

BAB III

METODOLOGI PENELITIAN

3.1 Gambaran Umum Objek Penelitian

Objek penelitian ini berfokus pada deteksi *phishing* pada *URL website* dengan merangkum identifikasi dan klasifikasi *URL* yang berpotensi sebagai situs *phishing*. Proses analisis melibatkan penelitian mendalam terhadap berbagai aspek *URL*, seperti struktur domain, konten halaman *web*, dan karakteristik khusus yang dapat menjadi tanda serangan *phishing*. Penelitian ini bertujuan untuk mengembangkan metode atau model yang dapat secara otomatis dapat membedakan *URL* yang sah dari yang berbahaya, sehingga dapat meminimalkan risiko pengguna terkena serangan *phishing* dan meningkatkan keamanan saat *browsing web*.

Pengidentifikasian pola-pola mencurigakan dalam *URL* menjadi fokus utama, termasuk fitur-fitur yang sering terkait dengan situs *phishing*, seperti *URL* menyesatkan, penggunaan domain palsu, dan karakteristik tertentu dalam kode halaman *web*. Dengan demikian, penelitian ini diarahkan untuk mengembangkan metode deteksi yang efektif dan efisien, yang pada akhirnya akan memberikan perlindungan yang lebih baik kepada pengguna internet dan memperkuat keseluruhan sistem perlindungan keamanan web.

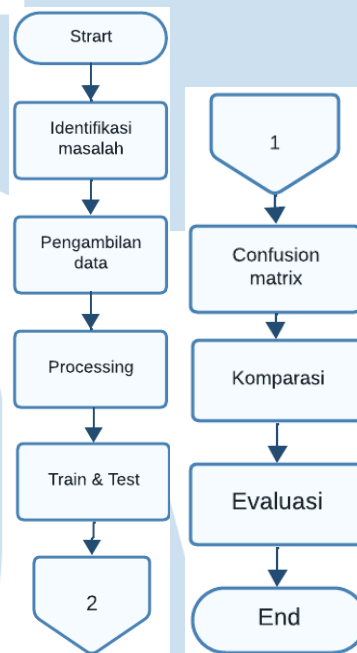
3.2 Metode Penelitian

3.2.1 Alur Penelitian

Alur penelitian adalah rencana kerja yang disusun untuk mengarahkan jalannya penelitian dan memastikan pelaksanaannya terstruktur disebut sebagai alur penelitian. Alur ini bertindak sebagai panduan yang mengatur proses penelitian. Berdasarkan gambar 3.1, penelitian ini dimulai dengan mengidentifikasi permasalahan yang ada dan semakin meningkat presentasinya di Indonesia yaitu mengenai *phishing* pada *url website*, oleh karena itu peningkatan keamanan *url website* menjadi tujuan utama dalam penelitian ini. Setelah mengidentifikasi isu yang

ada, langkah selanjutnya adalah mencari data phishing melalui situs dataset resmi yaitu Kaggle.

Berdasarkan *dataset* yang didapatkan, metode untuk menyelesaikan masalah ditetapkan dengan menggunakan algoritma *Random Forest*, *Support Machine Vectore*, dan *K-Nearest Neighbor* untuk menentukan hasil akurasi terakurat. Penetapan metode atau algoritma ini didasarkan pada studi literatur yang relevan dan disesuaikan dengan kebutuhan dan karakteristik data yang ada. Setelah menentukan pendekatan, proses pengolahan data dimulai dengan langkah-langkah yang diatur berdasarkan kerangka kerja *CRISP-DM*. Proses penelitian ini sampai pada tahap evaluasi yang nantinya dapat dilanjutkan oleh peneliti lain menjadi *website* maupun *mobile app*.



Gambar 3. 1 Flowchart Alur Penelitian

3.2.2 Metode Data Mining

Metodologi yang digunakan dalam penelitian ini mengadopsi model *data mining* dengan pendekatan metode klasifikasi. Dalam pelaksanaannya, penelitian ini menerapkan kerangka kerja *CRISP-DM* yang dikenal sebagai kerangka kerja pengembangan *Knowledge Discovery in Databases (KDD)*. Untuk melakukan analisis klasifikasi penelitian ini

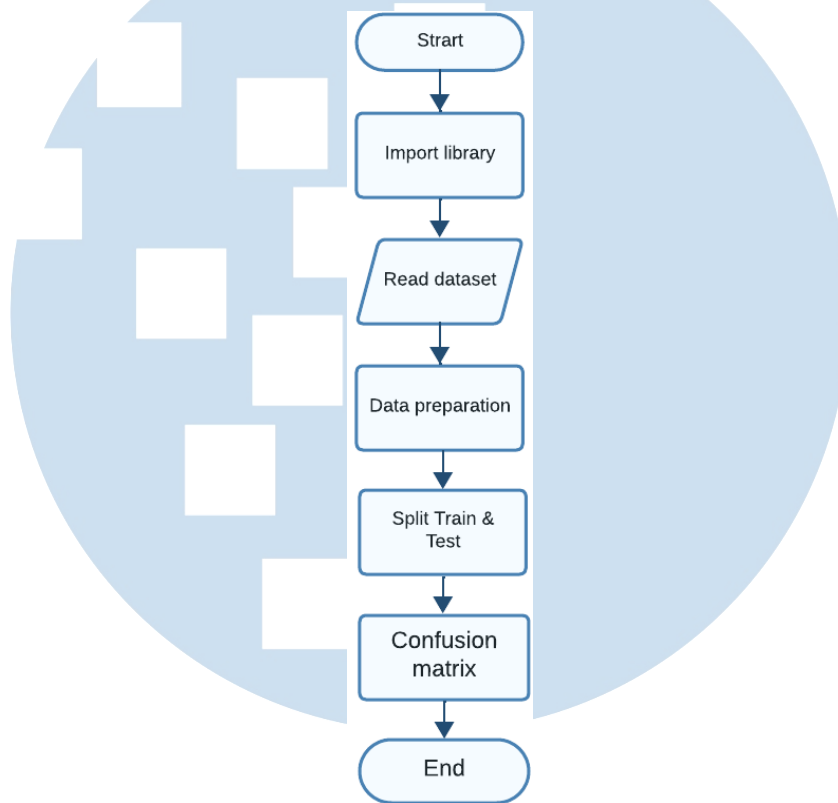
menggunakan *platform* Google Colab sebagai alat utama. *CRISP-DM* dipilih sebagai metode karena pertimbangan terhadap metode *data mining* lainnya.

Tabel 3. 1 Perbandingan *CRISP-DM* dan *KDD*

Indikator	<i>CRISP-DM</i>[38]	<i>Knowledge Discovery in Databases (KDD)</i>[38]
Fase	Enam tahapan yang dimulai dari pemahaman bisnis, pemahaman data, persiapan data, pemodelan, evaluasi, hingga implementasi.	Lima tahapan dimulai dari pemilihan, pra-pemrosesan, transformasi, penambahan data, hingga interpretasi/evaluasi.
Kekuatan	<ol style="list-style-type: none"> 1. Terdapat langkah akhir yang melibatkan tahap penerapan atau implementasi. 2. Proses yang terperinci dalam tahapannya. 3. Sesuai untuk proyek atau dataset yang besar. 4. Siklus berulang yang berkelanjutan. 5. Dilengkapi dengan dokumentasi dan contoh studi kasus yang serupa yang dapat mendukung jalannya proses. 	<ol style="list-style-type: none"> 1. Tidak terdapat proses akhir yang melibatkan tahap <i>deployment</i> (penerapan) namun di akhiri dengan tahap evaluasi. 2. Siklus pengulangan yang berkelanjutan. 3. Mendukung beragam model data mining, termasuk <i>Neural Network</i>. 4. Memerlukan pemahaman sebelumnya dalam <i>data mining</i>.
Limitasi	<ol style="list-style-type: none"> 1. Dimulai dengan memahami tujuan proyek melalui tahapan pemahaman bisnis. 2. Proses yang memakan waktu. 3. Tahap persiapan data dan pembuatan model berbeda dari <i>data mining</i> konvensional sehingga tidak termasuk dalam dokumen <i>CRISP-DM</i>. 	<ol style="list-style-type: none"> 1. Dimulai dengan tahap pemilihan data yang akan membentuk <i>dataset</i> yang akan digunakan dalam penelitian. 2. Memerlukan pemahaman sebelumnya dalam <i>data mining</i>.

Berdasarkan perbandingan dalam table diatas, penelitian ini memilih metode *CRISP-DM* karena prosesnya memiliki tahapan yang terstruktur dan jelas, yang dapat memberikan arahan yang kuat dalam jalannya penelitian. Pemilihan ini juga didasarkan pada penggunaan dataset yang

besar, di mana *CRISP-DM* memiliki panduan yang sesuai untuk melakukan proses data mining pada dataset yang besar, sehingga metode ini cocok untuk diadopsi dalam penelitian ini.



Gambar 3. 2 Flowchart Proses Pengolahan Data

3.3 Teknik Pengumpulan Data

Penelitian ini menggunakan *dataset* yang diambil dari Kaggle.com yang berasal dari Eswar Chand. Kaggle adalah *platform* daring yang menyediakan berbagai kompetisi data, sumber daya pembelajaran, dan kumpulan *dataset* untuk para ilmuwan data dan praktisi *machine learning*. *Dataset* yang digunakan merupakan data mengenai *phishing website detector*. Kaggle adalah *platform* daring yang menyediakan berbagai kompetisi data, sumber daya pembelajaran, dan kumpulan *dataset* untuk para ilmuwan data dan praktisi *machine learning*.

3.4 Variabel Penelitian

3.4.1 Variabel Independen

Variabel bebas, yang juga dikenal sebagai variabel independen, merupakan variabel yang diyakini menjadi penyebab atau faktor yang memiliki kemungkinan berdampak pada variabel lain dalam suatu penelitian atau eksperimen. Biasanya, variabel bebas ditunjukkan dengan simbol X. Secara umum, variabel bebas hadir terlebih dahulu dan dianggap mempengaruhi variabel lain yang terkait. Variable independen pada penelitian ini yaitu *Using IP, Long URL, Short URL, Symbol@, Redirecting, Prefix Suffix, Sub Domains, HTTPS, Domain Reg Len, Favicon, Non Std Port, HTTPS Domain URL, Request URL, Anchor URL, Links In Script Tags, Server From Handler, Info Email, Abnormal URL, Website Forwarding, Status Bar Cust, Disable Right Click, Using Popup Window, Iframe Redirection, Age of Domain, DNS Recording, Website Traffic, Page Rank, Google Index, Links Pointing to Page, Stats Report.*

3.4.2 Variabel Dependen

Variabel dependen adalah variabel yang terpengaruh atau variabel yang merupakan akibat dari variabel independen. Variabel dependen pada penelitian ini yaitu menderita mental disorder. Variabel tergantung atau variabel dependen adalah variabel yang dipengaruhi atau menjadi hasil dari variabel bebas. Dengan demikian, variabel ini bergantung pada variabel bebas dan besarnya dipengaruhi oleh perubahan dalam variabel independen. Koefisien perubahan dalam variabel independen memberikan gambaran seberapa besar perubahan yang mungkin terjadi dalam variabel dependen sebagai akibat dari perubahan dalam variabel independen. Variable dependen pada penelitian ini adalah class.

3.5 Teknik Analisis Data

3.5.1 *Business Understanding*

Pada tahap pertama, penelitian ini melakukan analisa terhadap kasus *phishing* di Indonesia serta meneliti data yang akan digunakan untuk

melihat kesesuaian data dengan penelitian. *Business Understanding* dalam konteks deteksi *URL web phishing* mencakup pemahaman yang mendalam tentang ancaman keamanan yang dihadapi oleh pengguna *web*. Hal ini mencakup pemahaman terhadap kebutuhan organisasi atau *platform* terkait keamanan *online*, serta dampak yang mungkin timbul akibat serangan *phishing*. Tujuannya adalah untuk mengidentifikasi aspek-aspek utama yang terkait dengan upaya pencegahan, deteksi, dan perlindungan terhadap serangan *phishing* pada *URL web*. Pada penelitian ini, tujuan utamanya adalah mendeteksi *URL website phishing* dengan mengadopsi metode klasifikasi *Random Forest*, *K-Nearest Neighbors*, dan *Support Vector Machine*.

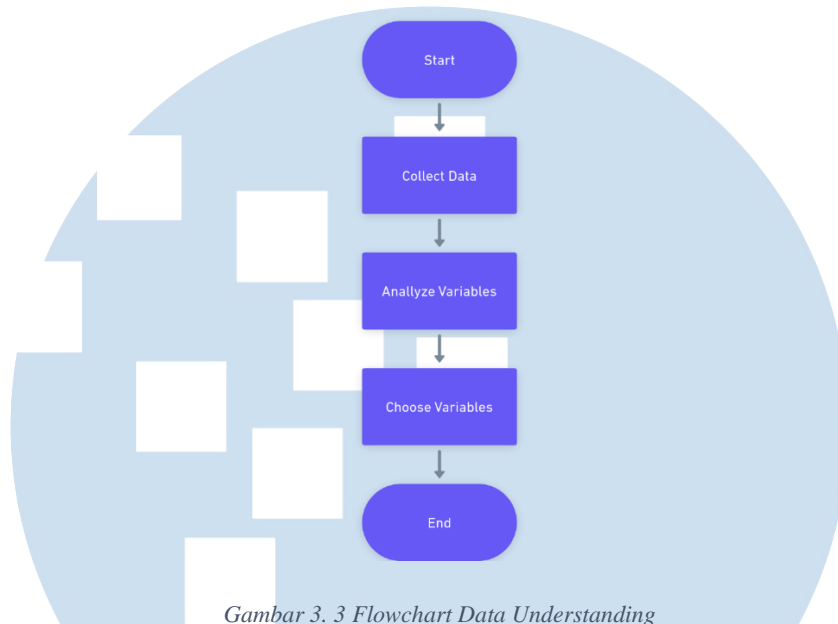
Tabel 3. 2 Penjelasan Atribut

No	Atribut	Keterangan
1	<i>Using IP</i>	Alamat IP sebagai domain.
2	<i>Long URL</i>	Panjangnya <i>url</i> .
3	<i>Short URL</i>	Pendeknya <i>url</i> .
4	<i>Symbol@</i>	Penggunaan symbol @.
5	<i>Redirecting</i>	Banyaknya pengalihan <i>website</i> yang dilakukan.
6	<i>Prefix Suffix</i>	Presentase penggunaan symbol “-“.
7	<i>Sub Domains</i>	Presentase penggunaan sub domain.
8	<i>HTTPS</i>	Alamat <i>https</i> .
9	<i>Domain Reg Len</i>	Batas berlakunya domain.
10	<i>Favicon</i>	Mempunyai <i>favicon</i> dari <i>link</i> eksternal.
11	<i>Non Std Port</i>	Tidak terdapat penggunaan <i>port</i> .
12	<i>HTTPS Domain URL</i>	Pemakaian <i>https</i> ke dalam bagaian domain pada <i>url</i> .

No	Atribut	Keterangan
13	<i>Request URL</i>	Presentase permintaan <i>url</i> eksternal dari keseluruhan (polinomial).
14	<i>Anchor URL</i>	Presentase penggunaan tag yang mengarah selain ke domain yang sama dari keseluruhan.

15	<i>Links In Script Tags</i>	Presentase pemakaian tag <i><link></i> , <i><script></i> , dan <i><meta></i> yang merujuk selain ke domain.
16	<i>Server From Handler</i>	Domain pemrosesan <i>server form handler</i> .
17	<i>Info Email</i>	Pemakaian fungsi “mail” dalam php.
18	<i>Abnormal URL</i>	Url yang tidak normal, kecocokan dengan web yang dituju.
19	<i>Website Forwarding</i>	Presentase <i>forwarding url</i> .
20	<i>Status Bar Cust</i>	Perubahan status bar ketika <i>mouseover</i> aktif.
21	<i>Disable Right Click</i>	Kondisi klik kanan pada <i>web</i> .
22	<i>Using Popup Window</i>	Presentase <i>popup window</i> dalam meminta <i>user</i> mengisi data.
23	<i>Iframe Redirection</i>	Penggunaan fungsi <i>iframe</i> .
24	<i>Age of Domain</i>	Umur domain.
25	<i>DNS Recording</i>	Terdapat catatan DNS pada domain.
26	<i>Website Traffic</i>	Presentase lalu lintas <i>web</i> dalam basis data Alexa.
27	<i>Page Rank</i>	Nilai <i>page rank website</i> .
28	<i>Google Index</i>	<i>Website</i> dalam indeks pencarian <i>google</i> .
29	<i>Links Pointing to Page</i>	Banyaknya link eksternal yang mengarah ke <i>website</i> .
30	<i>Stats Report</i>	Top <i>phishing</i> yang dibuat oleh beberapa pihak.
31	<i>Class</i>	Kelompok atribut.
32	<i>Id</i>	No atribut.

3.5.2 Data Understanding



Gambar 3. 3 Flowchart Data Understanding

Pada tahap kedua, dilakukan pengambilan data dari *platform* Kaggle. Proses ini dimulai dengan otomatisasi pada *platform* untuk mendapatkan semua tautan *listing* dan melakukan proses *scraping*. Kaggle dipilih sebagai sumber data karena *dataset* di Kaggle telah melalui proses kurasi dan pemrosesan awal. Ini berarti data-data tersebut telah diolah, divalidasi, dan disaring dari berbagai sumber untuk memastikan kualitas dan integritasnya. Proses ini menghasilkan data dengan 34 atribut yang diperoleh. Total data yang terdapat dalam dataset ini yaitu 11.000 lebih data. Atribut *dataset* tersebut terdiri dari 32 atribut sebagai berikut:

U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A

3.5.3 Data Preparation



Gambar 3. 4 Flowchart Data Preparation

Pada tahap ketiga, merupakan tahap *cleaning* dan *processing* data. Tahap ini dilakukan pengecekan *missing value* yang bertujuan untuk mengidentifikasi apakah terdapat nilai yang hilang atau tidak lengkap dalam *dataset* tersebut sehingga selanjutnya dapat dilakukan pengecekan deskripsi data, apakah data tersebut berupa numerik atau kategorikal, dan untuk menampilkan atribut apa saja yang terdapat dalam *dataset* tersebut. Dari hasil data *preparation* didapatkan 32 atribut dengan 11.000 data. Pada penelitian ini membagi data *preparation* menjadi dua bagian yaitu *data cleansing*, dan *split data*, 1) Saat melakukan pembersihan data, atribut yang tidak diperlukan akan dihapus, dan jika terdapat *missing value* pada salah satu atribut, baris tersebut akan dihapus secara keseluruhan, Dengan melakukan proses *cleansing*, kualitas data yang digunakan dalam analisis meningkat, menghasilkan informasi yang lebih akurat dan berkualitas. Proses ini memainkan peran penting dalam tahap *data mining* karena kualitas analisis bergantung pada kualitas data yang digunakan, 2) Proses pemisahan data akan membagi data menjadi dua bagian, yakni data *training* dan data *testing* dalam perbandingan 80:20, dengan 80% untuk *training* dan 20% untuk *testing*. Setelah menyelesaikan tahap persiapan data, langkah selanjutnya adalah pemodelan data.

3.5.4 Modeling

Pada tahap keempat ini dilakukan pemodelan, yaitu pemilihan model dan algoritma yang akan digunakan, selanjutnya akan dilakukan penerapan algoritma yang telah ditentukan. Pada penelitian ini menggunakan pemodelan data berupa klasifikasi dengan menampilkan nilai akurasi, presisi, *recall*, dan *f1 score*. Pada langkah ini, model data yang digunakan disesuaikan dengan kebutuhan untuk mencapai hasil yang diinginkan. Pada penelitian ini dipilih 3 algoritma dengan metode klasifikasi, yaitu terdapat algoritma *Random Forest*, *KNN*, dan *SVM*. Algoritma-algoritma ini dipilih berdasarkan evaluasi kelebihan, kekurangan, dan informasi dari penelitian sebelumnya terkait prediksi *URL website phishing* menggunakan *data mining*. Google colab dan *library* terkait akan dipakai untuk menyusun prediksi dalam pemodelan.

Tabel 3. 3 Kelebihan dan Kekurangan Algoritma

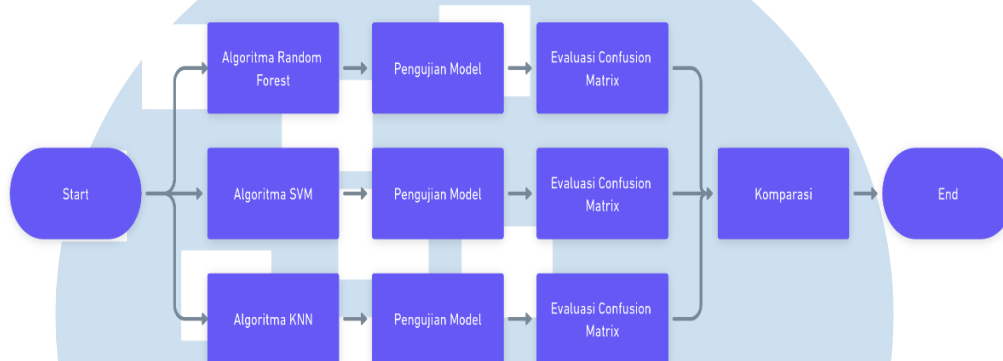
Algoritma	Kelebihan	Kekurangan
<i>Random Forest</i>	<ol style="list-style-type: none"> 1. Kekuatan dalam Klasifikasi dan Regresi: Mampu menangani baik masalah klasifikasi maupun regresi. 2. Pemrosesan Data yang Besar: Dapat bekerja dengan baik pada dataset yang besar dan memiliki kemampuan untuk menangani banyak fitur. 	<ol style="list-style-type: none"> 1. Kompleksitas Model: Dalam beberapa kasus, model <i>Random Forest</i> cenderung terlalu kompleks dan bisa sulit diinterpretasi. 2. Waktu Komputasi: Pada dataset yang sangat besar, waktu komputasi untuk membuat model bisa menjadi lambat.
<i>SVM</i>	<ol style="list-style-type: none"> 1. Efektif dalam Dimensi Tinggi: Cocok untuk <i>dataset</i> dengan dimensi tinggi dan bisa mengatasi <i>overfitting</i>. 2. Mampu Menangani <i>Non-linearitas</i>: Menggunakan <i>kernel trick</i> untuk menangani hubungan <i>non-linear</i> dalam data. 	<ol style="list-style-type: none"> 1. Sulit diatur: Memiliki beberapa parameter yang perlu diatur, yang dapat mempengaruhi performa model. 2. Kinerja Terbatas pada <i>Dataset</i> Besar: Waktu pelatihan <i>SVM</i> dapat memakan waktu pada dataset yang sangat besar.

Algoritma	Kelebihan	Kekurangan
<i>KNN</i>	<ol style="list-style-type: none"> 1. Sederhana dan Mudah Dimengerti: Konsep <i>KNN</i> relatif mudah dipahami. 2. Penerapan yang Fleksibel: Cocok untuk klasifikasi dan regresi, serta bisa digunakan dengan berbagai metrik jarak. 	<ol style="list-style-type: none"> 1. Pengaruh <i>Outlier</i>: Rentan terhadap pengaruh <i>outlier</i> karena keputusan berdasarkan sejumlah tetangga terdekat. 2. Kinerja yang Kurang di Dataset Besar: Pada dataset besar, waktu komputasi dan kinerja <i>KNN</i> bisa menjadi masalah.
<i>Naïve Bayes</i>	<ol style="list-style-type: none"> 1. Sederhana dan Cepat: Model yang sederhana dan efisien dalam pembelajaran. 2. Efektif pada Dataset Kecil: Bekerja dengan baik pada dataset kecil dan saat asumsi independensi berlaku. 	<ol style="list-style-type: none"> 1. Asumsi Independensi yang Sering Tidak Realistis: Asumsi bahwa semua fitur adalah independen tidak selalu berlaku pada data aktual. 2. Kinerja Terbatas pada Dataset yang Lebih Kompleks: Tidak begitu cocok untuk dataset yang sangat kompleks atau memiliki hubungan yang kompleks antara fitur.
<i>Decision Tree</i>	<ol style="list-style-type: none"> 1. Interpretasi yang Baik: Model <i>decision tree</i> mudah diinterpretasi oleh manusia. 2. Pemrosesan Data yang Efisien: Bisa mengatasi dataset dengan jumlah besar dan atribut yang beragam. 	<ol style="list-style-type: none"> 1. Kecenderungan <i>Overfitting</i>: <i>Decision tree</i> cenderung <i>overfit</i> pada data yang kompleks jika tidak diatur dengan baik. 2. Kehilangan Informasi: Pada saat pemisahan node, informasi bisa hilang jika tidak ada atribut lain yang dapat memecahkan data.

Pemilihan ketiga algoritma ini didapat dari studi-studi terdahulu di mana *Random Forest* telah sering diaplikasikan dalam memprediksi *URL website phishing*. Algoritma *Random Forest* dikenal memiliki kinerja

yang solid dan menunjukkan tingkat kesalahan yang rendah pada hasil prediksinya.

3.5.5 Evaluation



Gambar 3. 5 Flowchart Evaluation

Pada tahap kelima ini dilakukan implementasi menggunakan python. Penelitian ini juga memanfaatkan metrik A (Akurasi) P (*Precision*), R (*Recall*), dan F1 Score untuk menilai kinerja model klasifikasi yang dibangun. *Evaluation* adalah tahap penilaian terhadap model yang telah dibuat pada tahap pemodelan. Fokus utamanya adalah memastikan apakah hasil yang diperoleh sudah sejalan dengan tujuan bisnis yang telah ditetapkan sebelumnya. Studi ini memanfaatkan sejumlah metrik performa seperti waktu pemrosesan, akurasi, sensitivitas, dan presisi untuk mengevaluasi kinerja dari setiap model yang telah dibuat. Setelah itu, model mana yang menghasilkan prediksi terbaik akan ditentukan. Dari hasil evaluasi ini, data prediksi gangguan mental akan dianalisis bersama data demografi dalam tahap deployment yang akan datang. Evaluasi kinerja model klasifikasi berfokus pada pengukuran seberapa baik sistem mampu mengklasifikasikan data [18]. Salah satu cara untuk mengukur kinerja model klasifikasi adalah dengan menggunakan confusion matrix [18].

Dalam kerangka matriks kebingungan (*confusion matrix*), terdapat empat kemungkinan hasil klasifikasi, yakni TP (*True Positive*), FP (*False Positive*), TN (*True Negative*), dan FN (*False Negative*). Matriks kebingungan sering digunakan dalam menghitung tingkat ketepatan pada

konsep *data mining*. Berikut adalah tabel matriks kebingungan yang dimaksud:

Tabel 3. 4 Hasil True dan False

Prediksi	Hasil	
	1	0
1	TP (<i>True Positive</i>)	FN (<i>False Negative</i>)
0	FN (<i>False Negative</i>)	TN (<i>True Negative</i>)

Precision, dikenal juga sebagai nilai prediksi positif, merupakan persentase prediksi yang benar dan dihitung dengan $TP / (TP + FP)$. *Recall* adalah proporsi positif aktual dalam data yang diuji dan diwakili oleh $TP / (TP + FN)$, sedangkan *F-Measure* adalah nilai rata-rata dari *Precision* dan *Recall*, dihitung sebagai $2 (P \times R) / (P + R)$. Ketiga metrik ini memiliki rentang nilai antara 0 hingga 1. Dalam penelitian sebelumnya oleh Dongsong Zhang dan rekan-rekannya [2], mereka menduga bahwa model klasifikasi yang menggunakan fitur berbasis konten dan *URL* akan mengungguli model dengan fitur tradisional berdasarkan *Precision*, *Recall*, dan *F-Measure*. Sementara dalam penelitian ini, hipotesanya adalah bahwa fitur baru yang berfokus pada pendekatan konten dan *URL* dapat meningkatkan kinerja deteksi situs *phishing* jauh melampaui fitur dasar pada penelitian sebelumnya [2]

U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A