

# BAB I PENDAHULUAN

## 1.1. Latar Belakang

Di zaman digital seperti sekarang, keamanan siber telah menjadi fokus utama bagi individu maupun organisasi. Salah satu ancaman utama adalah serangan *phishing*. *Phishing* adalah suatu tindakan kriminal yang bertujuan untuk mendapatkan informasi rahasia dan penting seseorang, seperti nama *user*, *password* atau kata sandi, dan informasi sensitif lainnya. Tindakan ini dilakukan dengan cara menyediakan situs web palsu serta menyamar sebagai orang atau bisnis tepercaya yang mirip dengan aslinya (Wahyudi et al., 2022). *Phishing* tetap menjadi metode paling umum bagi penjahat siber untuk mendapatkan akses tidak sah ke informasi sensitif, menyumbang sekitar 39,6% dari semua ancaman email (Smith, 2024). Diperkirakan sekitar 3,4 miliar email *phishing* dikirim setiap hari pada tahun 2024 (Michalowski, 2024). Oleh karena itu, deteksi dini terhadap situs web *phishing* sangat penting untuk melindungi pengguna internet dari potensi kehilangan data pribadi atau keuangan.

Algoritma *machine learning* telah terbukti efektif dalam mendeteksi serangan *phishing* dengan menganalisis pola dan fitur yang terkait dengan situs web *phishing*. *Machine learning* adalah kumpulan metode untuk mempelajari pola dari data masa lalu (Sunge, 2022). Pada penelitian ini menggunakan metode *machine learning* untuk klasifikasi data. Klasifikasi data merupakan proses untuk menemukan model atau fungsi yang dapat menjelaskan sekaligus membedakan kelas data beserta konsepnya. (Putri & Wijayanto, 2022). Penelitian ini akan melakukan komparasi beberapa algoritma untuk menemukan akurasi tertinggi. Beberapa algoritma *machine learning* diimplementasikan untuk tujuan ini, termasuk *Naive Bayes*, *Simple Logistic*, *Random Forest*, dan *Decision Tree J48*. *Naive Bayes* dikenal dengan kesederhanaannya dan efisiensi dalam pengolahan data, sementara *Simple Logistic* menawarkan pendekatan linier yang dapat diinterpretasikan dengan baik. Di sisi lain, *Random Forest* dan *Decision Tree J48* menyediakan solusi berbasis pohon keputusan yang kuat dalam menangani data yang kompleks. Meski begitu, belum terdapat kesepakatan mengenai algoritma mana yang memiliki kinerja terbaik dalam mengklasifikasikan situs web sebagai

*phishing* atau aman. Penelitian ini bertujuan untuk membandingkan kinerja empat algoritma *machine learning* yang berbeda, yaitu *Naive Bayes*, *Simple Logistic*, *Random Forest*, dan *Decision Tree J48*, dalam mengklasifikasikan situs web sebagai *phishing* atau *aman* berdasarkan fitur-fiturnya. Kemudian setelah menemukan algoritma terbaik, penelitian ini akan meneliti fitur apa saja yang paling relevan dan dapat sebagai indikator yang penting dalam membedakan website *phishing* dan *aman*. Dengan melakukan perbandingan ini, kita dapat mengevaluasi keunggulan dan kelemahan masing-masing algoritma dalam konteks deteksi *phishing*. Dataset yang digunakan dalam penelitian ini mencakup 11.055 situs web yang diperoleh dari berbagai sumber, termasuk arsip *PhishTank*, arsip *MillerSmiles*, dan operator pencarian Google (R. Mohammad & McCluskey, 2012). Fitur-fitur yang berbeda telah diidentifikasi yang terkait dengan situs web yang sah dan mencurigakan, membentuk dasar untuk analisis dan perbandingan algoritma.

Hasil dari penelitian ini diharapkan dapat memberikan wawasan penting bagi praktisi keamanan siber dan peneliti dalam menentukan algoritma *machine learning* yang paling tepat untuk mendeteksi *phishing*. dan dapat menemukan fitur indikator yang paling relevan untuk mendeteksi *website phishing*. Dengan memahami kinerja relatif dari *Naive Bayes*, *Simple Logistic*, *Random Forest*, dan *Decision Tree J48*, kita dapat meningkatkan efektivitas sistem deteksi *phishing* dan mengurangi risiko serangan terhadap pengguna internet. Seiring dengan meningkatnya kompleksitas serangan *phishing* dan pentingnya deteksi dini, penelitian ini diharapkan dapat memberikan kontribusi yang berarti dalam memperluas pemahaman mengenai kinerja algoritma *machine learning* dalam mengatasi ancaman siber.

## **1.2. Perumusan Masalah**

Berdasarkan latar belakang yang telah disampaikan, rumusan masalah yang akan menjadi fokus utama penelitian ini adalah:

1. Bagaimana proses *preprocessing* data pada dataset untuk mengoptimalkan hasil klasifikasi?

2. Bagaimana perbandingan kinerja antara algoritma *Naive Bayes*, *Simple Logistic*, *Random Forest*, dan *Decision Tree J48* dalam mengklasifikasikan situs web sebagai *phishing* atau *aman*?
3. Apa saja metrik evaluasi untuk pemilihan algoritma yang optimal untuk deteksi situs web *phishing*?
4. Apa saja fitur-fitur yang paling penting dan relevan dalam membedakan situs web *phishing* yang dapat menjadi karakteristik yang sering kali menjadi ciri khas dari situs web *phishing*?
5. Bagaimana metode untuk mengetahui fitur-fitur apa saja yang merupakan karakteristik web *phishing* dari algoritma terbaik?

### 1.3. Batasan Masalah

Berdasarkan kajian terhadap latar belakang permasalahan, penelitian ini akan membatasi ruang lingkupnya pada:

1. Penelitian ini akan membatasi penggunaan empat algoritma *machine learning*, yaitu *Naive Bayes*, *Simple Logistic*, *Random Forest*, dan *Decision Tree J48*. Algoritma lainnya tidak akan dipertimbangkan dalam penelitian ini.
2. Penelitian ini akan menggunakan fitur-fitur atau atribut yang telah diidentifikasi relevan. Namun, hanya fitur-fitur yang tersedia dan dapat diukur dalam dataset yang digunakan yang akan dipertimbangkan. Pengembangan fitur baru tidak akan dilakukan dalam lingkup penelitian ini.
3. Penelitian ini akan menggunakan dataset yang terdiri dari 11.055 situs web yang berasal dari berbagai sumber, termasuk arsip *PhishTank*, arsip *MillerSmiles*, dan operator pencarian *Google*. Namun, penelitian ini tidak akan mengumpulkan atau menggunakan dataset tambahan dari sumber lainnya.
4. Penelitian ini akan menganalisis kinerja setiap algoritma dengan menggunakan metrik standar, antara lain akurasi, presisi, recall, *F1 Score*, *MCC (Matthews Correlation Coefficient)*, dan *ROC AUC*. Namun, penelitian ini tidak akan mempertimbangkan faktor lain seperti waktu komputasi atau *resources* yang diperlukan.

5. Penelitian ini akan menyimpulkan fitur-fitur yang paling relevan dalam klasifikasi situs web *phishing* berdasarkan algoritma terbaik yang telah dipilih.

#### **1.4. Tujuan**

Berdasarkan permasalahan yang telah diidentifikasi, tujuan penelitian ini adalah antara lain:

1. Penelitian ini akan membahas bagaimana metode *preprocessing* digunakan untuk menangani ketidakseimbangan kelas dalam dataset.
2. Penelitian ini bertujuan untuk membandingkan dan mengevaluasi kinerja empat algoritma *machine learning* yang berbeda, yaitu Naive Bayes, Simple Logistic, Random Forest, dan Decision Tree J48, dalam mengklasifikasikan situs web sebagai *phishing* atau *non phishing* berdasarkan fitur-fiturnya.
3. Penelitian ini bertujuan untuk mengevaluasi algoritma Naive Bayes, Simple Logistic, Random Forest, dan Decision Tree J48 menggunakan metrik evaluasi akurasi, presisi, *recall*, *F1 Score*, *MCC*, dan *ROC AUC*.
4. Penelitian ini juga bertujuan untuk mengidentifikasi fitur-fitur yang paling penting dan relevan dalam membedakan situs web *phishing*.
5. Penelitian ini membahas metode untuk mengidentifikasi fitur-fitur yang paling relevan dalam pengklasifikasian situs web *phishing* dengan menggunakan hasil dari algoritma terbaik yang telah diuji.

#### **1.5. Sistematika Penulisan**

##### **a. BAB I Pendahuluan**

Bab ini mencakup latar belakang penelitian, perumusan masalah, batasan masalah, tujuan penelitian, dan sistematika penulisan. Pendahuluan memberikan gambaran umum mengenai penelitian yang dilakukan, termasuk alasan pentingnya topik yang dibahas.

##### **b. BAB II Tinjauan Pustaka**

Pada bab ini, diuraikan berbagai literatur dan penelitian sebelumnya yang relevan dengan topik klasifikasi situs web *phishing*. Selain itu, bab ini juga menjelaskan landasan teori yang digunakan sebagai dasar dari penelitian,

termasuk algoritma yang dikaji, seperti Naive Bayes, Simple Logistic, Random Forest, dan Decision Tree J48.

c. BAB III Metodologi

Bab ini menguraikan tahapan-tahapan penelitian, dimulai dari pengumpulan dataset, deskripsi dataset, preprocessing data menggunakan teknik SMOTE, metode evaluasi dengan K-Fold Cross Validation, hingga penggunaan metrik evaluasi. Prosedur implementasi setiap algoritma juga diuraikan secara rinci.

d. BAB IV Hasil dan Pembahasan

Pada bab ini, disajikan hasil penelitian, termasuk analisis dataset, evaluasi kinerja algoritma, dan identifikasi fitur paling relevan berdasarkan algoritma terbaik. Bab ini juga membahas proses pembuatan sistem prediksi berbasis algoritma machine learning dan uji coba sistem yang telah dirancang.

e. BAB V Penutup

Bab ini menyajikan kesimpulan dari hasil penelitian serta saran untuk pengembangan penelitian di masa mendatang.



**HALAMAN INI SENGAJA DIKOSONGKAN**