

**SELEKSI FITUR BERBASIS INFORMATION GAIN PADA KINERJA  
ALGORITMA KLASIFIKASI UNTUK DETEKSI WEB PHISHING**

**SKRIPSI**



disusun oleh

**Anggun Wahyu Andriyanto**

**17.11.1732**

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

**SELEKSI FITUR BERBASIS INFORMATION GAIN PADA KINERJA  
ALGORITMA KLASIFIKASI UNTUK DETEKSI WEB PHISHING**

**SKRIPSI**

untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Anggun Wahyu Andriyanto**

**17.11.1732**

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

**PERSETUJUAN****SKRIPSI****SELEKSI FITUR BERBASIS INFORMATION GAIN PADA KINERJA  
ALGORITMA KLASIFIKASI UNTUK DETEKSI WEB PHISHING**

yang dipersiapkan dan disusun oleh

**Anggun Wahyu Andriyanto**

**17.11.1732**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 04 Januari 2021

**Dosen Pembimbing,**

**Erni Seniwati, S.Kom., M.Cs.**  
**NIK. 190302231**

**PENGESAHAN****SKRIPSI****SELEKSI FITUR BERBASIS INFORMATION GAIN PADA KINERJA  
ALGORITMA KLASIFIKASI UNTUK DETEKSI WEB PHISHING**

yang dipersiapkan dan disusun oleh

**Anggun Wahyu Andriyanto**

**17.11.1732**

telah dipertahankan di depan Dewan Penguji  
pada tanggal 15 Januari 2021

**Susunan Dewan Penguji**

**Nama Penguji**

**Tanda Tangan**

**Erni Seniwati, S.Kom, M.Cs**  
**NIK. 190302231**

**Lilis Dwi Farida, S.Kom, M.Eng**  
**NIK. 190302288**

**Irma Rofni Wulandari, S.Pd, M.Eng**  
**NIK. 190302329**

Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 26 Januari 2021

**DEKAN FAKULTAS ILMU KOMPUTER**

**Krisnawati, S.Si, M.T**  
**NIK. 190302038**

### 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.

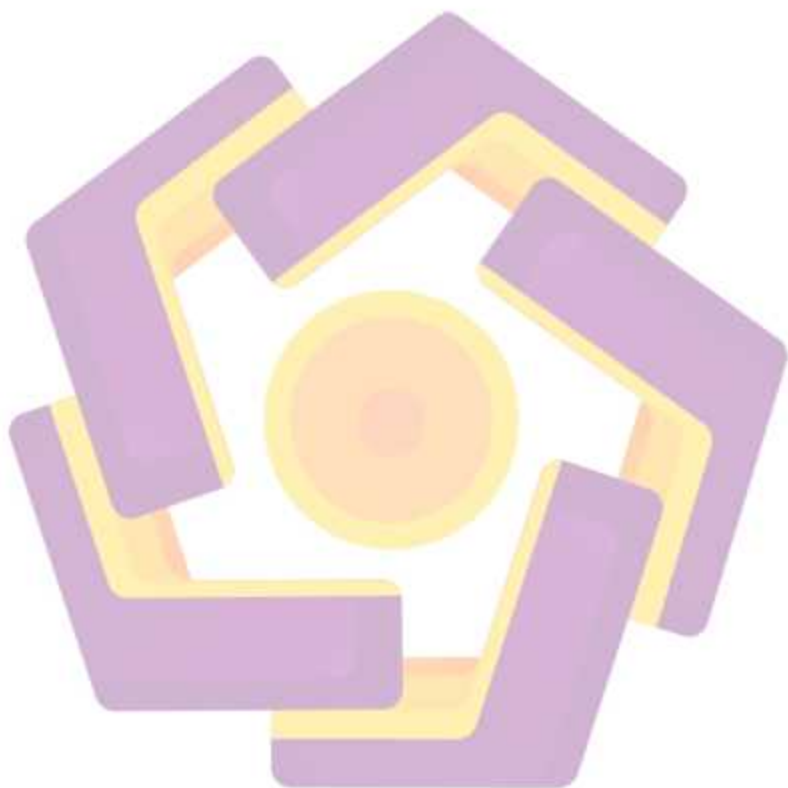
Bojolan, 2 Februari 2021



Anggun Wahyu Anandiyanto  
NIM. 17.11.1732

**MOTTO**

"Apa adanya dan adanya apa"  
(Tara Arts)



## PERSEMBAHAN

Puji syukur penulis panjatkan kepada Allah SWT, yang telah memberikan kesehatan, rahmat dan hidayah sehingga penulis diberikan kesempatan untuk menyelesaikan skripsi ini, sebagai salah satu syarat untuk mendapatkan gelar kesarjana. Walaupun jauh dari kata sempurna, namun penulis bangga telah mencapai pada titik ini, yang akhirnya skripsi ini bisa selesai diwaktu yang tepat.

Skripsi ini saya persembahkan untuk :

1. Bapak dan Ibu, terima kasih atas doa, semangat, pengorbanan, serta kasih sayang yang tidak pernah henti sampai saat ini.
2. Adikku Putri, terima kasih telah menjadi penyemangat dan tempat keluh kesah dalam proses pengerjaan skripsi.
3. Bapak Yoga Pristyanto, terima kasih telah memberikan ide dan gambaran dalam menentukan topik skripsi.
4. Ibu Erni Seniwati yang telah sabar membimbing saya karena sering menghilang tanpa kabar.
5. Warung Barokah atau yang sering disebut Mbok darmi, terima kasih karena dari PPM sampai mengambil skripsi selalu menjadi penyelamat perut dikala kelaparan dengan harga yang membumi.
6. Teman seperjuangan Rizqi dan Sulis dalam kelompok pendeteksi web phishing.
7. Buaji, Gita Mama, Nanda Okto, Mas Dayat, Yoga, Sella, Fery, Novita, Bahrul yang sudah membantu, memberikan masukan dan semangat dalam mengerjakan skripsi ini.
8. Semua teman-teman Informatika 12.
9. Kepada teman-teman, saudara dan dosen yang tidak bisa saya sebutkan satu persatu.



## KATA PENGANTAR

Assalamualaikum Wr. Wb.

Dengan rahmat Allah SWT Yang Maha Pengasih dan Penyayang, puji syukur penulis panjatkan kehadirat Allah SWT yang telah memberikan dan menganugrahkan kasih sayang, rezeki, dan kesehatan serta atas berkah, ridho dan hidayah-Nya, sehingga saya sebagai penulis dapat menyelesaikan skripsi dengan judul "Seleksi Fitur Berbasis Information Gain Pada Kinerja Algoritma Klasifikasi Untuk Deteksi Web Phishing". Penyusunan skripsi ini dibuat sebagai salah satu syarat kelulusan program sarjana jurusan Informatika, Fakultas Ilmu Komputer, Universitas AMIKOM Yogyakarta. Penulis menyadari selama menempuh pendidikan dan proses penyelesaian skripsi ini memperoleh bantuan dan dukungan dari berbagai pihak. Oleh karena itu, pada kesempatan kali ini penulis ingin mengucapkan rasa terima kasih yang mendalam dan tak terkira kepada:

1. Bapak Prof. Dr. M. Suyanto, M.M. selaku rektor Universitas AMIKOM Yogyakarta.
2. Ibu Krisnawati, S.Si, M.T. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
3. Bapak Sudarmawan, M.T. selaku Ketua Program Studi Informatika Universitas AMIKOM Yogyakarta.
4. Ibu Windha Mega Pradnya D, M. Kom. selaku Sekretaris Program Studi Informatika Universitas AMIKOM Yogyakarta.
5. Ibu Erni Seniwati, S.Kom, M.Cs. selaku dosen pembimbing, terima kasih banyak atas bimbingan yang telah diberikan dalam membimbing penulis sehingga skripsi ini dapat diselesaikan dengan baik.
6. Ibu Lilis Dwi Farida, S.Kom, M.Eng. selaku dosen penguji yang telah banyak memberikan masukan dan bimbingan dalam skripsi ini.
7. Ibu Irma Rofni Wulandari, S.Pd, M.Eng. selaku dosen penguji, terima kasih banyak atas bimbingan, kritik dan sarannya untuk perbaikan skripsi ini.
8. Seluruh Bapak dan Ibu dosen Fakultas Ilmu Komputer Jurusan Informatika yang telah memberikan bekal ilmu kepada penulis.



9. Bapak Yoga Pristyanto, S.Kom, M.Eng. selaku dosen mata kuliah Kecerdasan Buatan yang telah memberikan ide sebagai topik pada skripsi ini.
10. Seluruh staff dan karyawan Universitas AMIKOM Yogyakarta.

Bagi seluruh pihak yang tidak bisa penulis sebutkan namanya satu persatu, penulis mengucapkan rasa terima kasih atas segala doa dan dukungannya serta mohon maaf yang sebesar-besarnya. Semoga segala kebaikan, bantuan dana amal baik dari berbagai pihak tersebut diatas mendapat balasan yang setimpal dari Allah SWT dan penulis senantiasa berharap semoga skripsi yang dibuat ini dapat bermanfaat untuk berbagai pihak.

Wassalamualaikum Wr. Wb.

Boyolali, 25 Februari 2021

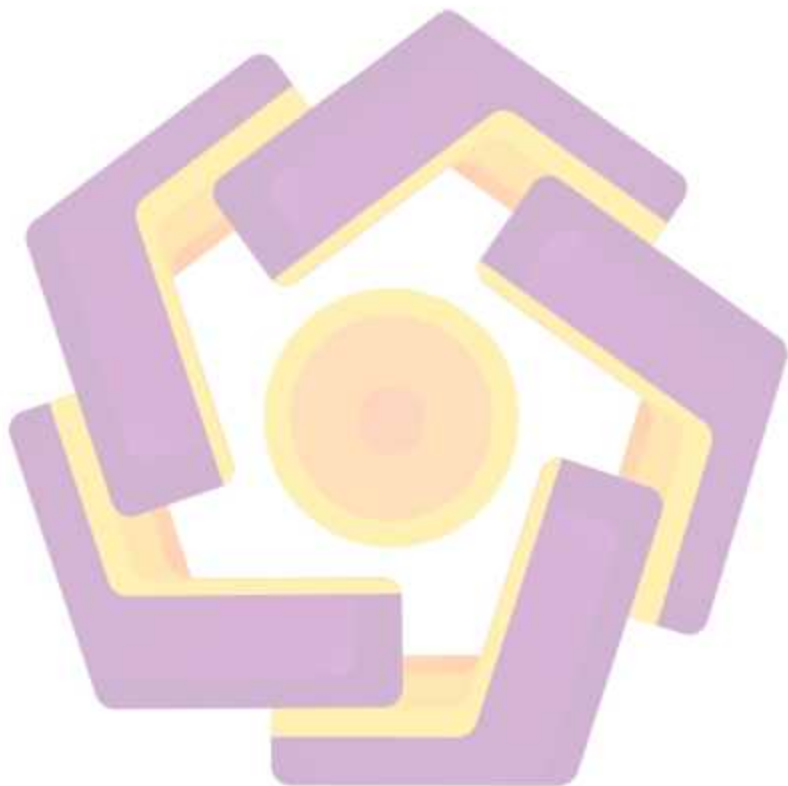
Penulis,  
Anggun Wahyu Andriyanto

## DAFTAR ISI

JUDUL.....	I
PERSETUJUAN.....	III
PENGESAHAN.....	IV
PERNYATAAN.....	IV
MOTTO.....	VI
PERSEMBAHAN.....	VII
KATA PENGANTAR.....	VIII
DAFTAR ISI.....	X
DAFTAR TABEL.....	XIII
DAFTAR GAMBAR.....	XV
<i>INTISARI</i> .....	XVII
<i>ABSTRACT</i> .....	XVII
<b>BAB I PENDAHULUAN</b> .....	<b>1</b>
1.1 LATAR BELAKANG.....	1
1.2 RUMUSAN MASALAH.....	3
1.3 BATASAN MASALAH.....	3
1.4 MAKSUD DAN TUJUAN PENELITIAN.....	4
1.4.1 Maksud Penelitian.....	4
1.4.2 Tujuan Penelitian.....	5
1.5 METODE PENELITIAN.....	5
1.5.1 Metode Pengumpulan Data.....	5
1.5.2 Tahap-Tahap Penelitian.....	5
1.6 SISTEMATIKA PENULISAN.....	6
<b>BAB II LANDASAN TEORI</b> .....	<b>8</b>
2.1 KAJIAN PUSTAKA.....	8
2.2 WEB PHISHING.....	20
2.3 KLASIFIKASI.....	22
2.4 FITUR DETEKSI.....	24
2.5 KINERJA DETEKSI DENGAN KLASIFIKASI.....	25
2.6 SELEKSI FITUR.....	26
2.7 INFORMATION GAIN.....	27
2.8 DECISION TREE.....	28
2.9 CART.....	30
2.10 NAIVE BAYES.....	31
2.11 NEAREST NEIGHBOR.....	33

2.12	BAHASA PEMROGRAMAN PYTHON .....	34
2.13	FRAMEWORK FLASK .....	35
<b>BAB III METODE PENELITIAN .....</b>		<b>25</b>
3.1	ALAT DAN BAHAN PENELITIAN.....	25
3.1.1	<i>Alat</i> .....	25
3.1.1.1	Perangkat Keras.....	25
3.1.1.2	Perangkat Lunak.....	26
3.2	ALUR PENELITIAN.....	26
3.3	STUDI KEPUSTAKAAN .....	27
3.4	PENGUMPULAN DATA .....	28
3.5	FITUR EKSTRAKSI .....	29
3.6	PRE-PROCESSING .....	35
3.7	PERANCANGAN MODEL KLASIFIKASI DETEKSI <i>WEB</i> PHISHING.....	36
3.7.1	<i>Kinerja Information Gain</i> .....	37
3.7.2	<i>Kinerja Model Klasifikasi</i> .....	38
3.8	UJI COBA .....	39
3.9	IMPLEMENTASI SISTEM .....	39
3.10	UJI COBA SISTEM.....	40
<b>BAB IV HASIL DAN PEMBAHASAN .....</b>		<b>42</b>
4.1	HASIL PENGUJIAN .....	42
4.1.1	<i>Hasil Pengujian Klasifikasi Tanpa Information Gain</i> .....	43
4.1.1.1	Hasil Pengujian Decision Tree (CART).....	45
4.1.1.2	Hasil Pengujian K-Nearest Neighbor .....	46
4.1.1.3	Hasil Pengujian Naïve Bayes.....	47
4.1.2	<i>Hasil Pengujian Klasifikasi Dengan Information Gain</i> .....	54
4.1.2.1	Seleksi Fitur dengan Information Gain.....	56
4.1.2.2	Hasil Pengujian Decision Tree (CART).....	61
4.1.2.3	Hasil Pengujian K-Nearest Neighbor .....	65
4.1.2.4	Hasil Pengujian Naïve Bayes.....	69
4.2	ANALISIS HASIL PENGUJIAN .....	70
4.2.1	<i>Accuracy (Akurasi)</i> .....	71
4.2.2	<i>Precision, Recall, dan F-Measure Phishing</i> .....	72
4.2.2.1	Decision Tree (CART) Phishing.....	74
4.2.2.2	K-Nearest Neighbor Phishing.....	76
4.2.2.3	Naïve Bayes Phishing.....	78
4.2.3	<i>Precision, Recall, dan F-Measure Non-Phishing</i> .....	76
4.2.3.1	Decision Tree (CART) Non-Phishing.....	78
4.2.3.2	K-Nearest Neighbor Non-Phishing .....	80
4.2.3.3	Naïve Bayes Non-Phishing.....	81
4.2.4	<i>Pemilihan Algoritma Klasifikasi</i> .....	79
4.3	SISTEM DETEKSI <i>WEB</i> PHISHING.....	79
4.4	UJI COBA SISTEM DETEKSI <i>WEB</i> PHISHING.....	82
4.5	HASIL UJI COBA SISTEM DETEKSI <i>WEB</i> PHISHING.....	83

BAB V PENUTUP.....	86
5.1 KESIMPULAN.....	86
5.2 SARAN.....	87
DAFTAR PUSTAKA .....	1



## DAFTAR TABEL

Tabel 2.1 Indikator URL phishing pada penelitian S. Caroline, Elijah .....	17
Tabel 2.2 Perbandingan Algoritma C4.5 Sebelum dan Sesudah Optimasi.....	19
Tabel 2.3 Persamaan dan Perbedaan Penelitian Sebelumnya.....	19
Tabel 3.1 Tabel Confusion Matrix.....	38
Tabel 4.1 Hyperparameter Value Train Test Split Algoritma Decision Tree (CART) .....	44
Tabel 4.2 Hyperparameter Algoritma Decision Tree (CART).....	44
Tabel 4.3 Confusion Matrix Algoritma Decision Tree (CART).....	45
Tabel 4.4 Kinerja Algoritma Decision Tree (CART).....	46
Tabel 4.5 Hyperparameter Value Train Test Split Algoritma K-Nearest Neighbor .....	48
Tabel 4.6 Hyperparameter Algoritma K-Nearest Neighbor.....	48
Tabel 4.7 Confusion Matrix Algoritma K-Nearest Neighbor.....	48
Tabel 4.8 Kinerja Algoritma K-Nearest Neighbor.....	49
Tabel 4.9 Hyperparameter Value Training Test Split Algoritma Naïve Bayes ....	51
Tabel 4.10 Hyperparameter Algoritma Naïve Bayes.....	52
Tabel 4.11 Confusion Matrix Algoritma Naïve Bayes.....	52
Tabel 4.12 Kinerja Algoritma Naïve Bayes.....	53
Tabel 4.13 Hyperparameter Variance Threshold.....	55
Tabel 4.14 Hyperparameter Information Gain.....	56
Tabel 4.15 Information Gain Tiap-Fitur.....	56
Tabel 4.16 Hasil Seleksi Fitur Menggunakan Nilai Information Gain.....	58
Tabel 4.17 Hyperparameter Value Train Test Split Algoritma Decision Tree (CART).....	60
Tabel 4.18 Hyperparameter Algoritma Decision Tree (CART).....	60
Tabel 4.19 Confusion Matrix Algoritma Decision Tree (CART).....	61
Tabel 4.20 Kinerja Algoritma Decision Tree (CART).....	61



Tabel 4. 21 Hyperparameter Value Train Test Split Algoritma K-Nearest Neighbor .....	64
Tabel 4. 22 Hyperparameter Algoritma K-Nearest Neighbor.....	64
Tabel 4.23 Confusion Matrix Algoritma K-Nearest Neighbor .....	64
Tabel 4.24 Kinerja Algoritma K-Nearest Neighbor.....	65
Tabel 4.25 Hyperparameter Value Train Test Split Algoritma Naïve Bayes .....	67
Tabel 4.26 Hyperparameter Algoritma Naïve Bayes.....	67
Tabel 4.27 Confusion Matrix Algoritma Naïve Bayes .....	68
Tabel 4.28 Kinerja Algoritma Naïve Bayes .....	68
Tabel 4.29 Hasil Kinerja Akurasi Algoritma Klasifikasi Sebelum dan Sesudah Diterapkan Information Gain.....	71
Tabel 4.30 Precision, Recall, dan F-Measure Phishing dari Decision Tree (CART) .....	72
Tabel 4.31 Precision, Recall, dan F-Measure Phishing dari K-Nearest Neighbor	74
Tabel 4.32 Precision, Recall, dan F-Measure Phishing dari Naïve Bayes.....	75
Tabel 4.33 Precision, Recall, dan F-Measure Non-Phishing dari Decision Tree (CART) .....	76
Tabel 4.34 Precision, Recall, dan F-Measure Non-Phishing dari K-Nearest Neighbor.....	77
Tabel 4.35 Precision, Recall, dan F-Measure Non-Phishing dari Naïve Bayes....	78
Tabel 4.36 Data Sampel URL Web Phishing dan Non-Phishing .....	82
Tabel 4.37 Data Sampel URL Web Phishing dan Non-Phishing .....	84

## DAFTAR GAMBAR

Gambar 2.1 Hasil Perbandingan Akurasi Algoritma Klasifikasi Deteksi Web Phishing [6].....	16
Gambar 2.2 Perbandingan akurasi sebelum dan sesudah seleksi fitur.....	18
Gambar 2.3 Proses Klasifikasi Pada Penelitian Jiawei Han [11].....	24
Gambar 3.1 Alur Penelitian.....	27
Gambar 3.2 Seleksi Fitur dan Model Klasifikasi Deteksi Web Phishing .....	37
Gambar 4.1 Perbandingan Akurasi Sebelum dan Sesudah Diterapkannya Information Gain.....	72
Gambar 4.2 Perbandingan Nilai Kinerja Decision Tree (CART) Phishing.....	73
Gambar 4.3 Perbandingan Nilai Kinerja kNN Phishing.....	74
Gambar 4.4 Perbandingan Nilai Kinerja Naïve Bayes Phishing .....	75
Gambar 4.5 Perbandingan Nilai Kinerja Decision Tree (CART) Non-Phishing..	77
Gambar 4.6 Perbandingan Nilai Kinerja kNN Non-Phishing.....	78
Gambar 4.7 Perbandingan Nilai Kinerja Naïve Bayes Non-Phishing .....	79
Gambar 4.8 <i>Interface</i> dari Sistem Deteksi Web Phishing.....	80
Gambar 4.9 Hasil Indikasi Jika Web Dideteksi Sebagai Non-Phishing .....	81
Gambar 4.10 Hasil Indikasi Jika Web Dideteksi Sebagai Phishing.....	81
Gambar 4.11 Alur Proses Sistem Deteksi Web Phishing .....	82



## INTISARI

Pencurian data pengguna dari internet melalui web palsu yang dirancang menyerupai asli disebut web phishing. Web phishing juga terjadi pada media sosial dan internet *banking*. Web phishing tersebut akan mengelabui pengguna internet untuk memasukkan informasi penting seperti *password* ataupun data bank, sehingga pelaku phishing bisa mendapatkan data tersebut.

Pada penelitian ini, untuk mencegah kriminalitas phishing sehingga mengurangi dan bisa menghindari kerugian situs phishing terhadap pengguna internet maka dibutuhkan suatu model klasifikasi untuk memprediksi web yang terindikasi phishing, menggunakan kinerja terbaik dari salah satu algoritma klasifikasi *decision tree* (CART), *k-nearest neighbor*, dan *naïve bayes*. Untuk meningkatkan kinerja dari algoritma klasifikasi maka dilakukan seleksi fitur menggunakan metode *information gain* untuk menyeleksi atribut yang paling berpengaruh dalam mendeteksi web phishing.

Berdasarkan hasil uji coba, penerapan algoritma klasifikasi *decision tree* (CART), *k-nearest neighbor*, dan *naïve bayes* dalam klasifikasi web phishing dihasilkan akurasi sebesar 96,56%, 95,07%, dan 90,14%, setelah dilakukan seleksi fitur dengan metode *information gain* akurasi yang dihasilkan 97,01%, 94,84%, dan 90,09%. Dari uji coba terjadi peningkatan akurasi pada *decision tree* (CART), penurunan yang tidak signifikan bahkan cenderung stabil pada *k-nearest neighbor* dan *naïve bayes* setelah dihilangkannya beberapa fitur. Dari hasil uji coba terbukti bahwa seleksi fitur dengan metode *information gain* mampu menghilangkan atribut redundan, lalu dihasilkan akurasi algoritma klasifikasi yang tidak jauh berbeda dengan data ketika atributnya masih lengkap dan bias diterapkan untuk mendeteksi web phishing.

**Kata Kunci:** web phishing, klasifikasi, deteksi phishing, *decision tree* (CART), *k-nearest neighbor*, *naïve bayes*, *information gain*.

## ABSTRACT

The theft of user data from the internet via a fake web designed to resemble the real thing or called web phishing. Web phishing also occurs in social media and internet banking. The web phishing will trick internet users into entering important information such as passwords or bank data, so that phishers can get the data.

In this study, to prevent phishing crime so as to reduce and avoid loss of phishing sites to internet users, a classification model is needed to predict phishing-indicated webs, using the best performance from one of the decision tree (CART) classification algorithms, *k*-nearest neighbor, and naïve bayes. To improve the performance of the classification algorithm, feature selection is carried out using the information gain method to select the most influential attributes in detecting web phishing.

Based on the test results, the application of the decision tree (CART), *k*-nearest neighbor, and naïve bayes classification algorithm in the web phishing classification resulted in an accuracy of 96.56%, 95.07%, and 90.14%, after selecting features with the information gain method resulted in the accuracy of 97.01%, 94.84%, and 90.09%. From the testing, there was an increase in the accuracy of the decision tree (CART), the decrease was not significant and even tended to be stable in the *k*-nearest neighbor and naïve bayes after removing some features. From the test results, it is proven that feature selection using the information gain method is able to eliminate redundant attributes and the resulting classification algorithm accuracy is not much different when the attributes are complete.

**Keyword:** website phishing, classification, phishing detection, decision tree (CART), *k*-nearest neighbor, naïve bayes, information gain.

