penulis dapat menyelesaikan skripsi yang berjudul "Perbandingan Logistic Regression Dan Random Forest Menggunakan Correlation-Based Feature Selection Untuk Deteksi Website Phising" dengan baik dan lancar.

Skripsi ini disusun sebagai salah satu syarat kelulusann program sarjana jurusan Informatika, Fakultas Ilmu Komputer, Universitas Amikom Yogyakarta. Penulis menyadari bahwa skripsi ini tidak mungkin terselesaikan tanpa aadanya dukungan, bantuan, bimbingan, dan nasihat dari berbagai pihak selama penyusunan skripsi ini. Pada kesempatan ini penulis menyampaikan terima kasih setulus-tulusnya kepada:

1. Bapak Prof. Dr. M. Suyanto, M.M. selaku Rektor Universitas AMIKOM Yogyakarta.
2. Bapak Hanif Al Fatta, S.Kom., M.Kom. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
3. Bapak Ali Mustopa, M.Kom selaku dosen pembimbing penulis yang selalu memberikan arahan dan bimbingan sehingga skripsi ini dapat diselesaikan dengan baik.
4. Bapak Rizqi Sukma Kharisma, M.Kom selaku dosen penguji yang telah memberikan kritik dan saran untuk perbaikan skripsi ini.
5. Bapak Akhmad Dahlan, M.Kom selaku dosen penguji yang telah memberikan masukan dalam skripsi ini.
6. Bapak Ibu dosen Universitas AMIKOM Yogyakarta yang telah memberikan bekal ilmu dan motivasi selama perkuliahan kepada penulis.
7. Para penulis sumber bacaan, jurnal, dan makalah yang penulis jadikan referensi dalam penulisan laporan skripsi ini.
8. Kedua orang tua, keluarga besar dan teman-teman yang senantiasa memberikan doa, dukungan dan semangat kepada penulis dalam menyelesaikan skripsi dan perkuliahan.

9. Sahabat-sahabatku ketika kuliah Anggun, Hafifah, Fay, Rizqi, Sulis, Tito, Rangga, Seyma, Fatimah dan semua anak IF12 2017 yang tidak bisa disebutkan semuanya.
10. Seluruh pihak yang telah membantu baik secara langsung maupun tidak langsung yang tidak bisa penulis sebutkan satu persatu.

Penulis mengucapkan terima kasih atas segala do'a dan dukungannya serta mohon maaf yang sebesar-besarnya. Semoga segala bantuan serta amal baik semua pihak diatas mendapat balasan yang setimpal dari Allah SWT. Penulis menyadari bahwa masih ada banyak kekurangan di dalam laporan ini. Namun penulis berharap laporan skripsi ini dapat memberikan manfaat kepada pembaca.

Akhir kata, penulis berharap laporan skripsi ini dapat berguna dan bermanfaat sebagai bahan kajian untuk mahasiswa yang akan melakukan penelitian selanjutnya.

Wassalaamu'alaikum warahmatullaahi wabarakaatuh

Banjarnegara, 04 Desember 2021

Penulis

Farida

DAFTAR ISI

JUDUL.....	ii
PERSETUJUAN.....	iii
PENGESAHAN.....	iv
PERNYATAAN.....	v
MOTTO.....	vi
PERSEMBAHAN.....	vii
KATA PENGANTAR.....	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR.....	xv
INTISARI.....	xvi
ABSTRACT.....	xvii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	3
1.3 Batasan Masalah.....	4
1.4 Maksud dan Tujuan Penelitian.....	4
1.5 Metode Penelitian.....	4
1.6 Sistematika Penulisan.....	5
BAB II LANDASAN TEORI.....	7
2.1 Kajian Pustaka.....	7
2.2 Web Phishing.....	9
2.3 Klasifikasi.....	11
2.4 Logistic Regression.....	13
2.5 Random Forest.....	17

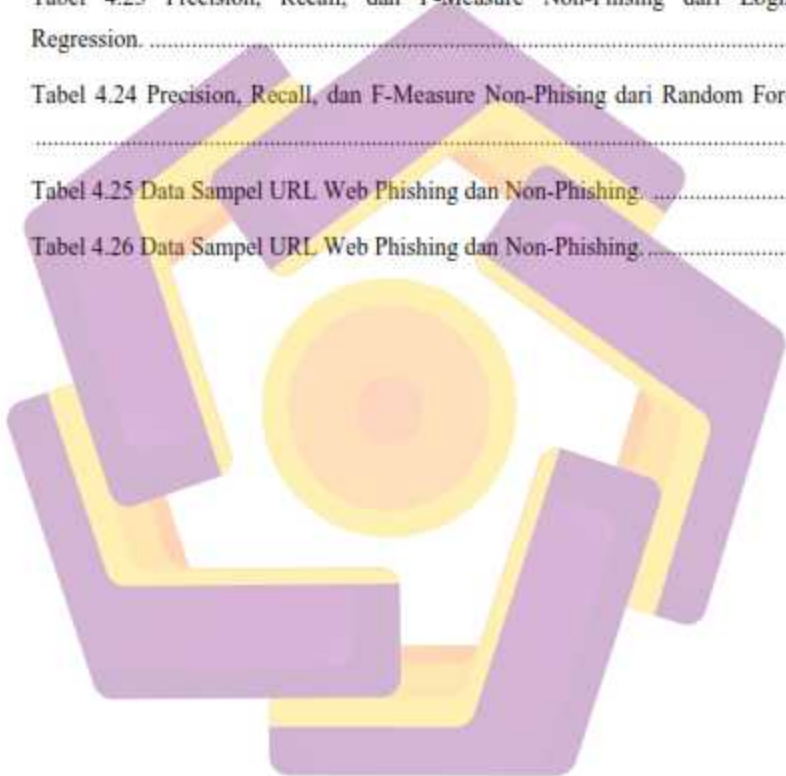
2.6	Seleksi Fitur	18
2.7	Correlation-Based Feature Selection	20
2.8	Confusion Matrix	21
2.9	Bahasa Pemrograman Python	22
BAB III METODE PENELITIAN.....		23
3.1	Alat dan Bahan Penelitian.....	23
3.2	Alur penelitian	24
3.2.1	Studi Literatur.....	24
3.2.2	Pengumpulan Data.....	25
3.2.3	<i>Pre-Processing</i>	28
3.2.4	Uji coba	29
3.2.5	Implementasi Sistem	29
3.2.6	Hasil dan Kesimpulan.....	29
BAB IV HASIL DAN PEMBAHASAN.....		31
4.1	Hasil Pengujian	31
4.1.1	Hasil Pengujian tanpa Correlation-Based Feature Selection.....	31
4.1.1.1	Hasil Pengujian Logistic Regression tanpa CFS	31
4.1.1.2	Hasil Pengujian Random Forest tanpa CFS.....	35
4.1.2	Hasil Pengujian dengan Correlation-Based Feature Selection	38
4.1.2.1	Seleksi fitur dengan Correlation-Based Feature Selection	39
4.1.2.2	Hasil Pengujian Logistic Regression dengan CFS	42
4.1.2.3	Hasil Pengujian Random Forest dengan CFS.....	46
4.2	Analisis Hasil Pengujian	49
4.2.1	<i>Accuracy</i> (Akurasi).....	50
4.2.2	<i>Precision, Recall, dan F-Measure</i> Phising	51
4.2.2.1	Logistic Regression Phising.....	51

4.2.2.2	Random Forest Phising.....	52
4.2.3	<i>Precision, Recall, dan F-Measure</i> Non-Phising.....	52
4.2.3.1	Logistic Regression Non-Phising.....	53
4.2.3.2	Random Forest Non-Phising.....	53
4.2.4	Pemilihan Algoritma Klasifikasi.....	54
4.3	Sistem Deteksi Web Phishing.....	55
4.3.1	Uji Coba Sistem Deteksi Website Phishing.....	57
4.3.2	Hasil Uji Coba Sistem Deteksi Website Phishing.....	58
BAB V PENUTUP		60
5.1	Kesimpulan.....	60
5.2	Saran.....	61
DAFTAR PUSTAKA		63

DAFTAR TABEL

Tabel 2.1 Perbandingan Penelitian.....	8
Tabel 2.2 Confusion Matrix.....	21
Tabel 3.1 Atribut dan Label Klasifikasi Website Phising.....	26
Tabel 4.1 Hyperparameter Value Train Test Split Algoritma Logistic Regression.	32
Tabel 4.2 Hyperparameter Algoritma Logistic Regression.....	32
Tabel 4.3 Confusion Matrix Algoritma Logistic Regression.....	32
Tabel 4.4 Kinerja Algoritma Logistic Regression.....	33
Tabel 4.5 Hyperparameter Value Train Test Split Algoritma Random Forest...35	
Tabel 4.6 Hyperparameter Algoritma Random Forest.....	35
Tabel 4.7 Confusion Matrix Algoritma Random Forest.....	36
Tabel 4.8 Kinerja Algoritma Random Forest.....	37
Tabel 4.9 Hyperparameter Variance Threshold.....	39
Tabel 4.10 Nilai CFS Setiap Fitur.....	40
Tabel 4.11 Hasil Seleksi Fitur Menggunakan CFS.....	41
Tabel 4.12 Hyperparameter Value Train Test Split Algoritma Logistic Regression	43
Tabel 4.13 Hyperparameter Algoritma Logistic Regression.....	43
Tabel 4.14 Confusion Matrix Algoritma Logistic Regression.....	44
Tabel 4.15 Kinerja Algoritma Logistic Regression.....	44
Tabel 4.16 Hyperparameter Value Train Test Split Algoritma Random Forest ..46	
Tabel 4.17 Hyperparameter Algoritma Random Forest.....	46
Tabel 4.18 Confusion Matrix Algoritma Random Forest.....	47
Tabel 4.19 Kinerja Algoritma Random Forest.....	48

Tabel 4.20 Hasil Kinerja Akurasi LR dan RF Sebelum dan Sesudah Diterapkan CFS	50
Tabel 4.21 Precision, Recall, dan F-Measure Phising dari Logistic Regression... ..	51
Tabel 4.22 Precision, Recall, dan F-Measure Phising dari Random Forest.	52
Tabel 4.23 Precision, Recall, dan F-Measure Non-Phising dari Logistic Regression.	53
Tabel 4.24 Precision, Recall, dan F-Measure Non-Phising dari Random Forest.	53
Tabel 4.25 Data Sampel URL Web Phishing dan Non-Phishing	57
Tabel 4.26 Data Sampel URL Web Phishing dan Non-Phishing	59



DAFTAR GAMBAR

Gambar 2.1 Proses Klasifikasi (Jiawei Han, 2011)	13
Gambar 3.1 Alur penelitian.....	30
Gambar 4.1 Mengecek dan menghapus fitur duplikat	39
No table of figures entries found.	
Gambar 4.2 Mencari dan menghapus fitur berkorelasi.....	41
Gambar 4.3 Perbandingan Akurasi Sebelum dan Sesudah Diterapkannya CFS.	50
Gambar 4.4 Perbandingan Nilai Kinerja Logistic Regression Phising.....	51
Gambar 4.5 Perbandingan Nilai Kinerja Random Forest Phising.....	52
Gambar 4.6 Perbandingan Nilai Kinerja Logistic Regression Non-Phising.....	53
Gambar 4.7 Perbandingan Nilai Kinerja Random Forest Non-Phising.....	54
Gambar 4.8 Interface dari Sistem Deteksi Web Phishing.....	55
Gambar 4.9 Hasil Indikasi Bukan Web Phising	56
Gambar 4.10 Hasil Indikasi Web Phising.....	56
Gambar 4.11 Alur Proses Sistem Deteksi Web Phishing	57

INTISARI

Dunia pada saat ini tengah mengalami perkembangan pada bidang teknologi dan komunikasi secara massif, terlebih pada masa pandemi saat ini yang mengharuskan kita semua belajar bahkan bekerja secara daring. Hal ini yang memicu banyaknya kejahatan di dunia internet. Salah satunya yaitu mencuri data pengguna internet melalui sebuah website palsu yang dibangun seperti asli atau disebut juga website phishing.

Pada penelitian ini, untuk mengatasi maraknya website phishing di dunia maya, maka diperlukan suatu model klasifikasi untuk mendeteksi website yang terindikasi phishing dengan menggunakan kinerja terbaik dari salah satu algoritma klasifikasi logistic regression dan random forest. Untuk meningkatkan kinerja algoritma klasifikasi dilakukan seleksi fitur menggunakan metode correlation-based feature selection (CFS) untuk menyeleksi atribut yang paling berpengaruh dalam mendeteksi web phishing.

Berdasarkan hasil uji coba, penerapan algoritma klasifikasi logistic regression dan random forest dalam klasifikasi web phishing dihasilkan akurasi sebesar 93,035 % dan 96,834 %, setelah dilakukan seleksi fitur dengan CFS akurasi yang dihasilkan menjadi 92,718 % dan 97,015 %. Dari uji coba terjadi peningkatan akurasi pada random forest sebesar 0,181 % dan terjadi penurunan yang tidak signifikan pada logistic regression. Hasil uji coba membuktikan bahwa seleksi fitur dengan CFS dapat menghilangkan atribut redundan dan dihasilkan akurasi algoritma klasifikasi yang tidak jauh berbeda ketika atribut lengkap.

Kata Kunci: website phishing, klasifikasi, logistic regression, random forest, correlation-based feature selection (CFS).

ABSTRACT

The world is currently experiencing mass developments of information technology, especially during the current pandemic which requires all of us to learn and even work online. This is what triggers a lot of crime in the internet world. One of them is stealing internet user data through a fake website that is built like the original or also called a phishing website.

In this research, in order to overcome the rise of phishing websites in cyberspace, a classification model is needed to detect phishing websites by using the best performance from one of the logistic regression and random forest classification algorithms. To improve classification performance, feature selection is carried out using the correlation-based feature selection (CFS) method to select the most influential attribute in detecting web phishing.

Based on the test results, the application of the logistic regression and random forest classification algorithm in the classification of web phishing resulted in an accuracy of 93.035% and 96.834%, after feature selection with CFS the resulting accuracy was 92.718% and 97.015%, respectively. From the trial, there was an increase in accuracy in random forest by 0.181% and an insignificant decrease in logistic regression. The test results prove that feature selection with CFS can eliminate redundant attributes and the resulting classification algorithm accuracy is not much different when the attributes are complete.

Keyword: *website phishing, classification, logistic regression, random forest, correlation-based feature selection (CFS).*