



Artikel Penelitian

## Naïve Bayes dan Confusion Matrix untuk Efisiensi Analisa Intrusion Detection System Alert

Muhammad Kamil Suryadewiansyah<sup>a</sup>, Teja Endra Eng Tju<sup>b,\*</sup>

<sup>a, b</sup> Universitas Budi Luhur, Jl. Ciledug Raya, Petukangan Utara, Jakarta Selatan, 12260, DKI Jakarta, Indonesia

### INFORMASI ARTIKEL

#### Sejarah Artikel:

Diterima Redaksi: 14 Juli 2022

Revisi Akhir: 29 Agustus 2022

Diterbitkan Online: 31 Agustus 2022

### KATA KUNCI

Rule-based,  
Likelihood,  
Security,  
Probability,  
F-1 score

### KORESPONDENSI

E-mail: [teja.endraengtju@budiluhur.ac.id](mailto:teja.endraengtju@budiluhur.ac.id)

### A B S T R A C T

Banyaknya *malware* menyebabkan IDS (*Intrusion Detection System*) dituntut menyesuaikan diri semakin kompleks sehingga mahal dan membebani perusahaan yang menggunakannya. Sistem yang berbasis teknologi *Host-based* IDS dan *Signature-based* IDS sudah banyak digunakan namun hanya mampu mendeteksi serangan yang sudah diketahui sebelumnya, untuk memperbaiki kinerjanya perlu dilakukan analisa pada data *log* berdasarkan *alert* yang diberikan. Teknik klasifikasi *Naïve Bayes* digunakan untuk membantu meningkatkan efisiensi dan efektifitas analisa tersebut. Penelitian ini dilakukan dengan mengambil empat langkah bagian dari metodologi SKKNI (Standar Kompetensi Kerja Nasional Indonesia) No.299 tahun 2020, *Artificial Intelligence*, sub bidang *Data Science*, yaitu *data understanding*, *data preparation*, *modeling*, dan *model evaluation*. *Dataset* dari penyedia layanan IDS sebanyak 575 data yang dibagi menjadi 515 data latihan dan 60 data uji. Hasil evaluasi data uji dengan *confusion matrix* diperoleh pengukuran metrik *accuracy* 0,87, *recall* 0,89, *precision* 0,83, dan *F-Measure* 0,86. Adanya FP (*False Positive*) dan FN (*False Negatif*), keduanya sangat penting bagi pengguna IDS untuk meningkatkan kualitas layanan kepada pelanggan dan mengurangi resiko akibat adanya intrusi. FP dan FN menjadi fokus dalam melakukan analisa *log alert* dari IDS sehingga tidak perlu menganalisa keseluruhan data, berdampak memberikan hasil 85% lebih efektif dan berkontribusi pada efisiensi tenaga dan waktu bagi tim keamanan suatu perusahaan pengguna IDS. Selain itu didapat bahwa sekitar 50% data IDS adalah intrusi atau gangguan lainnya.

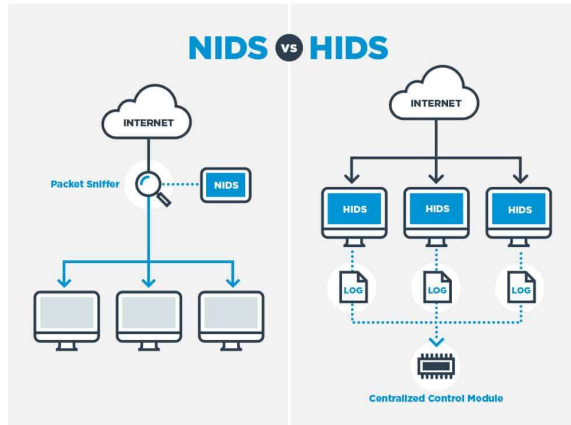
## 1. PENDAHULUAN

IDS (*Intrusion Detection System*) terus berkembang sejalan semakin kompleks, rumit, dan banyaknya *malware*. Teknologi IDS pada umumnya terdiri dari dua yaitu HIDS (*Host-based* atau *Host-level* IDS), yang lebih sederhana serta ekonomis, dan NIDS (*Network-based* atau *Network-level* IDS) yang memerlukan perangkat keras di jaringan [1,2]. HIDS mendeteksi pemakaian tidak sah, abnormal, aktivitas jahat pada *host* (suatu *server* komputer), sedangkan NIDS mendeteksi penyerangan atau intrusi pada jaringan. HIDS merupakan aplikasi yang dipasang pada sistem autonomous (komputer) untuk mendeteksi penyalahgunaan yang biasanya tercatat pada *log* sistem operasi meliputi *system files*, proses, utilisasi sumber daya, hak akses.

NIDS mencatat paket dari trafik jaringan, diekstraksi dengan *protocol analysis tools* pada *header* dari paket tersebut berdasarkan beberapa parameter untuk mendeteksi aktivitas jahat dan bisa dipasang pada *backbone* jaringan, *server*, *switches* *routers*, dan *gateways* [2,3]. Di dalam konfigurasi jaringan komputer, perbedaan antara HIDS dan NIDS ditunjukkan pada Gambar 1. Tampak pada NIDS terdapat perangkat keras *Packet Sniffer* yang disisipkan di tengah jaringan komputer, sedangkan pada HIDS *log* yang berisi *alert* dihasilkan dari aplikasi atau perangkat lunak yang terpasang pada *host* atau *server*.

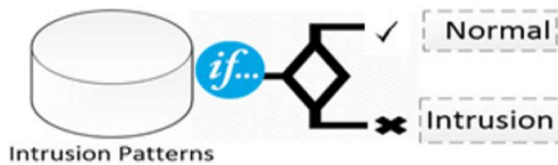
Metode deteksi IDS secara luas dibedakan menjadi dua kategori yaitu SIDS (*Signature-based* IDS), berdasarkan pola yang sudah ada sebelumnya serta diterapkan secara *rule-based*, dan AIDS (*Anomaly-based* IDS) dengan menggunakan statistik atau

machine learning [1]. SIDS telah terbukti efektif mendeteksi serangan tanpa banyak kesalahan, deteksi dilakukan berdasarkan pola dari serangan yang sudah diketahui dan disimpan di dalam basis data informasi, seperti pada Gambar 2, sehingga tidak mampu mendeteksi jenis serangan yang belum diketahui sebelumnya. AIDS mampu mendeteksi serangan yang sudah ataupun belum diketahui berdasarkan profil atau model statistik dari data histori, dengan mencari anomali atau kejadian yang tidak wajar dari kebiasaan [4,5].



Gambar 1. Konfigurasi Jaringan NIDS dan HIDS [6] (Sumber: Cooper, 2022)

Defenxor sebagai penyedia produk keamanan teknologi informasi yang salah satunya adalah IDS berbasis HIDS dan SIDS menghadapi tantangan untuk memberikan inovasi dengan biaya rendah. Adanya keluhan dari pengguna IDS bahwa log sistem memberikan alert yang sangat banyak dan diperiksa secara manual perlu diatasi. Untuk menghadapi tantangan dan keluhan tersebut, penelitian ini dilakukan sehingga proses analisa bisa dilakukan dengan lebih efektif dan efisien menggunakan machine learning dengan teknik klasifikasi Naive Bayes.



Gambar 2. Konsep Kerja SIDS [1] (Sumber: Khraisat et al, 2019)

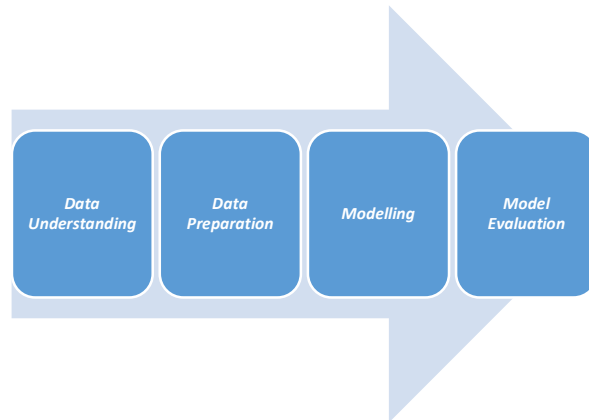
Peneliti sebelumnya terkait alert IDS dengan klasifikasi Naive Bayes antara lain dilakukan dengan: KDD99 dataset dan teknik PCA [7], AIDS dan teknik Correlation-Based Feature Selection [8], NSL-KDD dataset [9], zero probability [10], NSL-KDD dataset dan teknik PCA ditambah SVM [11], KNN sebagai perbandingan [12], diskritisasi variabel [13], Gaussian Naive Bayes dan Kyoto dataset [14], HND dan KDD99 dataset [15], MNBIDS dan KDD99 dataset [16].

Berdasarkan penjelasan di atas, penelitian ini perlu dilakukan sebagai penerapan metode Naive Bayes dasar untuk keperluan efisiensi dan efektifitas kerja tim keamanan jaringan suatu perusahaan dalam menganalisa data alert IDS. Kontribusi kebaruan berupa penerapan metode dengan data alert dari log

pada perusahaan pengelola IDS. Dengan mempertahankan mesin IDS berbasis HIDS dan SIDS yang ekonomis karena merupakan teknik yang paling sederhana, dibantu dengan teknik klasifikasi Naive Bayes dan evaluasi Confusion Matrix diharapkan memberikan efektifitas dan efisiensi dalam analisa alert dari IDS tersebut.

## 2. METODE

Penelitian menggunakan dataset dari log alert IDS sebanyak 575 data dengan metodologi menggunakan SKKNI (Standar Kompetensi Kerja Nasional Indonesia) No. 299 tahun 2020, bidang keahlian Artificial Intelligence, sub bidang Data Science [17]. Dari tujuh fungsi utama (business understanding, data understanding, data preparation, modeling, model evaluation, deployment, dan evaluation) digunakan empat saja yang relevan dengan kegiatan penelitian ini yaitu data understanding, data preparation, modeling, dan model evaluation [18], seperti ditunjukkan pada Gambar 3.

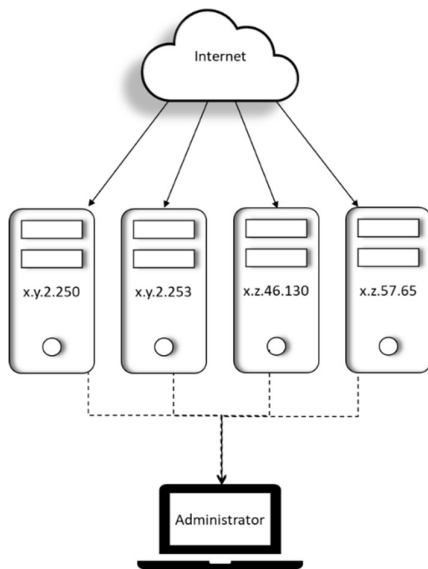


Gambar 3. Metodologi Penelitian

### 2.1. Data Understanding

Data diambil dari lokasi pengguna IDS dengan konfigurasi jaringan pada Gambar 4. Terdiri dari empat server yang terpasang HIDS, log dari setiap server diambil oleh administrator yang juga sebagai centralized control module SIDS.

Setiap data terdiri dari lima atribut yaitu Alert, Severity Level, IP Category, IP Destination, dan Event. Setiap atribut memiliki variabel nominal atau kategorikal. Alert berisi Yes berarti terjadi intrusi, dan No berarti aman. Severity Level dikategorikan Critical dan Medium saja. IP Category dibedakan Internal, External, Unknown yang merupakan alamat IP pengakses berasal. IP Destination terdiri dari empat alamat IP server yang dituju dengan beberapa angka disamarkan dengan huruf x, y, z karena alasan kerahasiaan dan keamanan. Event memiliki enam jenis kejadian yang terdeteksi oleh IDS yaitu 4264 (logon success), 4625 (logon failed), 4634 (logout), 4662 (Active Directory Enumeration), 4782 (pass the hash), dan XSS (cross site scripting) [19]. Pada Tabel 1 ditampilkan dua puluh contoh data dari dataset yang digunakan.



Gambar 4. Konfigurasi Jaringan Sumber Data

Tabel 1. Tampilan Dua Puluh Contoh Data dari Dataset.

No	Event	IP Destination	IP Category	Severity Level	Alert
1	4662	x.y.2.250	Internal	Critical	No
2	4662	x.z.46.130	Internal	Medium	Yes
3	XSS	x.y.2.250	Unknown	Critical	Yes
4	4662	x.z.57.65	Internal	Medium	No
5	4662	x.z.57.65	Internal	Medium	No
6	4662	x.z.57.65	Eksternal	Medium	No
7	XSS	x.z.46.130	Unknown	Medium	No
8	4624	x.y.2.253	Internal	Critical	No
9	XSS	x.y.2.253	Unknown	Critical	Yes
10	4634	x.y.2.253	Internal	Critical	No
11	4782	x.z.2.253	Internal	Critical	No
12	4624	x.z.2.253	Unknown	Critical	No
13	XSS	x.z.46.130	Internal	Medium	Yes
14	4782	x.z.2.250	Internal	Critical	No
15	XSS	x.z.46.130	Eksternal	Medium	Yes
16	4634	x.y.2.250	Internal	Critical	No
17	4625	x.y.57.65	Internal	Medium	Yes
18	4625	x.z.46.130	Internal	Medium	No
19	4625	x.z.2.253	Unknown	Critical	Yes
20	4624	x.z.57.65	Internal	Medium	No
...	...	...	...	...	...

### 2.2. Data Preparation

Dataset dengan 575 data secara acak (*random*) dibagi menjadi data latih (*training*) sebanyak 515 dan data uji (*test*) sebanyak 60 [20]. Atribut *Alert* dalam analisa dipakai sebagai acuan ketepatan IDS, oleh karena itu dalam penelitian digunakan sebagai supervisi dengan kelas *Yes* dan *No*. Selanjutnya dari *dataset* latih dihitung semua kejadian yang muncul pada tiap-tiap variabel dan dikelompokkan berdasarkan kelas *Yes* dan *No*. Hasil pengelompokan dan perhitungan jumlah muncul pada data latih ditunjukkan pada Tabel 2.

Tabel 2. Jumlah Kemunculan Setiap Variabel pada Data Latih.

Atribut	Variabel	Jumlah Muncul	
		Yes	No
<i>Alert</i> (supervisi)	<i>Yes (Y)</i>	264	-
	<i>No (N)</i>	-	251
<i>Severity Level (Sx)</i>	<i>Critical (SC)</i>	147	118
	<i>Medium (SM)</i>	117	133
	<i>Internal (CI)</i>	57	208
<i>IP Category (Cx)</i>	<i>External (CE)</i>	89	25
	<i>Unknown (CU)</i>	118	18
<i>IP Destination (Dx)</i>	x.y.2.250 ( <i>D1</i> )	65	55
	x.y.2.253 ( <i>D2</i> )	82	63
	x.z.46.130 ( <i>D3</i> )	58	60
	x.z.57.65 ( <i>D4</i> )	59	73
<i>Event (Ex)</i>	4624 ( <i>E1</i> )	30	49
	4625 ( <i>E2</i> )	33	15
	4634 ( <i>E3</i> )	12	42
	4662 ( <i>E4</i> )	50	115
	4782 ( <i>E5</i> )	17	11
	XSS ( <i>E6</i> )	122	19

### 2.3. Modeling

Berdasarkan tipe data yang bersifat nominal ataupun kategorikal serta setiap variabel bisa dianggap independen maka dipilih teknik klasifikasi dengan model *Naïve Bayes* [21]. Selain itu jumlah data tidak terlalu banyak sehingga lebih sederhana, mudah, dan cepat diimplementasikan [22].

Metode *Naïve Bayes* berdasarkan pada teorema yang dibuat oleh Thomas Bayes pada pertengahan abad 18 [23]. Teorema ini menggunakan hipotesa dari data yang diberikan termasuk ke dalam kelas tertentu dengan menghitung probabilitas data tersebut mengacu pada hasil perhitungan probabilitas data histori pada setiap kelas. Jika diterapkan dengan atribut pada Tabel 2, maka diperoleh persamaan (1) dan (2).

$$P(Y|Sx, Cx, Dx, Ex) = \frac{P(Sx, Cx, Dx, Ex|Y) P(Y)}{P(Sx, Cx, Dx, Ex)} \tag{1}$$

$$P(N|Sx, Cx, Dx, Ex) = \frac{P(Sx, Cx, Dx, Ex|N) P(N)}{P(Sx, Cx, Dx, Ex)} \tag{2}$$

Simbol *x* di dalam persamaan (1) dan (2) menunjukkan variabel yang muncul pada setiap kejadian (pada data dengan atribut terkait). Tampak bahwa *denominator* (penyebut/pembagi) pada kedua persamaan tersebut sama dan karena tujuannya adalah membandingkan hasil perhitungan untuk mendapatkan prediksi *Yes* atau *No* (berdasarkan nilai yang lebih besar) maka cukup dihitung *numerator* (pembilang) sebagai *likelihood* (hipotesa kemungkinan) [23], ditunjukkan pada persamaan (3) dan (4).

$$P(Sx, Cx, Dx, Ex|Y) P(Y) = P(Sx|Y) P(Cx|Y) P(Dx|Y) P(Ex|Y) P(Y) \tag{3}$$

$$P(Sx, Cx, Dx, Ex|N) P(N) = P(Sx|N) P(Cx|N) P(Dx|N) P(Ex|N) P(N) \tag{4}$$

**2.4. Model Evaluation**

Hasil prediksi dengan *Naïve Bayes* yaitu *Yes* atau *No* dibandingkan dengan isian *Alert* yang sudah ada di setiap data. Berdasarkan penjelasan di *data understanding*, didefinisikan jika *Yes* adalah kelas positif dan *No* adalah kelas negatif. Dengan demikian, jika hasil prediksi dan *Alert*, sama-sama *Yes* maka disebut *True Positive* (TP), apabila sama-sama *No* maka disebut *True Negative* (TN). Jika hasil prediksi *Yes* dan *Alert* berisi *No* maka disebut *False Positive* (FP), apabila sebaliknya disebut *False Negative* (FN). Selanjutnya, hasil perbandingan dihitung frekuensinya dan disajikan dalam bentuk *confusion matrix* [23–25] seperti pada Tabel 3.

Metrik yang digunakan adalah *accuracy*, *precision*, *recall*, dan *F-Measure*. *Accuracy* mengukur tingkat ketepatan prediksi yang bernilai sama dengan *Alert* yaitu TP dan TN terhadap hasil keseluruhan. *Precision* mengidentifikasi frekuensi prediksi sesuai dengan *Alert* yaitu sama-sama positif atau TP terhadap seluruh prediksi positif. *Recall* mengidentifikasi frekuensi prediksi sesuai dengan *Alert* yaitu sama-sama positif atau TP terhadap seluruh *Alert* positif. *F-Measure* dikenal juga sebagai *F-score* atau *F-1 score* merupakan rerata harmonis dari *precision* dan *recall* dengan nilai terbaik adalah 1 (*precision* dan *recall* sempurna) dan nilai terburuk adalah 0 [23,26].

**3. HASIL**

Perhitungan probabilitas setiap variabel pada atributnya dengan syarat kelas *Yes* dan *No* berdasarkan jumlah muncul atau frekuensi dari Tabel 2 sedemikian sehingga hasilnya ditampilkan pada Tabel 4.

Tabel 3. Confusion Matrix.

		Prediksi		
		Yes	No	
Alert	Yes	TP	FN	$Recall = \frac{TP}{(TP+FN)}$
	No	FP	TN	
		$Precision = \frac{TP}{(TP+FP)}$		$Accuracy = \frac{(TP+TN)}{(TP+TN+FP+FN)}$
		$F - Measure = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$		

Tabel 4. Parameter Model Hasil Perhitungan Probabilitas.

Atribut	Variabel	Probabilitas	
		Yes	No
A	Y	$P(Y) = 0,512621$	-
	N	-	$P(N) = 0,487379$
Sx	SC	$P(SC Y) = 0,556818$	$P(SC N) = 0,470120$
	SM	$P(SM Y) = 0,443182$	$P(SM N) = 0,529880$
Cx	CI	$P(CI Y) = 0,215909$	$P(CI N) = 0,828685$
	CE	$P(CE Y) = 0,337121$	$P(CE N) = 0,099602$
Dx	CU	$P(CU Y) = 0,446969$	$P(CU N) = 0,071713$
	D1	$P(D1 Y) = 0,246212$	$P(D1 N) = 0,219124$
Dx	D2	$P(D2 Y) = 0,310606$	$P(D2 N) = 0,250996$
	D3	$P(D3 Y) = 0,219697$	$P(D3 N) = 0,239044$
Dx	D4	$P(D4 Y) = 0,223485$	$P(D4 N) = 0,290837$
	E1	$P(E1 Y) = 0,113636$	$P(E1 N) = 0,195219$
Ex	E2	$P(E2 Y) = 0,125000$	$P(E2 N) = 0,059761$
	E3	$P(E3 Y) = 0,045455$	$P(E3 N) = 0,167331$
Ex	E4	$P(E4 Y) = 0,189394$	$P(E4 N) = 0,458167$
	E5	$P(E5 Y) = 0,064394$	$P(E5 N) = 0,043825$
Ex	E6	$P(E6 Y) = 0,462121$	$P(E6 N) = 0,075697$

Parameter pada Tabel 4 digunakan sebagai referensi untuk menghitung *likelihood Yes* dengan persamaan (3) dan *No* dengan persamaan (4) pada setiap data untuk mendapatkan hasil prediksi [27]. Pada Tabel 5 diberikan contoh hasil perhitungan *likelihood* dengan menggunakan dua puluh data dari Tabel 1.

Tabel 5. Contoh Pehitungan *Likelihood* dan Hasil Prediksi.

No	Ex	Dx	Cx	Sx	Alert	Prediksi	CM
1	E4	D1	CI	SC	No	No	TN
Y	$P(E4 Y)$	$P(D1 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,002874	
N	$P(E4 N)$	$P(D1 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,019063</b>	✓
2	E4	D3	CI	SM	Yes	No	FN
Y	$P(E4 Y)$	$P(D3 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,002041	
N	$P(E4 N)$	$P(D3 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,023439</b>	✓
3	E6	D1	CU	SC	Yes	Yes	TP
Y	$P(E6 Y)$	$P(D1 Y)$	$P(CU Y)$	$P(SC Y)$	$P(Y)$	<b>0,014516</b>	✓
N	$P(E6 N)$	$P(D1 N)$	$P(CU N)$	$P(SC N)$	$P(N)$	0,000273	
4	E4	D4	CI	SM	No	No	TN
Y	$P(E4 Y)$	$P(D4 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,002076	
N	$P(E4 N)$	$P(D4 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,028517</b>	✓
5	E4	D4	CI	SM	No	No	TN
Y	$P(E4 Y)$	$P(D4 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,002076	
N	$P(E4 N)$	$P(D4 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,028517</b>	✓
6	E4	D4	CE	SM	No	No	TN
Y	$P(E4 Y)$	$P(D4 Y)$	$P(CE Y)$	$P(SM Y)$	$P(Y)$	0,003242	
N	$P(E4 N)$	$P(D4 N)$	$P(CE N)$	$P(SM N)$	$P(N)$	<b>0,003428</b>	✓
7	E6	D3	CU	SM	No	Yes	FP
Y	$P(E6 Y)$	$P(D3 Y)$	$P(CU Y)$	$P(SM Y)$	$P(Y)$	<b>0,010309</b>	✓
N	$P(E6 N)$	$P(D3 N)$	$P(CU N)$	$P(SM N)$	$P(N)$	0,000335	
8	E1	D2	CI	SC	No	No	TN
Y	$P(E1 Y)$	$P(D2 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,002175	
N	$P(E1 N)$	$P(D2 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,009304</b>	✓
9	E6	D2	CU	SC	Yes	Yes	TP
Y	$P(E6 Y)$	$P(D2 Y)$	$P(CU Y)$	$P(SC Y)$	$P(Y)$	<b>0,018313</b>	✓
N	$P(E6 N)$	$P(D2 N)$	$P(CU N)$	$P(SC N)$	$P(N)$	0,000312	
10	E3	D2	CI	SC	No	No	TN
Y	$P(E3 Y)$	$P(D2 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,000870	
N	$P(E3 N)$	$P(D2 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,007975</b>	✓

11	E5	D2	CI	SC	No	No	TN
Y	$P(E5 Y)$	$P(D2 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,001233	
N	$P(E5 N)$	$P(D2 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,002089</b>	✓
12	E1	D2	CU	SC	No	Yes	FP
Y	$P(E1 Y)$	$P(D2 Y)$	$P(CU Y)$	$P(SC Y)$	$P(Y)$	<b>0,004503</b>	✓
N	$P(E1 N)$	$P(D2 N)$	$P(CU N)$	$P(SC N)$	$P(N)$	0,000805	
13	E6	D3	CI	SM	Yes	Yes	TP
Y	$P(E6 Y)$	$P(D3 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	<b>0,004980</b>	✓
N	$P(E6 N)$	$P(D3 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	0,003873	
14	E5	D1	CI	SC	No	No	TN
Y	$P(E5 Y)$	$P(D1 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,000977	
N	$P(E5 N)$	$P(D1 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,001823</b>	✓
15	E6	D3	CE	SM	Yes	Yes	TP
Y	$P(E6 Y)$	$P(D3 Y)$	$P(CE Y)$	$P(SM Y)$	$P(Y)$	<b>0,007776</b>	✓
N	$P(E6 N)$	$P(D3 N)$	$P(CE N)$	$P(SM N)$	$P(N)$	0,000466	
16	E3	D1	CI	SC	No	No	TN
Y	$P(E3 Y)$	$P(D1 Y)$	$P(CI Y)$	$P(SC Y)$	$P(Y)$	0,000690	
N	$P(E3 N)$	$P(D1 N)$	$P(CI N)$	$P(SC N)$	$P(N)$	<b>0,006962</b>	✓
17	E2	D4	CI	SM	Yes	No	FN
Y	$P(E2 Y)$	$P(D4 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,001370	
N	$P(E2 N)$	$P(D4 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,003720</b>	✓
18	E2	D3	CI	SM	No	No	TN
Y	$P(E2 Y)$	$P(D3 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,001347	
N	$P(E2 N)$	$P(D3 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,003057</b>	✓
19	E2	D2	CU	SC	Yes	Yes	TP
Y	$P(E2 Y)$	$P(D2 Y)$	$P(CU Y)$	$P(SC Y)$	$P(Y)$	<b>0,004953</b>	✓
N	$P(E2 N)$	$P(D2 N)$	$P(CU N)$	$P(SC N)$	$P(N)$	0,000247	
20	E1	D4	CI	SM	No	No	TN
Y	$P(E1 Y)$	$P(D4 Y)$	$P(CI Y)$	$P(SM Y)$	$P(Y)$	0,001246	
N	$P(E1 N)$	$P(D4 N)$	$P(CI N)$	$P(SM N)$	$P(N)$	<b>0,012151</b>	✓
...	...	...	...	...	...	...	...

Hasil perhitungan keseluruhan pada 515 data latih dan pada 60 data uji dengan menghitung banyaknya TP, TN, FP, dan FN disajikan pada Tabel 6 berupa *confusion matrix* dan metrik pengukuran *accuracy*, *precision*, *recall*, dan *F-Measure*. Tampak antara data latih dan data uji memberikan angka relatif sama.

Tabel 6. Hasil Pengukuran Data Latih dan Data Uji.

		Prediksi		
		Yes	No	
<b>Data Latih</b>				
<b>Alert</b>	<b>Yes</b>	TP = 239	FN = 25	<b>Recall = 0,91</b>
	<b>No</b>	FP = 52	TN = 199	
		<b>Precision = 0,82</b>		<b>Accuracy = 0,85</b>
		<b>F – Measure = 0,86</b>		
<b>Data Uji</b>				
<b>Alert</b>	<b>Yes</b>	TP = 25	FN = 3	<b>Recall = 0,89</b>
	<b>No</b>	FP = 5	TN = 27	
		<b>Precision = 0,83</b>		<b>Accuracy = 0,87</b>
		<b>F – Measure = 0,86</b>		

#### 4. PEMBAHASAN

Dari Tabel 6, hasil pengolahan data tampak jumlah intrusi (dari total TP data latih dan data uji) sebanyak 264, kurang lebih setengah dari data yang ada kemungkinan besar adalah intrusi. Hal ini menunjukkan bahwa IDS sangat penting untuk perusahaan dan cukup banyak lalu lintas di Internet adalah gangguan atau tidak ada manfaatnya.

Adanya FN dan FP, sesuai dengan karakteristik IDS dan tujuan bisnis bahwa keduanya cukup penting dan kritikal. Jika FN besar yang berarti hasil prediksi tidak terjadi intrusi sedangkan pada *Alert* terjadi intrusi maka transaksi akan diblokir dan menyebabkan banyak pelanggan kecewa karena mungkin saja transaksi tersebut bukan suatu intrusi. Jika FP besar yang berarti hasil prediksi terjadi intrusi sedangkan pada *Alert* tidak terjadi intrusi maka transaksi akan diteruskan dan menyebabkan banyak kerugian karena mungkin saja transaksi tersebut adalah intrusi.

Hasil metrik *accuracy*, *precision*, *recall*, dan *F-Measure* dari data latih dan data uji tidak jauh berbeda dengan skor cukup bagus. Untuk meningkatkan skor tersebut bisa dilakukan dengan menurunkan jumlah FN dan FP. Hal ini sejalan dengan tujuan penelitian yaitu untuk meningkatkan efisiensi dan efektifitas analisa *log alert* IDS dengan fokus pada FN dan FP saja karena secara probabilitas telah dibantu dengan metode klasifikasi *Naïve Bayes*. Dari total 575 data yang dianalisa secara manual berkurang menjadi 85 data saja (sekitar 15%) yaitu jumlah keseluruhan FN dan FP dari data latih dan uji.

#### 5. KESIMPULAN

Dari hasil dan pembahasan dapat disimpulkan bahwa dengan teknik klasifikasi *Naïve Bayes* dan dengan *Confusion Matrix* didapat sekitar 50% data transaksi yang melalui IDS adalah intrusi atau transaksi lain yang tidak ada manfaat dan cenderung sebagai gangguan. FN dan FP memiliki peran penting untuk meningkatkan kualitas layanan serta mitigasi resiko suatu transaksi, dengan fokus pada analisa pada FN dan TP akan

memberikan 85% lebih efektif sehingga efisien waktu dan tenaga tim keamanan suatu perusahaan pengguna IDS.

Pengguna IDS perlu senantiasa melakukan pemutakhiran *rule* yang digunakan untuk meningkatkan nilai *accuracy* dan *F-Measure*. Akan lebih baik jika dibuat aplikasi berbasis penelitian ini agar memudahkan dan otomatisasi dalam perhitungan dan pemutakhiran model secara berkala.

#### DAFTAR PUSTAKA

- [1] A. Khraisat, I. Gondal, P. Vamplew, and J. Kamruzzaman, "Survey of intrusion detection systems: techniques, datasets and challenges," *Cybersecurity*, vol. 2, no. 1, p. 20, Dec. 2019, doi: [10.1186/s42400-019-0038-7](https://doi.org/10.1186/s42400-019-0038-7).
- [2] J. Liu, K. Xiao, L. Luo, Y. Li, and L. Chen, "An intrusion detection system integrating network-level intrusion detection and host-level intrusion detection," in *2020 IEEE 20th International Conference on Software Quality, Reliability and Security (QRS)*, Dec. 2020, pp. 122–129. doi: [10.1109/QRS51102.2020.00028](https://doi.org/10.1109/QRS51102.2020.00028).
- [3] M. Kumar and A. K. Singh, "Distributed Intrusion Detection System using Blockchain and Cloud Computing Infrastructure," in *2020 4th International Conference on Trends in Electronics and Informatics (ICOEI)(48184)*, Jun. 2020, pp. 248–252. doi: [10.1109/ICOEI48184.2020.9142954](https://doi.org/10.1109/ICOEI48184.2020.9142954).
- [4] R. Malani, A. B. W. Putra, and M. Rifani, "Implementation of the Naive Bayes Classifier Method for Potential Network Port Selection," *International Journal of Computer Network and Information Security*, vol. 12, no. 2, pp. 32–40, Apr. 2020, doi: [10.5815/ijcnis.2020.02.04](https://doi.org/10.5815/ijcnis.2020.02.04).
- [5] A. Alazab, M. Hobbs, J. Abawajy, A. Khraisat, and M. Alazab, "Using response action with intelligent intrusion detection and prevention system against web application malware," *Information Management & Computer Security*, vol. 22, no. 5, pp. 431–449, Nov. 2014, doi: [10.1108/IMCS-02-2013-0007](https://doi.org/10.1108/IMCS-02-2013-0007).
- [6] S. Cooper, "Intrusion Detection Systems Explained: 14 Best IDS Software Tools Reviewed," May 06, 2022. <https://www.comparitech.com/net-admin/network-intrusion-detection-tools/> (accessed Jul. 28, 2022).
- [7] B. S. Sharmila and R. Nagapadma, "Intrusion Detection System using Naive Bayes algorithm," in *2019 IEEE International WIE Conference on Electrical and Computer Engineering (WIECON-ECE)*, Nov. 2019, pp. 1–4. doi: [10.1109/WIECON-ECE48653.2019.9019921](https://doi.org/10.1109/WIECON-ECE48653.2019.9019921).
- [8] S. Anwar, F. Septian, and R. D. Septiana, "Jurnal Teknologi Sistem Informasi dan Aplikasi Klasifikasi Anomali Intrusion Detection System (IDS) Menggunakan Algoritma Naïve Bayes Classifier dan Correlation-Based Feature Selection," *Jurnal Teknologi Sistem Informasi dan Aplikasi*, vol. 2, no. 4, pp. 135–140, Oct. 2019, [Online]. Available: <http://openjournal.unpam.ac.id/index.php/JTSI/index>
- [9] A. Prasetyo, L. Affandi, and D. Arpandi, "IMPLEMENTASI METODE NAIVE BAYES UNTUK INTRUSION DETECTION SYSTEM (IDS)," *Jurnal Informatika Polinema*, vol. 4, no. 4, pp. 280–284, Aug. 2018.

[10] Y. I. Kurniawan, F. Razi, N. Nofiyati, B. Wijayanto, and M. L. Hidayat, "Naive Bayes modification for intrusion detection system classification with zero probability," *Bulletin of Electrical Engineering and Informatics*, vol. 10, no. 5, pp. 2751–2758, Oct. 2021, doi: [10.11591/eei.v10i5.2833](https://doi.org/10.11591/eei.v10i5.2833).

[11] T. Wisanwanichthan and M. Thammawichai, "A Double-Layered Hybrid Approach for Network Intrusion Detection System Using Combined Naive Bayes and SVM," *IEEE Access*, vol. 9, pp. 138432–138450, 2021, doi: [10.1109/ACCESS.2021.3118573](https://doi.org/10.1109/ACCESS.2021.3118573).

[12] A. D. Afifaturahman, F. Maulana, and S. Artikel, "Perbandingan Algoritma K-Nearest Neighbour (KNN) dan Naive Bayes pada Intrusion Detection System (IDS)," *INNOVATION IN RESEARCH OF INFORMATICS*, vol. 3, no. 1, pp. 17–25, 2021, [Online]. Available: <http://innovatics.unsil.ac.id>

[13] I. N. T. Wirawan and I. Eksistyanto, "PENERAPAN NAIVE BAYES PADA INTRUSION DETECTION SYSTEM DENGAN DISKRITISASI VARIABEL," *Jurnal Ilmiah Teknologi Informasi*, vol. 13, no. 2, pp. 182–189, Jul. 2015.

[14] A. J. Meerja, A. Ashu, and A. Rajani Kanth, "Gaussian Naïve Bayes Based Intrusion Detection System," in *Advances in Intelligent Systems and Computing*, vol. 1182 AISC, 2021, pp. 150–156. doi: [10.1007/978-3-030-49345-5\\_16](https://doi.org/10.1007/978-3-030-49345-5_16).

[15] L. Koc, T. A. Mazzuchi, and S. Sarkani, "A network intrusion detection system based on a Hidden Naïve Bayes multiclass classifier," *Expert Systems with Applications*, vol. 39, no. 18, pp. 13492–13500, Dec. 2012, doi: [10.1016/j.eswa.2012.07.009](https://doi.org/10.1016/j.eswa.2012.07.009).

[16] K. S. Bhosale, M. Nenova, and G. Iliev, "Modified Naive Bayes Intrusion Detection System (MNBIDS)," in *2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS)*, Dec. 2018, pp. 291–296. doi: [10.1109/CTEMS.2018.8769248](https://doi.org/10.1109/CTEMS.2018.8769248).

[17] Kementerian Ketenagakerjaan Republik Indonesia, "SKKNI Keahlian Artificial Intelligence (Data Science)," 2020. [https://skkni.kemnaker.go.id/tentang-skkni/dokumen?area=data\\_science&limit=20&page=1](https://skkni.kemnaker.go.id/tentang-skkni/dokumen?area=data_science&limit=20&page=1) (accessed Jul. 10, 2022).

[18] T. E. E. Tju, D. S. Maylawati, G. Munawar, and S. Utomo, "Prediction of the COVID-19 Vaccination Target Achievement with Exponential Regression," *JISA (Jurnal Informatika dan Sains)*, vol. 4, no. 2, pp. 179–182, Dec. 2021, doi: [10.31326/jisa.v4i2.1051](https://doi.org/10.31326/jisa.v4i2.1051).

[19] R. F. Smith, "Windows Security Log Encyclopedia." <https://www.ultimatewindowssecurity.com/securitylog/encyclopedia/default.aspx?i=j> (accessed Jul. 28, 2022).

[20] B. Barz and J. Denzler, "Do We Train on Test Data? Purging CIFAR of Near-Duplicates," *Journal of Imaging*, vol. 6, no. 6, p. 41, Jun. 2020, doi: [10.3390/jimaging6060041](https://doi.org/10.3390/jimaging6060041).

[21] G. I. Webb, "Naïve Bayes," in *Encyclopedia of Machine Learning and Data Mining*, Boston, MA: Springer US, 2016, pp. 1–2. doi: [10.1007/978-1-4899-7502-7\\_581-1](https://doi.org/10.1007/978-1-4899-7502-7_581-1).

[22] Z. Zhang, "Naïve Bayes classification in R," *Annals of Translational Medicine*, vol. 4, no. 12, pp. 241–241, Jun. 2016, doi: [10.21037/atm.2016.03.38](https://doi.org/10.21037/atm.2016.03.38).

[23] P. Bhatia, *Data Mining and Data Warehousing*. Cambridge University Press, 2019. doi: [10.1017/9781108635592](https://doi.org/10.1017/9781108635592).

[24] O. Caelen, "A Bayesian interpretation of the confusion matrix," *Annals of Mathematics and Artificial Intelligence*, vol. 81, no. 3–4, pp. 429–450, Dec. 2017, doi: [10.1007/s10472-017-9564-8](https://doi.org/10.1007/s10472-017-9564-8).

[25] E. Conrad, Seth Misener, and Joshua Feldman, *Eleventh Hour CISSP®*. Elsevier Science, 2016.

[26] J. Xu, Y. Zhang, and D. Miao, "Three-way confusion matrix for classification: A measure driven view," *Information Sciences*, vol. 507, pp. 772–794, Jan. 2020, doi: [10.1016/j.ins.2019.06.064](https://doi.org/10.1016/j.ins.2019.06.064).

[27] Z. Zhao and X. Wang, "Multi-segments Naïve Bayes classifier in likelihood space," *IET Computer Vision*, vol. 12, no. 6, pp. 882–891, Sep. 2018, doi: [10.1049/iet-cvi.2017.0546](https://doi.org/10.1049/iet-cvi.2017.0546).

## NOMENKLATUR

$P(Y Sx,Cx,Dx,Ex)$	Probabilitas <i>Yes</i> untuk kemunculan suatu kejadian dengan variabel <i>Sx, Cx, Dx, dan Ex</i>
$P(N Sx,Cx,Dx,Ex)$	Probabilitas <i>No</i> untuk kemunculan suatu kejadian dengan variabel <i>Sx, Cx, Dx, dan Ex</i>
$P(Sx,Cx,Dx,Ex Y)$	Probabilitas dengan variabel dari <i>Sx, Cx, Dx, dan Ex</i> di kelas <i>Yes</i>
$P(Sx,Cx,Dx,Ex N)$	Probabilitas dengan variabel dari <i>Sx, Cx, Dx, dan Ex</i> di kelas <i>No</i>
$P(Sx,Cx,Dx,Ex)$	Probabilitas kemunculan suatu kejadian dengan variabel dari <i>Sx, Cx, Dx, dan Ex</i>
$P(Sx Y)$	Probabilitas variabel dari <i>Sx</i> di kelas <i>Yes</i>
$P(Sx N)$	Probabilitas variabel dari <i>Sx</i> di kelas <i>No</i>
$P(SC Y)$	Probabilitas SC di dalam <i>Sx</i> di kelas <i>Yes</i>
$P(SC N)$	Probabilitas SC di dalam <i>Sx</i> di kelas <i>No</i>
$P(SM Y)$	Probabilitas SM di dalam <i>Sx</i> di kelas <i>Yes</i>
$P(SM N)$	Probabilitas SM di dalam <i>Sx</i> di kelas <i>No</i>
$P(Cx Y)$	Probabilitas variabel dari <i>Cx</i> di kelas <i>Yes</i>
$P(Cx N)$	Probabilitas variabel dari <i>Cx</i> di kelas <i>No</i>
$P(CI Y)$	Probabilitas <i>CI</i> di dalam <i>Cx</i> di kelas <i>Yes</i>
$P(CI N)$	Probabilitas <i>CI</i> di dalam <i>Cx</i> di kelas <i>No</i>
$P(CE Y)$	Probabilitas <i>CE</i> di dalam <i>Cx</i> di kelas <i>Yes</i>
$P(CE N)$	Probabilitas <i>CE</i> di dalam <i>Cx</i> di kelas <i>No</i>
$P(CU Y)$	Probabilitas <i>CU</i> di dalam <i>Cx</i> di kelas <i>Yes</i>
$P(CU N)$	Probabilitas <i>CU</i> di dalam <i>Cx</i> di kelas <i>No</i>
$P(Dx Y)$	Probabilitas variabel dari <i>Dx</i> di kelas <i>Yes</i>
$P(Dx N)$	Probabilitas variabel dari <i>Dx</i> di kelas <i>No</i>
$P(D1 Y)$	Probabilitas <i>D1</i> di dalam <i>Dx</i> di kelas <i>Yes</i>
$P(D1 N)$	Probabilitas <i>D1</i> di dalam <i>Dx</i> di kelas <i>No</i>
$P(D2 Y)$	Probabilitas <i>D2</i> di dalam <i>Dx</i> di kelas <i>Yes</i>
$P(D2 N)$	Probabilitas <i>D2</i> di dalam <i>Dx</i> di kelas <i>No</i>
$P(D3 Y)$	Probabilitas <i>D3</i> di dalam <i>Dx</i> di kelas <i>Yes</i>
$P(D3 N)$	Probabilitas <i>D3</i> di dalam <i>Dx</i> di kelas <i>No</i>
$P(D4 Y)$	Probabilitas <i>D4</i> di dalam <i>Dx</i> di kelas <i>Yes</i>
$P(D4 N)$	Probabilitas <i>D4</i> di dalam <i>Dx</i> di kelas <i>No</i>
$P(Ex Y)$	Probabilitas variabel dari <i>Ex</i> di kelas <i>Yes</i>
$P(Ex N)$	Probabilitas variabel dari <i>Ex</i> di kelas <i>No</i>
$P(E1 Y)$	Probabilitas <i>E1</i> di dalam <i>Ex</i> di kelas <i>Yes</i>
$P(E1 N)$	Probabilitas <i>E1</i> di dalam <i>Ex</i> di kelas <i>No</i>
$P(E2 Y)$	Probabilitas <i>E2</i> di dalam <i>Ex</i> di kelas <i>Yes</i>
$P(E2 N)$	Probabilitas <i>E2</i> di dalam <i>Ex</i> di kelas <i>No</i>
$P(E3 Y)$	Probabilitas <i>E3</i> di dalam <i>Ex</i> di kelas <i>Yes</i>

$P(E3 N)$	Probabilitas $E3$ di dalam $Ex$ di kelas $No$
$P(E4 Y)$	Probabilitas $E4$ di dalam $Ex$ di kelas $Yes$
$P(E4 N)$	Probabilitas $E4$ di dalam $Ex$ di kelas $No$
$P(E5 Y)$	Probabilitas $E5$ di dalam $Ex$ di kelas $Yes$
$P(E5 N)$	Probabilitas $E5$ di dalam $Ex$ di kelas $No$
$P(E6 Y)$	Probabilitas $E6$ di dalam $Ex$ di kelas $Yes$
$P(E6 N)$	Probabilitas $E6$ di dalam $Ex$ di kelas $No$
$P(Y)$	Probabilitas kelas $Yes$
$P(N)$	Probabilitas kelas $No$
$Sx$	<i>Severity Level</i> yang bisa berisi $SC$ atau $SM$
$SC$	<i>Severity Level = Critical</i>
$SM$	<i>Severity Level = Medium</i>
$Cx$	<i>IP Category</i> yang berisi $CI$ , $CE$ , atau $CU$
$CI$	<i>IP Category = Internal</i>
$CE$	<i>IP Category = External</i>
$CU$	<i>IP Category = Unknown</i>
$Dx$	<i>IP Destination</i> yang berisi $D1$ , $D1$ , $D3$ , $D4$
$D1$	<i>IP Destination = x.y.2.250</i>
$D2$	<i>IP Destination = x.y.2.253</i>
$D3$	<i>IP Destination = x.z.46.130</i>
$D4$	<i>IP Destination = x.z.57.65</i>
$Ex$	<i>Event</i> yang berisi $E1$ , $E2$ , $E3$ , $E4$ , $E5$ , atau $E6$
$E1$	<i>Event = 4624</i>
$E2$	<i>Event = 4625</i>
$E3$	<i>Event = 4634</i>
$E4$	<i>Event = 4662</i>
$E5$	<i>Event = 4782</i>
$E6$	<i>Event = XSS</i>

## BIODATA PENULIS

### MUHAMMAD KAMIL SURYADEWIANSYAH



Lahir di Jakarta pada akhir tahun 2000, saat ini sebagai mahasiswa tahap akhir di Universitas Budi Luhur, Fakultas Teknologi Informasi, Program Studi Teknik Informatika dengan peminatan Cyber Security. Terpilih sebagai mahasiswa kelas unggulan Teknik Informatika angkatan 2019, mencoba menggabungkan antara *cyber security* dengan *machine learning*.