

BAB 3

METODE PENELITIAN

Metode yang digunakan pada penelitian ini yaitu klasifikasi dengan menggunakan algoritma C4.5 dan membuat pohon keputusan (*decision tree*) untuk menentukan tingkat akurasi dalam mengklasifikasikan serangan pada dataset. Data *log event*, *packet capture*, dan data lainnya yang menggambarkan serangan pada *intrusion detection system* dikumpulkan dari sumber sumber yang sesuai. Dataset NSL-KDD yang diperoleh dari *Kaggle.com*, dipilih untuk digunakan dalam penelitian ini. Data ini memiliki 41 atribut dan berisi informasi tentang serangan pada sistem IDS.

3.1 BAHAN PENELITIAN

Data yang digunakan dalam penelitian ini yaitu data sekunder NSL-KDD (*Network Security Layer – Knowledge Discovery in Database*). Data ini diambil dari *Kaggle.com* dan merupakan modifikasi dari dataset KDD Cup'99, yang berasal dari University of California, Irvine (UCI). Para peneliti di Institut Keamanan Siber Kanada (CIC) di Universitas New Brunswick (UNB) membuat kumpulan dataset NSL-KDD pada tahun 2009 dan telah mengalami beberapa pembaruan terbaru yang dirilis pada tahun 2019. Mereka mengambil data KDD Cup'99 dan mengatasi kekurangannya, termasuk redundansi kelas yang tidak seimbang, sehingga menghasilkan kumpulan data yang lebih baik untuk digunakan. Dataset ini digunakan untuk mengklasifikasikan serangan pada *intrusion detection system*. Ada beberapa alasan mengapa dataset NSL-KDD digunakan untuk algoritma pohon keputusan C4.5 dalam penelitian *intrusion detection system*:

1. Kesesuaian untuk algoritma C4.5

Pohon keputusan C4.5 baik digunakan untuk pengambilan keputusan yang selaras dengan karakteristik kumpulan dataset NSL-KDD. Format terstruktur yang berisi berbagai fitur, menyediakan pola yang sesuai dengan algoritma C4.5 untuk membangun aturan keputusan.

2. Format Standar

Kumpulan dataset NSL-KDD memiliki format yang terstruktur dengan baik dan terstandarisasi, sehingga mudah digunakan dengan berbagai alat dan algoritma pembelajaran mesin.

3. Data berlabel

Kumpulan dataset NSL-KDD diberi label, artinya setiap rekaman koneksi jaringan dikategorikan sebagai jenis normal atau jenis serangan. Data berlabel ini sangat penting untuk melatih model pembelajaran mesin untuk mengidentifikasi serangan.

Dataset NSL-KDD dan algoritma C4.5 membentuk kombinasi yang baik untuk penelitian *intrusion detection*. Kesesuaian kumpulan data untuk pohon keputusan, kemampuan interpretasi, dan kemampuan pemilihan fitur menjadikan pilihan yang sangat baik untuk mengevaluasi model IDS berbasis C4.5. Dengan pemrosesan awal data yang cermat dan penyesuaian parameter, C4.5 dapat secara efektif mengekstrak pola dari kumpulan dataset NSL-KDD dan membangun aturan keputusan yang kuat untuk mengidentifikasi serangan. Dataset NSL-KDD memiliki atribut yang menunjukkan perilaku dari setiap koneksi jaringan. Jenis layanan, protokol, jumlah byte yang dikirim, dan informasi lainnya dapat ditemukan di atribut ini. Data ini mencoba mensimulasikan skenario serangan jaringan komputer dengan menggabungkan berbagai macam serangan. Variasi serangan dan koneksi jaringan normal disimpan dalam data ini yang dapat digunakan sebagai referensi untuk deteksi serangan. Dalam dataset NSL-KDD setiap koneksi termasuk dalam salah satu dari dua kategori koneksi normal atau mencurigakan. Data tersebut mencakup jenis serangan yang dibagi menjadi lima kategori, yaitu:

1. Normal adalah lalu lintas jaringan yang normal dan tidak menunjukkan perilaku anormal.
2. *Denial of Services (DoS)* adalah serangan komputer dengan tujuan untuk menghentikan dan mengganggu layanan yang telah disediakan.
3. *Probe* adalah serangan yang dimaksudkan untuk mengumpulkan data dari jaringan atau sistem komputer seperti *scanning port* yang mencoba menemukan kelemahan sistem yang ada.

4. *User to Root (U2R)* adalah serangan dengan tujuan untuk meningkatkan akses user menjadi super user ketika pelaku mempunyai akun dalam sistem pada level user.
5. *Remote to Local (R2L)* adalah serangan jarak jauh ke lokasi dengan tujuan untuk mengakses sistem dalam jaringan menggunakan komputer lain.

Dataset NSL-KDD memiliki 41 fitur yang dapat digunakan untuk klasifikasi serangan pada IDS. Data ini bersumber dari *www.researchgate.net* yang ditunjukkan pada tabel 3.1

Tabel 3.1 41 Fitur Dataset NSL-KDD

No	Nama Atribut	Deskripsi
1	<i>Duration</i>	Panjang waktu
2	<i>Protocol type</i>	Penggunaan protokol koneksi seperti TCP, UDP, dan ICMP
3	<i>Service</i>	Layanan tujuan yang digunakan dalam jaringan
4	<i>Src_byte</i>	Jumlah <i>byte</i> data yang dikirim dari satu sumber ke tujuan dalam satu koneksi
5	<i>Dst_byte</i>	Jumlah <i>byte</i> data yang dikirim dari tujuan untuk sumber di satu koneksi
6	<i>Flag</i>	Status koneksi (normal atau <i>error</i>)
7	<i>Land</i>	Variabel ini bernilai 1 jika alamat IP sumber dan <i>port</i> tujuan sama. Jika tidak, variabel ini bernilai 0
8	<i>Wrong_fragment</i>	Jumlah nomor untuk fragment yang salah dalam koneksi
9	<i>Urgent</i>	Jumlah paket yang diperlukan untuk koneksi
10	<i>Hot</i>	Banyaknya indikator konten panas seperti: mengakses sistem direktori, membuat dan menjalankan program

No	Nama Atribut	Deskripsi
11	<i>Num_failed_logins</i>	Menghitung jumlah upaya login yang tidak berhasil
12	<i>Logged_in</i>	Login status ditandai dengan 1 berhasil dan 0 sebaliknya
13	<i>Num_compromised</i>	Nomor kondisi yang dapat dipromosikan
14	<i>Root_shell</i>	Jika nilainya 1, <i>shell root</i> dapat diakses, dan jika 0 sebaliknya
15	<i>Su_attempted</i>	Jika nilainya 1, percobaan perintah <i>su root</i> dapat digunakan, dan jika 0 sebaliknya
16	<i>Num_root</i>	Mengakses <i>root</i> yang signifikan atau jumlah kegiatan yang dilakukan sebagai <i>root</i> dalam koneksi
17	<i>Num_file_creations</i>	Jumlah operasi yang dilakukan untuk membuat file dalam satu koneksi
18	<i>Num_shells</i>	Banyaknya jumlah <i>shell prompt</i>
19	<i>Num_access_shells</i>	Nomor yang terlibat dalam operasi akses kontrol file
20	<i>Num_outbound_cmds</i>	Banyaknya perintah outbound yang ada dalam sesi ftp
21	<i>Is_hot_login</i>	Jika nilainya 1, login termasuk daftar <i>hot</i> , tetapi jika nilainya 0 maka tidak
22	<i>Is_guest_login</i>	Jika nilainya 1 maka login <i>guest</i> , tetapi jika nilainya 0 maka tidak
23	<i>Count</i>	Jumlah koneksi ke <i>host</i> tujuan yang sama dalam dua detik sebelumnya
24	<i>Serror_rate</i>	Persentase koneksi pada <i>flag</i> seperti <i>s0</i> , <i>s1</i> , <i>s2</i> atau <i>s3</i> yang dikumpulkan ke dalam <i>count</i>

No	Nama Atribut	Deskripsi
25	<i>Error_rate</i>	Persentase koneksi yang telah diaktifkan pada koneksi <i>flag</i> REJ yang dikumpulkan ke dalam <i>count</i>
26	<i>Same_srv_rate</i>	Persentase koneksi ke <i>service</i> yang sama, yang dikumpulkan pada koneksi <i>count</i>
27	<i>Diff_srv_rate</i>	Persentase koneksi ke berbagai <i>service</i> yang dikumpulkan pada koneksi <i>count</i>
28	<i>Srv_count</i>	Jumlah koneksi ke <i>service</i> yang sama dalam 2 detik terakhir
29	<i>Srv_serror_rate</i>	Persentase koneksi yang terhubung diantara koneksi <i>flag</i> s0, s1, s2, dan s3 yang dikumpulkan pada <i>srv_count</i>
30	<i>Srv_error_rate</i>	Jumlah koneksi yang telah mengaktifkan <i>flag</i> REJ, dihitung diantara koneksi <i>srv_count</i>
31	<i>Srv_diff_host_rate</i>	Persentase koneksi ke berbagai mesin tujuan dan dikumpulkan diantara koneksi <i>srv_count</i>
32	<i>Dst_host_count</i>	Nomor yang terhubung ke tujuan <i>host</i> IP <i>address</i>
33	<i>Dst_host_srv_count</i>	Nomor koneksi dengan <i>port</i> yang sama
34	<i>Dst_host_same_srv_count</i>	Jumlah koneksi ke layanan yang sama, dikumpulkan diantara koneksi <i>dst_host_count</i>
35	<i>Dst_host_diff_srv_count</i>	Persentase koneksi ke berbagai <i>services</i> , dikumpulkan diantara koneksi <i>dst_host_count</i>

No	Nama Atribut	Deskripsi
36	<i>Dst_host_same_src_port_rate</i>	Jumlah koneksi ke sumber <i>port</i> yang sama, yang dikumpulkan diantara koneksi <i>dst_host_srv_count</i>
37	<i>Dst_host_srv_diff_host_rate</i>	Jumlah koneksi ke mesin tujuan yang berbeda, dan dikumpulkan diantara koneksi <i>dst_host_srv_count</i>
38	<i>Dst_host_serror_rate</i>	Jumlah koneksi yang telah diaktifkan <i>flag</i> s0, s1, s2, dan s3 dikumpulkan diantara koneksi <i>dst_host_count</i>
39	<i>Dst_host_srv_serror_rate</i>	Jumlah koneksi yang telah diaktifkan <i>flag</i> s0, s1, s2, dan s3 dikumpulkan diantara koneksi <i>dst_host_srv_count</i>
40	<i>Dst_host_rerror_rate</i>	Jumlah koneksi yang telah diaktifkan <i>flag</i> REJ, yang dikumpulkan diantara koneksi <i>dst_host_count</i>
41	<i>Dst_host_srv_rerror_rate</i>	Jumlah koneksi yang telah diaktifkan diantara koneksi <i>dst_host_srv_count</i>

3.2 ALAT PENELITIAN

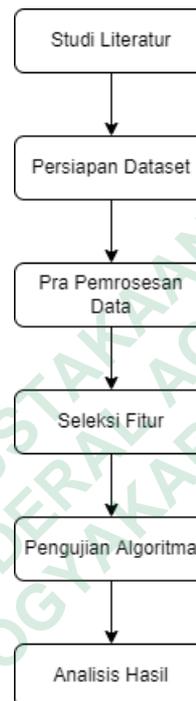
Untuk mendukung pelaksanaan penelitian, diperlukan alat-alat berupa piranti-piranti yang berguna untuk kebutuhan analisis dan keberhasilan penelitian untuk memproses atau memperlakukan bahan penelitian. Komputer dengan spesifikasi cukup digunakan dalam penelitian ini untuk menjalankan sistem operasi serta software yang dapat terhubung ke internet. Berikut adalah sistem operasi dan aplikasi yang penulis gunakan meliputi:

1. Sistem Operasi: Windows 10
2. Processor: Intel Celeron tipe N4000 CPU@2.8 Ghz
3. RAM: 4 GB
4. Storage: SSD 512 GB

5. Software: Weka, Google Colab

3.3 JALAN PENELITIAN

Jalan penelitian mencakup penjelasan lengkap dan mendalam tentang proses yang diambil selama pelaksanaan penelitian dan pengumpulan data. Jalan penelitian diperlihatkan pada gambar 3.1



Gambar 3.1 Jalan Penelitian

Berikut adalah penjelasan tentang alur jalannya penelitian:

1. Studi Literatur

Dalam tahap ini dilakukan literatur *review* untuk mendapatkan informasi dan acuan dari buku atau jurnal penelitian yang relevan mengenai implementasi *data mining* untuk klasifikasi serangan pada sistem deteksi intrusi.

2. Persiapan dataset

Pengumpulan data *intrusion detection system* yang akan digunakan untuk penelitian. Data ini berupa log kegiatan yang terjadi dalam jangka waktu tertentu pada sistem jaringan. Dataset yang digunakan adalah NSL-KDD yang berupa file CSV yang telah berlabel.

3. Pra Pemrosesan Data

Setelah data dikumpulkan, proses *preprocessing* harus dilakukan untuk menghilangkan data yang tidak diperlukan atau tidak valid, mengubah format data. Setelah itu, data difilter dengan mengubah atribut yang akan digunakan.

4. Seleksi Fitur

Setelah *preprocessing*, dilakukan pemilihan fitur yang relevan untuk klasifikasi serangan. Fitur ini biasanya disebut atribut dan merupakan input yang digunakan untuk algoritma *data mining*. Dalam penelitian sebelumnya, atribut seperti *protocol type*, *service*, *flag*, *src bytes*, dan *dst bytes* adalah atribut yang berkaitan dengan klasifikasi serangan pada *intrusion detection system* yang sering digunakan.

5. Pengujian Algoritma C4.5

Algoritma ini merupakan pembagian kelas yang menggunakan pohon keputusan untuk mengklasifikasikan data ke dalam kelas tertentu. Setelah data diproses, maka data akan dibuat menjadi input. Algoritma C4.5 akan membuat pohon keputusan untuk mengklasifikasikan serangan pada *intrusion detection system*. Setelah pohon keputusan dibuat, maka perlu dilakukan pengujian akurasi. Akurasi algoritma C4.5 dapat dihitung dengan menggunakan persamaan *precision*, *recall*, dan *accuracy*.

6. Analisis Hasil

Setelah pengujian selesai, maka perlu dilakukan penilaian dan penyelesaian tentang hasil pengujian. Penilaian terdiri atas bagaimana kinerja dan tingkat akurasi klasifikasi dengan algoritma C4.5 pada *intrusion detection system*.