

Klasifikasi Traffic Network Dengan Menggunakan Naive Bayes dan Feature Selection Dalam Pemilihan Atribut

by Alfin Ramdhani

Submission date: 24-Jul-2019 02:34PM (UTC+0700)

Submission ID: 1154575955

File name: BAB_1-6.pdf (2.06M)

Word count: 4482

Character count: 53853

PENDAHULUAN

1.1 Pendahuluan

Bab ini membahas mengenai latar belakang penelitian Klasifikasi *Traffic Network* dengan Menggunakan *Naïve Bayes* dan *Feature Selection* dalam pemilihan atribut yang akan dibahas secara umum dan singkat pada latar belakang.

1.2 Latar Belakang

Kebutuhan koneksi internet semakin hari semakin meningkat, trafik internet pun meningkat. Dengan trafik yang semakin tinggi, maka akses/koneksi internet akan semakin berat/lambat. Sehingga perlu diketahui bagaimana pola trafik internet yang ada selama ini. Pola tersebut berguna untuk dijadikan dasar kebijakan manajemen koneksi internet untuk saat sekarang dan waktu yang akan datang, bermanfaat juga untuk mengetahui ada tidaknya pola yang tidak wajar yang bisa jadi mengarah ke serangan dari luar yang semakin membebani jaringan dalam aspek keamanan (Manshaei et al., 2013). Selain itu, pola yang didapatkan bisa menunjukkan aktifitas pengguna sehari-hari seperti apa, yaitu aplikasi internet apa saja yang mayoritas dimanfaatkan oleh pengguna selama ini. Hal tersebut berkaitan dengan tujuan utama dan prioritas dari ketersediaan internet. Sehingga jangan sampai, internet lebih banyak dimanfaatkan untuk hal-hal di luar tujuan utamanya (Sigit, 2016).

Jadi, dengan adanya deskripsi permasalahan diatas, maka kita harus mengklasifikasikan lalu lintas jaringan, untuk memonitor kapan lalu lintas sedang sibuk dan ketika lalu lintas sedikit aktivitas. Lihat jenis kegiatan apa yang sering dilakukan seseorang ketika terhubung ke internet dengan *ranking system* (*browsing, chatting, videocalling, streaming, dll*). Serta mengatasi *bottle-*

necking atau pengiriman data yang tertunda beberapa saat. Salah satu metode yang akan digunakan dalam pengklasifikasian adalah metode *Naive Bayes*.

Naive Bayes adalah bagian dari Jaringan Syaraf Tiruan. Klasifikasi Naive Bayes akan mengklasifikasikan jumlah catatan trafik jaringan, dan juga menggunakan kombinasi metode *Feature Selection* dalam penyeleksian atribut pada data trafik jaringan.

Naive Bayes (NB) adalah salah satu metode paling awal yang digunakan untuk klasifikasi trafik internet, yang sederhana dan cukup efektif untuk mengklasifikasi peluang (Sigit, 2016). Serta karena performa dan kecepatan yang tinggi dalam proses klasifikasi, dan mudah untuk menghasilkan probabilitas posterior data yang di tes terhadap kelasnya (Schlosser et al., 2009).

Sedangkan *Feature Selection* adalah metode tambahan yang terdiri dari *Entropy* dan *Information Gain*. Kedua fitur tersebut memungkinkan penulis dalam proses penyeleksian data trafik yang digunakan seperti, pengurangan atribut yang tidak banyak mendukung dalam penelitian ini, mendapatkan nilai bobot dari setiap atribut, dan pemeringkatan atribut yang banyak muncul di data trafik.

Penelitian sebelumnya yang dilakukan oleh (Okililas & Tasmi, 2017) menggunakan metode DPI dalam mengelompokkan jenis trafik, sedang penelitian yang dilakukan oleh (Bujlow, Carela-Español, & Barlet-Ros, 2015) membandingkan tool DPI dalam mengelompokkan jenis trafik internet.

Solusi dalam mengklasifikasi trafik pada network telah banyak dilakukan dengan menghasilkan solusi yang aktif dan pasif sebagai solusi yang ditawarkan, seperti penelitian yang dilakukan oleh (Molina et al., 2012) menyatakan bahwa management network sebagai media pendukung dalam kasus identifikasi paket data dengan pendekatan *Operation, Administration, Maintenance & Provisioning* dan juga penelitian yang dilakukan oleh (Zhang et al., 2016) mereka berhasil mengenali pola-pola paket dengan baik, namun sistem masih bersifat pasif sehingga tidak ada control trafik yang keluar masuk.

Dhote (2015) pada penelitian survei seleksi fitur untuk klasifikasi trafik internet menyatakan, seleksi fitur dapat membantu memahami data, mengurangi perhitungan, mengurangi efek *curse of dimensionality*, meningkatkan kinerja dan mengurangi waktu komputasi. Sedangkan Aliakbarian (2013) menyatakan ekstraksi mampu membuat data ekstrak baru dengan menghilangkan korelasi yang menghasilkan klasifikasi trafik tersebut optimal. Berdasarkan hal tersebut, maka penelitian ini akan membangun perangkat lunak untuk klasifikasi data trafik internet dengan menggunakan metode *Naive Bayes* dan

Feature Selection (entropy, information gain) sebagai pemilihan atribut, sehingga dihasilkan klasifikasi data dimensi tinggi dengan akurasi yang baik.

1.3 Rumusan Masalah

Rumusan permasalahan yang diselesaikan dalam penelitian ini adalah sebagai berikut:

1. Bagaimana mekanisme *Naive Bayes* dan *Feature Selection* untuk klasifikasi data?
2. Bagaimana hasil dari implementasi *Naive Bayes* dan *Feature Selection* pada klasifikasi data trafik internet?

1.4 Tujuan Penelitian

Tujuan penelitian ini adalah mengetahui nilai atribut yang banyak muncul pada data trafik dengan *Confusion Matrix* sesuai dengan *class* yang telah ditetapkan pada *Naive Bayes* dan *Feature Selection*.

1.5 Manfaat Penelitian

Manfaat dari penelitian ini adalah sebagai berikut:

1. Memahami *Naive Bayes Classification* dan *Feature Selection* sebagai metode klasifikasi data trafik;
2. Mampu menerapkan teknik klasifikasi dan seleksi *Naive Bayes Classification* dan *Feature Selection* pada klasifikasi data trafik;

3. Hasil penelitian dapat digunakan untuk referensi dalam penelitian lain yang sejenis yang menggunakan seleksi fitur *entropy* dan *information gain* dalam metode *Naïve Bayes*.

1.6 Batasan Masalah

Batasan masalah pada penelitian ini adalah sebagai berikut:

1. Data yang digunakan berupa data trafik internet yang diunduh dari situs *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk/>) dalam bentuk .xls.
2. Metode klasifikasi yang digunakan adalah *Naïve Bayes*.
3. Evaluasi kualitas pengklasifikasian dilakukan dengan *Confusion Matrix*.
4. Validasi data trafik menggunakan *10-fold Cross Validation*.

1.7 Sistematika Penulisan

Penyusunan skripsi ini disusun dengan sistematika penulisan sebagai berikut:

BAB I. PENDAHULUAN

Bab ini membahas mengenai latar belakang, perumusan masalah, tujuan penelitian, manfaat penelitian, batasan masalah, dan sistematika penulisan.

BAB II. KAJIAN LITERATUR

Pada bab ini akan membahas dasar-dasar teori yang akan digunakan dalam penelitian, seperti pengetahuan dasar tentang klasifikasi dan metode yang akan digunakan dalam proses klasifikasi data trafik internet.

BAB III. METODOLOGI PENELITIAN

Pada bab ini akan dibahas mengenai unit penelitian, tahapan yang akan dilaksanakan pada penelitian ini, tahapan proses secara umum, metode pengembangan perangkat lunak, teknik pengujian dan manajemen proyek penelitian.

BAB IV. PENGEMBANGAN PERANGKAT LUNAK

Pada bab ini membahas mengenai analisis dan perancangan perangkat lunak yang akan digunakan sebagai alat penelitian. Dimulai dari pengumpulan dan analisis kebutuhan, rancangan dan konstruksi perangkat lunak serta pengujian untuk memastikan semua kebutuhan pengembangan perangkat lunak sesuai dengan dengan kebutuhan. Penyusunan pada bab ini memiliki kerangka penulisan dengan fase-fase dan elemen-elemen pengembangan perangkat lunak bersifat berorientasi objek.

BAB V. HASIL DAN ANALISA PENELITIAN

Pada bab ini diuraikan hasil pengujian berdasarkan langkah-langkah yang telah direncanakan. Tabel hasil pengujian serta analisisnya disajikan sebagai basis dari kesimpulan yang akan diambil dalam penelitian ini.

BAB VI. KESIMPULAN DAN SARAN

Pada bab ini berisi kesimpulan dari semua uraian-uraian pada bab-bab sebelumnya dan juga saran-saran yang diharapkan berguna untuk pengembangan selanjutnya.

1.8 Kesimpulan

Penelitian mengenai pengklasifikasian data trafik jaringan akan dilakukan dengan metode *Naïve Bayes* dan *Feature Selection* dalam pemilihan atribut. Tujuannya adalah untuk mengembangkan perangkat lunak yang mampu mengklasifikasi data trafik internet.

BAB II

TINJAUAN PUSTAKA

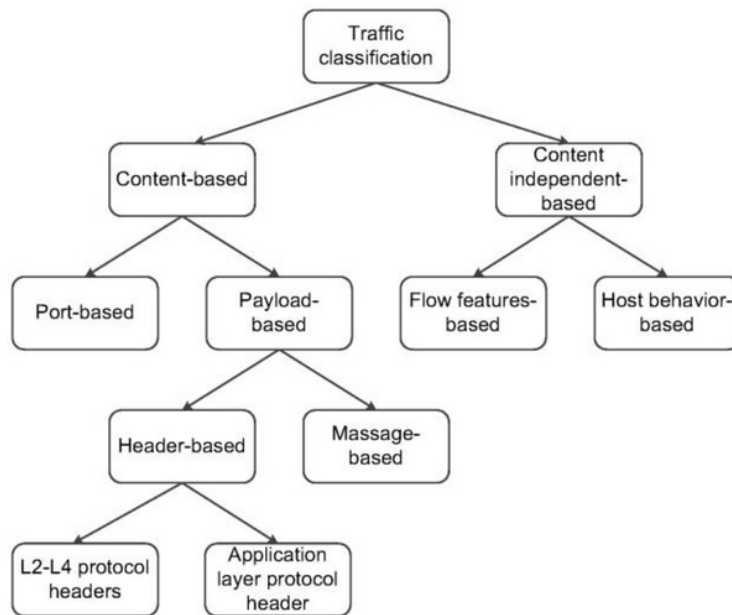
2.1 Pendahuluan

Pada bab ini berisi tentang teori-teori dasar yang berhubungan dengan penelitian. Dalam bab ini dijelaskan secara singkat mengenai klasifikasi trafik, *Naïve Bayes*, *Entropy* dan *Information Gain* dalam *Feature Selection*, *Confusion Matrix*, dan *Cross Validation*.

Pada akhir bab ini disertakan penelitian terdahulu yang relevan dengan penelitian ini dan kesimpulan.

2.2 Klasifikasi Trafik

Penelitian sebelumnya oleh (Siqueira et al., 2007) menyatakan bahwa klasifikasi trafik adalah satu metode yang digunakan untuk mengoptimalkan *bandwidth* dengan tujuan koneksi internet yang handal dan stabil. Metode-metode klasifikasi trafik dapat dikelompokkan dalam metode *Port-Based*, *Payload Based*, *Statistical Classification*. Pada gambar II-1 menunjukkan dalam proses klasifikasi trafik ada dua metode yaitu *content-based* dan *content independent-based*, dimana di bagian *content-based* terdiri dari dua metode yaitu (*port-based* dan *payload-based*) sedangkan dibagian *Content Independent-based* terdiri dari *Flow features-based* dan *Host behavior-based*.



Gambar II-1. Proses Klasifikasi Trafik (Siqueira et al., 2007)

2.2.1 Naïve Bayes

Menurut (Dhivya & P.Shanmugaraja, 2015) *Naïve Bayes* adalah algoritma yang berbasis *Bayesian theorem* serta menggunakan perhitungan probabilitas dalam menentukan kelas. Algoritma ini disebut juga sebagai mode fitur independen karena nilai dari atribut pada sebuah kelas tidak tergantung pada nilai atribut yang lain.

Persamaan dari Naïve Bayes sebagai berikut:

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} \quad (\text{II-1})$$

Keterangan:

- X = Kriteria suatu kasus berdasarkan masukan
- C_i = Kelas solusi pola ke-
 i , dimana i adalah jumlah label kelas
- $P(C_i|X)$ = Probabilitas kriteria masukan X dengan label ke-
las C_i
- $P(C_i)$ = Probabilitas label kelas C

Naïve bayes merupakan salah satu pembelajaran *supervised* dalam *machine learning*, artinya dalam tahap pembelajaran dibutuhkan data training untuk mengambil keputusan. Pada proses klasifikasi akan dihitung nilai probabilitas dari masing-masing atribut dalam sebuah kelas. Atribut yang mempunyai nilai probabilitas yang tinggi akan dijadikan label dari kelas tersebut.

Pada proses klasifikasi *Bayesian theorem* diperlukan beberapa penjelasan untuk menentukan kelas yang baik bagi sample yang akan dianalisis, oleh karena itu persamaan *Bayesian theorem* untuk klasifikasi ditampilkan pada persamaan (II-2).

$$p(C|F_1, \dots, F_n) = \frac{p(C)p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)} \quad (\text{II-2})$$

Variabel C menjelaskan kelas, variabel $F_1 \dots F_n$ merupakan karakteristik-

karakteristik petunjuk yang dibutuhkan untuk melakukan klasifikasi. Pe

rsamaan (II-

2) menjelaskan bahwa peluang masuknya sampel dengan karakteristik tertentu dalam kelas C (*posterior*) adalah peluang munculnya kelas C (*prior*) dikali dengan peluang kemunculan karakteristik-

karakteristik sampel pada kelas C (*likelihood*), dibagi dengan peluang kemunculan karakteristik-karakteristik sampel secara global (*evidence*).

2.3 Feature Selection Untuk Pemilihan Atribut

Fungsi dari proses seleksi atribut ini adalah pertama untuk menentukan ranking dari setiap atribut, yang kedua adalah untuk menghilangkan atribut yang tidak relevan. Pada penelitian oleh (Slocum, 2012) menyatakan bahwa pada dasarnya *entropy* digunakan untuk menafsirkan ketidakpastian dari beberapa atribut dari sebuah dataset. Semakin tinggi *entropy* sebuah atribut maka nilai ketidakpastian semakin tinggi, sedangkan penelitian oleh (Novaković, Strbac, & Bulatović, 2011) untuk mencari nilai *entropy* menggunakan persamaan (II-3) sebagai berikut:

$$H(Y) = - \sum_{y \in Y} (y) 2(p(y)) \quad (\text{II-3})$$

Dimana, (py) adalah fungsi probabilitas marginal untuk nilai variabel acak Y. Jika nilai-nilai Y yang diukur dalam dataset S dibagi dengan nilai-nilai fitur kedua X, dan *entropy* Y terhadap partisi dipengaruhi oleh X kurang dari *entropy* Y sebelum proses partisi, maka ada relasi antar fitur Y dan X, maka persamaan *entropy* Y setelah mengamati X ditunjukkan pada persamaan II.4.

$$H(Y|X) = - \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log_2 p(y|x)$$

(II-4)

$p(y|x)$ adalah probabilitas bersyarat dari y terhadap x . Entropy sebagai kriteria tidak valid dalam *training dataset* S , maka dapat didefinisikan ukuran merupakan informasi tambahan tentang Y yang disediakan oleh X yang mewakili jumlah dimana nilai entropy Y menurun. Bagian ini dikenal sebagai *Information Gain (IG)*.

$$IG = H(Y) - H(Y|X) = H(X) - H(X|Y) \quad (II-5)$$

Dari persamaan II-

5, informasi yang diperoleh tentang Y setelah mengamati X sama dengan informasi yang diperoleh tentang X setelah mengamati Y .

2.4 Cross Validation

Cross Validation merupakan teknik yang dapat digunakan untuk melakukan validasi terhadap model klasifikasi. Teknik ini dapat digunakan untuk menguji performa dan mengukur seberapa akurat model klasifikasi yang telah dibuat. Salah satu teknik *cross validation* adalah *k-fold cross validation*. *K-fold cross validation* adalah metode validasi di mana dataset dibagi menjadi k bagian, dimana standar yang biasa digunakan untuk memperoleh estimasi kesalahan terbaik adalah 10 bagian atau *ten-fold cross validation* (Gorunescu, 2011).

2.5 Confusion Matrix

Classifier Accuracy Measures (Han dan Kamber, 2006 : 360) adalah metode klasifikasi yang dilakukan berdasarkan tingkat akurasi model dalam melakukan prediksi. Hal ini dilakukan karena keakuratan dalam mengolah data merupakan salah satu hal yang penting.

Metode yang digunakan untuk menguji tingkat akurasi model klasifikasi ini adalah *Confusion matrix*. Dalam metode ini memberikan keputusan yang diperoleh dalam *training* dan *testing*, *Confusion Matrix* memberikan penilaian *performance* klasifikasi berdasarkan objek dengan benar atau salah. *Confusion matrix* berisi informasi aktual (*actual*) dan prediksi (*predicted*) pada sistem klasifikasi.

Tabel II-1. Contoh Tabel Confusion Matrix

Classification	Predicted Class		
		Class = Yes	Class = No
Actual Class	Class = Yes	a(true positive-TP)	b(false negative-FN)
	Class = No	c(false positive-FP)	d(true negative-TN)

Keterangan:

- True Positive (TP) = Data dengan kelas positif diklasifikasikan ke positif.

- True Negative (TN) = Data dengan kelas negatif diklasifikasikan ke negatif.
- False Positive (FP) = Data dengan kelas negatif diklasifikasikan ke positif.
- False Negative (FN) = Data dengan kelas positif diklasifikasikan ke negatif.

Dimana setelah itu akan dilakukan perhitungan akurasi, precision, dan recall sebagai berikut:

$$\text{akurasi} = \frac{TP+TN}{TP+TN+FP+FN} \quad (\text{II-6})$$

$$\text{Precision} = \frac{TP}{TP+FN} \quad (\text{II-7})$$

$$\text{Recall} = \frac{TP}{TP+FP} \quad (\text{II-8})$$

2.6 Penelitian Lain Yang Relevan

Pada bagian ini dipaparkan beberapa penelitian yang telah dilakukan oleh beberapa peneliti lain. Dimaksudkan untuk memperkuat penalaran dan rasionalitas keterlibatan sejumlah variabel pada penelitian ini. Selain itu juga difungsikan sebagai pendapat ilmiah yang dipadukan dengan hasil kajian pustaka untuk membangun kerangka berpikir peneliti dalam kaitannya dengan masalah yang sedang diteliti.

2.6.1. Li Jun, Zhang Shunyi, Lu Yanqing, Zhang Zailong (2007)

Li Jun, Zhang Shunyi, Lu Yanqing, dan Zhang Zailong melakukan penelitian berjudul *Internet Traffic Classification Using Machine Learning*, Penelitian tersebut dilakukan pada Post and Telecommunications, Nanjing University, Nanjing, China dan Zhejiang Wanli University, Ningbo, China. Penelitian mengenai pengaruh teknik reduksi dimensi dan klasifikasi dilakukan oleh Jun et al. (2014) yang menerapkan teknik reduksi dimensi yaitu *genetic algorithm*(GA) pada metode klasifikasi Naïve Bayes, Naïve Bayes Tree, C4.5 Decision Tree, Random Forests, dan KNN. Penelitian dilakukan dua kali, di mana pertama hanya melakukan klasifikasi terhadap data full set dan kedua menggunakan data yang sudah menggunakan reduksi. Hasil penelitian tersebut reduksi fitur GA tidak terlalu mempengaruhi akurasi metode klasifikasi yang ada, namun meningkatkan *modelling time* dan *training time*.

2.6.1. K. Keerthi Vasan dan B. Surendiran (2016)

K. Keerthi Vasan dan B. Surendiran melakukan penelitian tentang *Dimensionality reduction using Principal Component Analysis for network intrusion detection*. Penelitian tersebut dilakukan di Department of Computer Science, National Institute of Technology Puducherry, Karaikal, India. Penerapan kombinasi *Principal Component Analysis* (PCA) untuk teknik reduksi dimensi dan teknik klasifikasi data in

truksi jaringan, yang dilakukan oleh (Vasan & Surendiran 2016) didapatkan dimensi ideal yakni 10 dimensi pertama. Penelitian ini menggunakan dua patokan kumpulan data jaringan, KDD Cup dan UNB IS CX dengan hasil rasio reduksi dimensi masing-masing adalah 0,24 dan 0,36. Sedangkan akurasi klasifikasi sebesar 99,7% dan 98,8%, hampir sama dengan akurasi yang diperoleh dengan menggunakan data asli 41 fitur untuk KDD dan 28 fitur untuk IS CX. Hasil penelitian ini membuktikan teknik reduksi dimensi dapat digunakan untuk membentuk jumlah dimensi ideal serta menghasilkan pengelompokan data yang baik.

2.8 Kesimpulan

Penelitian ini akan menerapkan metode *Naïve Bayes* untuk klasifikasi data trafik internet, Entropy dan Information Gain sebagai fitur seleksi, data akan divalidasi menggunakan *10-fold cross validation*, dan dievaluasi menggunakan *Confusion Matrix*.

¹ BAB III

METODOLOGI PENELITIAN

3.1 Pendahuluan

Bab ini akan menjelaskan unit dan tahapan penelitian yang diimplementasikan, metodologi penelitian serta manajemen proyek penelitian. Tahapan penelitian dijadikan sebagai acuan pada setiap fase pengembangan dan memberikan sebuah solusi untuk rumusan masalah dan mencapai tujuan penelitian.

3.2 Unit Penelitian

Dalam penelitian ini yang digunakan sebagai unit penelitian adalah situs ² *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk/>) yang menyediakan layanan dataset untuk penelitian.

¹ 3.3 Data

3.3.1. Jenis dan Sumber Data

Data yang digunakan sebagai obyek penelitian adalah jenis data sekunder berupa data trafik internet, yang bersumber dari situs ² *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk/>). Dataset yang digunakan terdiri 19384 data (*record*), 248 fitur, termasuk fitur kelas yang didefinisikan menggunakan teks (*String*), dimana nama-nama fitur dapat dilihat pada lampiran

I. Terdapat 10 kelas pada data trafik internet yaitu, Services, WWW, Attack, P2P, FTP-Pasv, Multimedia, FTP-Control, FTP-Data, dan Interactive.

3.3.2. Metode Pengumpulan Data

Data didapatkan secara manual dengan mengunduh data trafik internet di situs ² *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk/>) dalam bentuk file .arff, kemudian menggunakan weka untuk mengubah tipe data menjadi .csv

3.3 Tahapan Penelitian

Untuk mendapatkan hasil penyeleksian atribut yang paling banyak muncul dalam data trafik yang digunakan dengan *Feature Selection* dan pengklasifikasian menggunakan *Naïve Bayes*, maka penelitian ¹ dengan tahapan-tahapan yang akan dijelaskan pada subbab 3.4.1 sampai dengan 3.4.6.

3.3.1. Menetapkan Kerangka Kerja

Data akan di klasifikasi dengan *Naïve Bayes* dan akan diseleksi dengan *Feature Selection*.

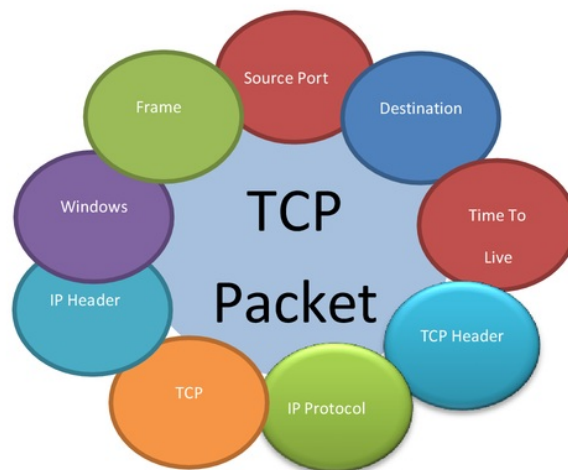
a. Praproses

Tahap praproses merupakan tahapan dalam mengelola data masukan. Tahapan pertama dalam praproses pada perangkat lunak ini adalah *Data*

Cleaning. Dalam penelitian ini, proses *Data Cleaning* akan membersihkan data yang tidak lengkap fiturnya (*missing value*) dengan cara mengganti nilai-nilai yang hilang. Nilai-nilai yang hilang jika bersifat kontinu akan diganti dengan rata-rata fitur, dan jika nilai yang hilang bersifat *yes/no* akan diganti dengan modus dari fitur tersebut.

b. Feature Extraction

Pada tahap ini adalah proses bagaimana dapat mengambil informasi-informasi atau atribut-atribut dari *data traffic* yang dibutuhkan seperti ditunjukkan pada gambar III-1. Pengambilan atau proses mengekstrak *file* menggunakan program *weka*. *Capture paket* sangat dibutuhkan dalam penelitian ini karena paket-paket yang lewat akan diproses. *Extract Packet* ini adalah proses untuk mendapatkan informasi-informasi dan nilai-nilai atribut dari data trafik.



Gambar III.1 Arsitektur atribut TCP Header (Ford et al., 2011)

c. Klasifikasi Naïve Bayes dan Feature Selection

¹ Kelas atau label yang digunakan pada klasifikasi data trafik internet adalah Mail, Services, WWW, Attack, P2P, FTP-Pasv, Multimedia, FTP-Control, FTP-Data, dan Interactive. Proses klasifikasi dilakukan oleh metode *Naïve Bayes*, dan penyeleksian atribut dilakukan oleh *entropy* dan *information gain* yang akan menghasilkan keluaran perankingan jenis atribut yang paling banyak keluar dari data trafik.

d. Evaluasi dan validasi menggunakan *Confusion Matrix* dan *10-fold Cross Validation*

Confusion Matrix digunakan untuk mencari nilai *true positive*, *true negative*, *false positive*, dan *false negative*. Jika label yang digunakan adalah Mail, Services, WWW, Attack, P2P, FTP-Pasv, Multimedia, FTP-Control, FTP-Data, dan Interactive, maka pada *confusion matrix* dihitung nilai *true Mail*, *true Services*, *true WWW*, *true Attack*, *true P2P*, *true FTP-pasv*, *true Multimedia*, *true FTP-control*, *true FTP-data*, *true Interactive*, *false Mail*, *false Services*, *false WWW*, *false Attack*, *false P2P*, *false FTP-pasv*, *false Multimedia*, *false FTP-control*, *false FTP-data*, dan *false Interactive*. Nilai-nilai tersebut akan digunakan untuk menghitung akurasi, *precision*, dan *recall*. Model klasifikasi divalidasi menggunakan *10-fold cross validation* ¹ dimana *dataset* dibagi menjadi 10 bagian dengan satu bagian digunakan sebagai data pengujian dan sisanya digunakan sebagai data pelatihan yang dilakukan sebanyak 10 kali dengan

bagian data pengujian yang berbeda-beda. Tiap pengujian diukur menggunakan *confusion matrix*.

3.3.2. Menetapkan Kriteria Pengujian

Pembahasan mengenai tahapan ini akan dijelaskan pada bab IV. Pada tahapan pengujian awal penelitian, data terlebih dahulu akan diklasifikasi dengan menggunakan metode *Naïve Bayes* dan diseleksi oleh *Entropy* dan *information gain* dengan validasi data menggunakan *10-fold cross validation*. Proses klasifikasi akan dilanjutkan dengan perhitungan mencari nilai akurasi, *precision*, dan *recall*, serta hasil akhir berupa data per kelas dalam bentuk tabel oleh *Confusion Matrix*.

3.3.3. Menetapkan Format Data Pengujian

Hasil perhitungan dalam pengujian, akan dijelaskan sebagai berikut:

Hasil pengujian klasifikasi dengan tabel *confusion matrix* pada tiap bagian data akan digambarkan dalam tabel III-1

Tabel III-1. Rancangan Tabel *Confusion Matrix* Untuk Setiap Hasil Pengujian

Cross validation		Prediksi			
	Label	MAIL	SERVICE	...	INTERACTIVE
	MAIL				
	SERVICE				
	...				
	INTERACTIVE				

Hasil pengujian klasifikasi akan digambarkan pada tabel III-2.

Tabel III-2. Rancangan Tabel Hasil Pengujian Untuk Setiap Hasil Klasifikasi

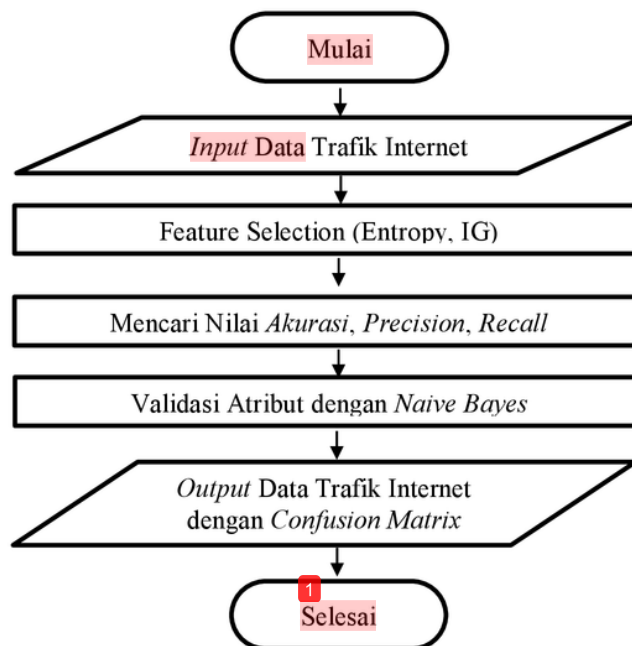
Cross Validation	Akurasi	Precision	Recall
1			
....			
10			
Rata-rata			

3.3.4 Menentukan Alat yang Digunakan dalam Pelaksanaan Penelitian

Untuk melaksanakan penelitian mengenai klasifikasi data trafik menggunakan Naïve Bayes dan Feature Selection dalam pemilihan atribut

¹ dibutuhkan alat penelitian. Oleh karena itu, penulis akan mengembangkan sebuah perangkat lunak yang dapat melakukan proses klasifikasi data trafik internet.

3.3.5 ¹ Melakukan Pengujian Penelitian



Gambar III-2. Tahapan Pengujian Penelitian

3.3.6 Melakukan Analisa Hasil Pengujian dan Membuat Kesimpulan

Untuk mengetahui hasil penelitian mengenai klasifikasi data trafik menggunakan *Naive Bayes* dan *Feature Selection* dalam pemilihan atribut, ¹ maka cara yang dilakukan adalah dengan membandingkan hasil pengujian akurasi, *precision*, dan *recall* dari Metode *Naive Bayes* dan *Feature Selection*. Analisis

¹ hasil pengujian dapat dilihat pada tabel III-3. ¹ Setelah mendapatkan hasil analisis pengujian penelitian, maka langkah selanjutnya adalah membuat kesimpulan penelitian yang akan dijelaskan pada bab V.

3.4 Metode Pengembangan Perangkat Lunak

Metodologi yang diterapkan dalam pengembangan perangkat lunak sebagai alat penelitian tugas akhir ini berorientasi pada objek menggunakan metode *Rational Unified Process* (RUP). Secara umum, langkah-langkah yang akan dilakukan pada pengembangan perangkat lunak adalah fase inepsi, elaborasi, konstruksi, dan transisi.

3.5 ¹ Manajemen Proyek Penelitian

Manajemen proyek merupakan perencanaan aktivitas penelitian dari tahap inisialisasi masalah sampai dengan pada tahap kesimpulan dari penelitian. Adapun kegiatan-kegiatan yang berlangsung selama penelitian dapat dilihat dalam *Work Breakdown Structure* (WBS) pada lampiran IV.

3.6 Kesimpulan

Penelitian ini menggunakan data trafik internet dari situs ² *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk/>). Data tersebut akan melalui tahap praproses sebelum di proses. Kemudian data terlebih dahulu akan diklasifikasi dengan metode *Naïve Bayes*. Lalu data diseleksi oleh *Feature Selection* dengan *Entropy* dan *Information Gain*, sehingga akan diperoleh hasilnya dengan membuat peringkat setiap atribut yang banyak muncul.

BAB IV

PENGEMBANGAN PERANGKAT LUNAK

4.1 Pendahuluan

Pada bab III disebutkan bahwa diperlukan sebuah alat berupa perangkat lunak yang digunakan dalam pelaksanaan penelitian ini, maka penulis mengembangkan perangkat lunak dengan metode pemrograman berorientasi obyek berdasarkan panduan *Rational Unified Process*. Di dalam *Rational Unified Process* terdapat empat fase pengembangan perangkat lunak yaitu fase insepsi, elaborasi, konstruksi, dan transisi. Setiap fase terdiri dari pemodelan bisnis, kebutuhan, analisis dan desain, implementasi, dan pengujian. Pada bab ini dibahas proses pengembangan perangkat lunak yang digunakan sebagai alat penelitian.

4.2 Fase Insepsi

Tahapan pertama dalam pengembangan perangkat lunak ialah melakukan identifikasi terhadap sistem yang dikembangkan. Aktivitas yang dilakukan pada fase ini meliputi penentuan *user requirement* dan fungsionalitas perangkat lunak pada permodelan bisnis, mengumpulkan data penelitian pada kebutuhan, membuat *use-case diagram* pada analisis dan desain, mendokumentasikan *user requirement*, fungsionalitas perangkat lunak, dan *use-*

case diagram pada implementasi, serta memastikan *user requirement* dan fungsionalitas perangkat lunak valid pada pengujian.

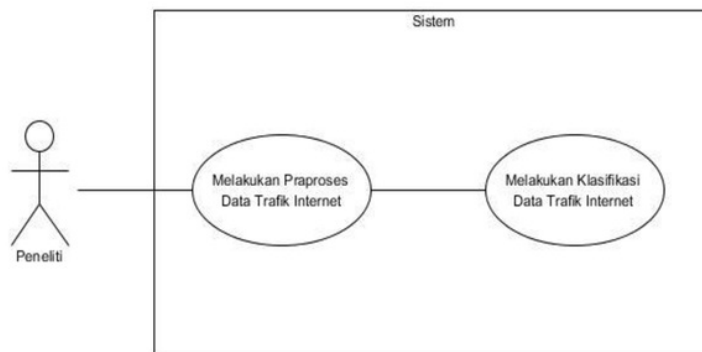
4.2.1 Permodelan Bisnis

Andersen dan Feamster (2006) menyatakan trafik internet menghasilkan sejumlah besar data yang membantu pengelola jaringan untuk mengelola dan mempelajari berbagai karakteristik trafik internet. Untuk mempermudah mengidentifikasi data trafik internet, trafik internet dapat dikelompokkan berdasarkan isi trafik internet tersebut. Namun pengelompokan secara manual menimbulkan permasalahan sebagai berikut:

1. Dibutuhkan ketelitian yang tinggi untuk menganalisis satu persatu yang melewati trafik;
2. Membutuhkan waktu yang lama untuk mengelompokkan;
3. Dibutuhkan banyak sumber daya manusia.

Dari permasalahan tersebut, dibutuhkan pengelompokan otomatis terhadap data trafik internet dengan melakukan klasifikasi. Pada gambar IV-

1 menunjukkan proses pengelompokkan data trafik internet menggunakan klasifikasi.



Gambar IV – 1. Diagram *Use Case Current Existing*

Namun, jika jumlah data yang akan diklasifikasi berjumlah sangat banyak, maka algoritma klasifikasi biasa tidak mampu menghasilkan performa klasifikasi yang optimal karena jumlah data yang banyak akan menghasilkan fitur data yang sangat banyak. Tidak semua fitur relevan digunakan pada proses klasifikasi dan juga dapat menurunkan performa klasifikasi. Untuk itu, dibutuhkan teknik reduksi dimensi untuk memaksimalkan kinerja algoritma klasifikasi.

Perangkat lunak yang dibangun merupakan perangkat lunak berbasis desktop yang mampu melakukan klasifikasi data trafik internet secara otomatis. Masukan pada perangkat lunak yang dikembangkan berupa berkas bertipe *.xls* yang berisikan konten data trafik internet. Keluaran yang dihasilkan berupa tabel berisi data trafik internet yang telah terkelompok berdasarkan kelasnya.

4.2.2 Kebutuhan Sistem

Berdasarkan permodelan bisnis yang telah dilakukan, spesifikasi kebutuhan perangkat lunak yang akan dikembangkan memiliki fitur prapengolahan dan klasifikasi.

4.2.2.1 Fitur Prapengolahan

Perangkat lunak dilengkapi fitur prapengolahan yang digunakan oleh pengguna untuk menyiapkan data trafik internet agar dapat di reduksi dan di klasifikasi. Prapengolahan yang akan dilakukan adalah *replace missing value*.

4.2.2.2 Fitur Klasifikasi

Fitur berikutnya adalah klasifikasi data trafik internet. Selanjutnya proses perhitungan akan berlanjut pada tahapan pengklasifikasian trafik internet menggunakan *Naïve Bayes* dari data trafik.

4.2.2.3 Fitur Seleksi

Fitur ini adalah fitur yang digunakan untuk menseleksi atribut data trafik yang banyak dengan atribut yang hanya berpengaruh dalam melakukan penelitian ini karena tidak semua atribut dari dataset yang dipakai cocok dengan topik dari penelitian.

4.2.2.4 Fitur Evaluasi dan Validasi

Fitur ini merupakan fitur yang digunakan untuk memvalidasi dataset menggunakan metode 10-*fold cross validation* dan mengukur akurasi, *precision*, dan recall menggunakan *confusion matrix*.

Untuk dapat merealisasikan fitur-fitur tersebut, perangkat lunak harus memenuhi kebutuhan fungsional dan non-fungsional. Kebutuhan fungsional menjelaskan kebutuhan atau fasilitas utama perangkat lunak yang dibangun pada tabel IV-1.

1. Sedangkan, kebutuhan non fungsional menjelaskan kebutuhan atau fasilitas yang tidak wajib dimiliki oleh perangkat lunak, dalam artian hanya merupakan pelengkap agar perangkat lunak lebih baik kinerjanya yang dapat dilihat pada tabel IV-2.

Tabel IV – 1 Kebutuhan Fungsional

No.	Kebutuhan
1	Perangkat lunak dapat memuat data trafik internet
2	Perangkat lunak dapat melakukan praproses data trafik internet
3	Perangkat lunak dapat melakukan proses klasifikasi data trafik internet

Tabel IV – 2 Kebutuhan Non Fungsional

No.	Kebutuhan
1	Perangkat lunak dapat menampilkan pesan kesalahan jika terdapat aksi pengguna yang salah dan konfirmasi tindakan pengguna.
2	Perangkat lunak memiliki antarmuka yang mudah dimengerti dan digunakan oleh pengguna.

4.2.3 Analisis dan Desain

Kegiatan yang dilakukan pada tahapan ini adalah menganalisis kebutuhan perangkat lunak, data, prapengolahan, klasifikasi dengan *Naive Bayes*, dan evaluasi menggunakan *confusion matrix*.

4.2.3.1 Analisis Kebutuhan Perangkat Lunak

Dari pemodelan bisnis yang telah dijabarkan, untuk menyelesaikan permasalahan yang terjadi diperlukan perangkat lunak yang mampu mengelompokkan data trafik internet secara otomatis. Untuk itu, perangkat lunak harus memiliki kemampuan sebagai berikut:

1. Melakukan prapengolahan berupa *replace missing value* terhadap data trafik internet;
2. Melakukan klasifikasi data trafik internet secara otomatis menggunakan *Naive Bayes* dan *Feature Selection*;

3. Melakukan validasi terhadap data trafik internet menggunakan *10-fold cross validation*;
4. Melakukan evaluasi terhadap hasil klasifikasi menggunakan *confusion matrix*.

Pengembangan perangkat lunak dimulai dengan mengumpulkan data trafik internet yang diunduh dari situs *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk>). Selanjutnya file data trafik internet disimpan dalam bentuk .xls. Selanjutnya berkas tersebut diimpor ke perangkat lunak lalu akan dilakukan prapengolahan berupa *replace missing value*. Tahap selanjutnya ialah melakukan klasifikasi *Naive Bayes*.

4.2.3.2 Analisis Prapengolahan

Tahap prapengolahan yang dilakukan pada data trafik internet yaitu *replace missing value*. Contoh data trafik internet yang digunakan sebanyak 7 buah, masing-masing memiliki 4 fitur.

Tabel IV – 3 Contoh Data Trafik Internet

F1	F2	F3	F4	KELAS
55	1	1	3	MAIL
57	1	1	2	MAIL
49	?	1	3	MAIL

49	1	?	3	MAIL
5	?	1	0	WWW
5	1	0	1	WWW
6	?	1	0	WWW

Berikut langkah-langkah prapengolahan yang dilakukan pada ketujuh data dapat dilihat dibawah ini:

1. Dari data yang telah dimasukkan, diperiksa tiap fiturnya;
2. Ditemukan, fitur F2 terdapat 3 missing value dan F3 terdapat 1 *missing value* dari 7 data trafik internet;
3. Hitung rata-rata fitur yang terdapat missing value kemudian ubah missing value tersebut dengan rata-rata fitur tersebut.

Tabel IV – 4 Hasil Prapengolahan

F1	F2	F3	F4	KELAS
55	1	1	3	MAIL
57	1	1	2	MAIL
49	1	1	3	MAIL

49	1	0.83	3	MAIL
5	1	1	0	WWW
5	1	0	1	WWW
6	1	1	0	WWW

4.2.3.3 Desain Perangkat Lunak

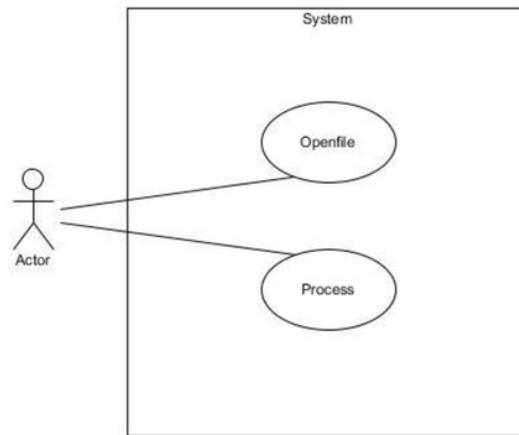
Desain perangkat lunak dideskripsikan dengan diagram *use case* dan diagram aktivitas.

1. *Use Case*

Pada subbab ini dijelaskan gambaran fungsionalitas perangkat lunak yang dibangun dengan menggunakan pemodelan *Use Case*.

a) Diagram *Use Case*

Diagram *Use Case* menjelaskan secara umum kegiatan yang dilakukan oleh aktor terhadap perangkat lunak yang dapat dilihat pada gambar IV-2.



Gambar IV-2. Diagram Use Case

b) Tabel Definisi Aktor

Dalam penelitian ini, seluruh pengguna menjadi aktor yang dijelaskan pada tabel IV –

5. Definisi dari aktor dapat dilihat pada tabel IV-5.

Tabel IV-5 Definisi Aktor Use Case

Nomor	Aktor	Definisi
1	Peneliti	Peneliti adalah orang yang berhubungan dengan perangkat lunak aplikasi untuk memasukkan data, melatih algoritma NB, melakukan seleksi fitur menggunakan Feature Selection, melakukan klasifikasi, dan menampilkan hasil dari pengujian yang dilakukan.

c) Tabel Definisi Use Case

Definisi *Use case* yang dijelaskan pada tabel IV-

6 merupakan penjelasan dari kerja perangkat lunak secara spesifik pada gambar IV-

2 disimbolkan dengan sebuah notasi lingkaran yang agak lonjong atau berbentuk oval.

Tabel IV – 6 Definisi *Use Case*

Nomor	Use Case	Deskripsi
1	Melakukan Praproses Data Trafik Internet	Aktivitas memasukkan data trafik internet dalam format .xls yang kemudian dilakukan praproses <i>replace missing value</i> .
2	Melakukan Klasifikasi dengan NB	Kegiatan ini digunakan untuk melakukan klasifikasi data trafik internet yang telah diproses menggunakan metode NB.

d) Skenario *Use Case*

Skenario adalah merupakan urutan spesifik dari aksi dan interaksi antara aktor dan sistem. Berikut ini adalah skenario dari *use case* yang telah didefinisikan sebelumnya.

1. Skenario *Use Case* Melakukan Klasifikasi Dengan *Naive Bayes*

No : 001

Nama Use Case: Melakukan Klasifikasi Dengan *Naive Bayes* dan *Feature Selection*

Aktor : Pengguna

Tujuan : Mengetahui kelompok data trafik internet

Deskripsi: *Use Case* ini digunakan untuk melakukan klasifikasi dengan metode *Naive Bayes* dan *Feature Selection*.

Kondisi Awal: Terdapat data trafik internet yang siap diproses.

Kondisi Akhir: Hasil klasifikasi data internet, nilai akurasi, hasil seleksi atribut data.

Skenario *use case* menambahkan data trafik internet dijelaskan pada tabel IV-7

Tabel IV-

7 Skenario *Use Case* Melakukan Klasifikasi dengan *Naive Bayes* dan *Feature Selection*.

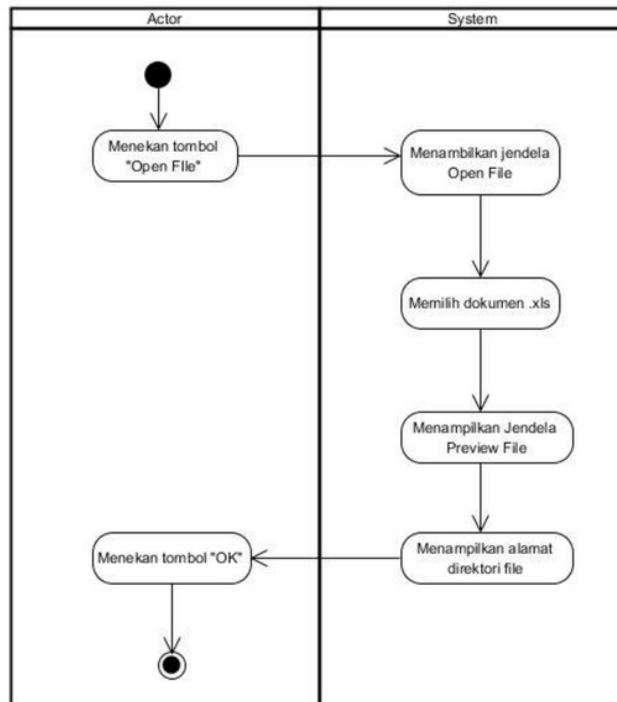
Aktor	Sistem
1. Pengguna menekan tombol "Open File"	
	2. Menampilkan kotak dialog yang menunjukkan direktori penyimpanan data trafik internet
3. Pengguna memilih data yang memiliki format .xls untuk diproses	
	4 Menampilkan jendela yang menunjukkan preview file yang akan diklasifikasi
6. Pengguna menekan tombol "Process"	
	7. Melakukan <i>10 fold Cross Validation</i>
	8. Klasifikasi menggunakan NB

	9. Penyeleksian atribut dengan <i>Entropy</i> dan <i>IG</i>
	10. Menampilkan hasil klasifikasi dan nilai akurasi, <i>precision</i> , <i>recall</i> , serta nilai hasil seleksi atribut dari <i>Confusion Matrix</i> dalam tabel

2. Diagram Aktivitas

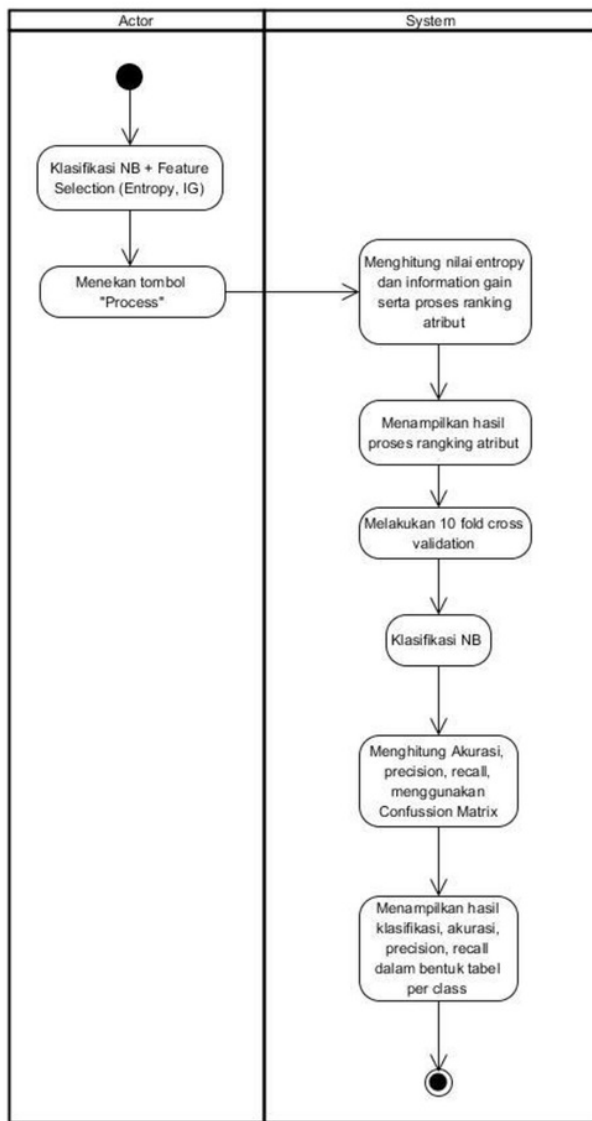
Diagram aktivitas menggambarkan aliran kerja atau aktivitas dari sebuah sistem atau proses bisnis. Diagram aktivitas melakukan praproses dapat dilihat pada gambar IV-

3, diagram aktivitas melakukan klasifikasi dengan Naive Bayes dan Feature Selection dapat dilihat pada gambar IV-4.



Gambar IV-

3. Diagram Aktivitas *Use Case* melakukan Praproses Data Trafik Internet



Gambar IV-

4. Diagram Aktivitas Use Case Melakukan Klasifikasi dengan *Naive Bayes* dan *Feature Selection*

4.3 Fase Elaborasi

Tahapan kedua dalam pengembangan perangkat lunak adalah melakukan identifikasi terhadap sistem yang dikembangkan. Aktivitas yang dilakukan mencakup perancangan data, perancangan antarmuka, identifikasi kebutuhan, perumusan kebutuhan pengujian, pemodelan diagram sequence, dan pembuatan dokumentasi.

4.3.1 Permodelan Bisnis

Pada subbab ini akan dibahas mengenai perancangan perangkat lunak yang akan dibangun. Perancangan dilakukan berdasarkan hasil analisis yang telah dilakukan pada fase insepisi. Perancangan yang dibahas pada subbab ini meliputi perancangan data dan perancangan antar muka.

4.3.1.1 Perancangan Data

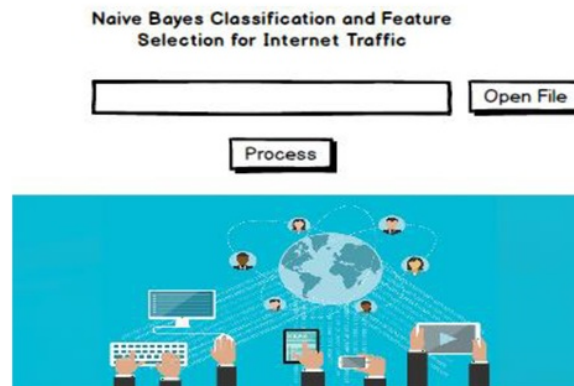
Perangkat lunak yang akan dibangun memiliki kemampuan klasifikasi terhadap data. Adapun data yang akan melalui proses klasifikasi adalah data trafik internet yang disimpan dalam file berformat .xls.

4.3.1.2 Perancangan Antar Muka

Pada subbab ini membahas tentang perancangan antar muka dari perangkat lunak yang dibangun. Adapun rancangan antarmuka Menu Utama digambarkan pada gambar IV-5, rancangan antarmuka *File Preview* digambarkan pada gamb

ar IV-

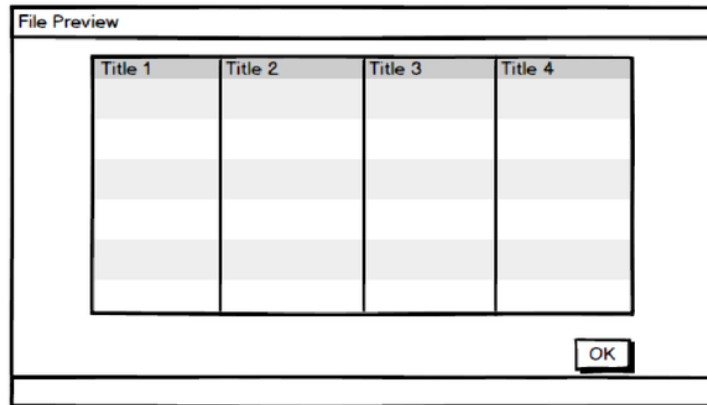
6, dan rancangan antarmuka *Output View* pada gambar IV-7.



Gambar IV-5. Rancangan Antarmuka Menu Utama

Gambar IV-

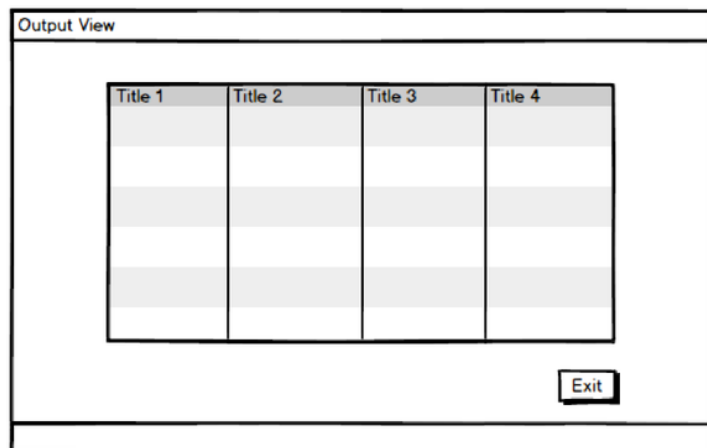
5 merupakan halaman utama yang berfungsi untuk melakukan klasifikasi data trafik internet. Terdapat empat bagian pada antarmuka ini yaitu tombol Open File untuk membuka file data .xls, *text field* untuk menampilkan direktori data trafik, tombol Process untuk memproses klasifikasi data trafik, dan tabel yang menampilkan output data per kelas.



Gambar IV-6. Rancangan antarmuka *File Preview*

Gambar IV-

6 merupakan tampilan yang berfungsi untuk menampilkan isi file .xls yang akan diproses. Terdapat dua bagian yaitu tabel yang akan menampilkan isi file tersebut dan tombol “OK” untuk mengkonfirmasi apakah file yang diuji sudah sesuai.



Gambar IV-7. Rancangan antarmuka *Output View*

Gambar IV-

7 merupakan tampilan yang berfungsi untuk menampilkan nilai output yang telah diproses. Terdapat dua bagian yaitu tabel yang akan menampilkan nilai tersebut dan tombol “Exit” untuk keluar dari jendela *Output View*.

4.3.2 Kebutuhan Sistem

Pada subbab ini dibahas mengenai kebutuhan sistem dari perangkat lunak yang dibangun berdasarkan hasil analisis dan perancangan pada tahap selanjutnya. Untuk membangun perangkat lunak dalam penelitian ini dibutuhkan perangkat keras (*hardware*), perangkat lunak (*software*) dan bahasa pemrograman. Bahasa pemrograman yang digunakan untuk implementasi perangkat lunak adalah Java, Perangkat keras yang digunakan pada tahap pengembangan dan penelitian ini adalah laptop dengan spesifikasi sebagai berikut:

1. Laptop merk ASUS GL553VE;
2. Processor Intel(R) Core(TM) i7-7700HQ CPU @ 2.80GHz;
3. RAM 8 GB;
4. Hard Disk 1 TB.

Sedangkan perangkat lunak yang digunakan untuk implementasi yaitu :

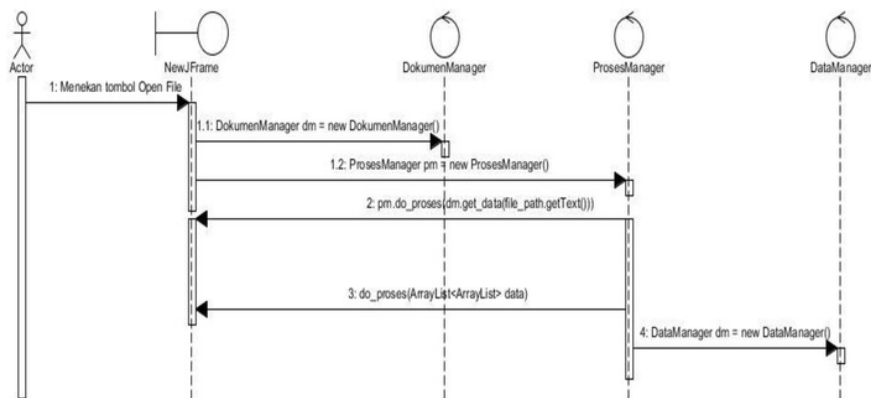
1. Windows 10 Home Single Language 64-bit;
2. Compiler Netbeans IDE 8.2;

3. Visual Paradigm UML Enterprise Edition V8.0.

4.3.3 Diagram Sequence

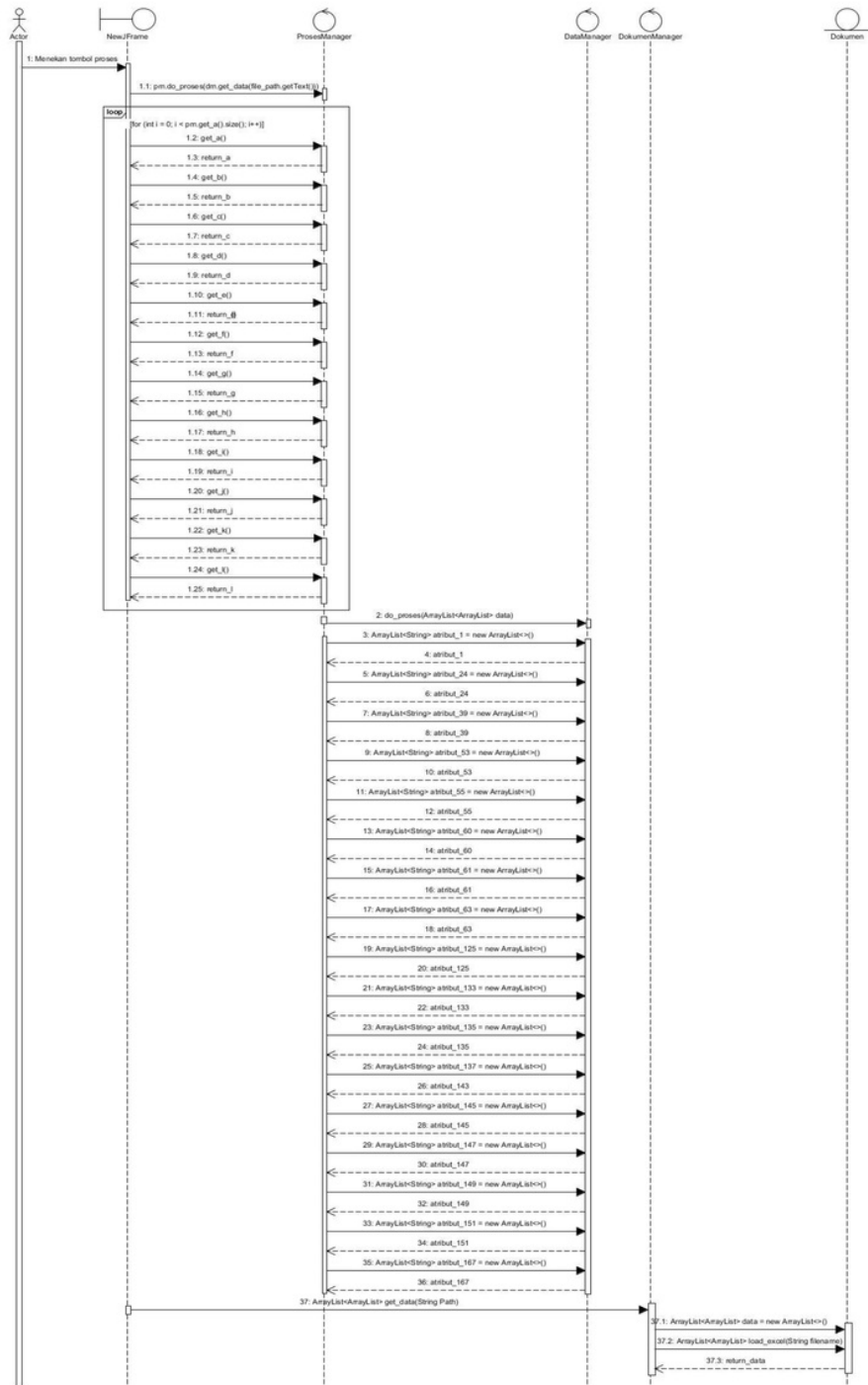
Diagram sequence adalah diagram yang menggambarkan kolaborasi dinamis antara sebuah objek. Berdasarkan use case yang ditentukan, diagram sequence yang dibentuk pada pengembangan perangkat lunak ini berjumlah 2 buah. Diagram *sequence* melakukan praproses Data Trafik Internet dapat dilihat pada gambar IV-

8, diagram sequence melakukan klasifikasi dengan Naive Bayes dan Feature Selection dapat dilihat pada gambar IV-9.



Gambar IV-

8. *Sequence Diagram* Melakukan Praproses Data Trafik Internet



Gambar IV-

9. *Sequence Diagram* Melakukan Klasifikasi dengan *Naïve Bayes*

4.4 Fase Konstruksi

Fase konstruksi berfokus pada pengembangan perangkat lunak baik komponen utama maupun fitur-fitur pendukung dengan melakukan sederet iterasi. Di setiap iterasi terdapat proses analisis, desain, implementasi, dan pengujian. Dalam proses pengembangannya dapat menggunakan konstruksi paralel agar mempercepat hasil perangkat lunak. Hasil yang diharapkan dari fase ini adalah sebuah produk perangkat lunak yang siap digunakan oleh end-user, yaitu sebuah produk perangkat lunak yang dapat digunakan sebagai alat penelitian.

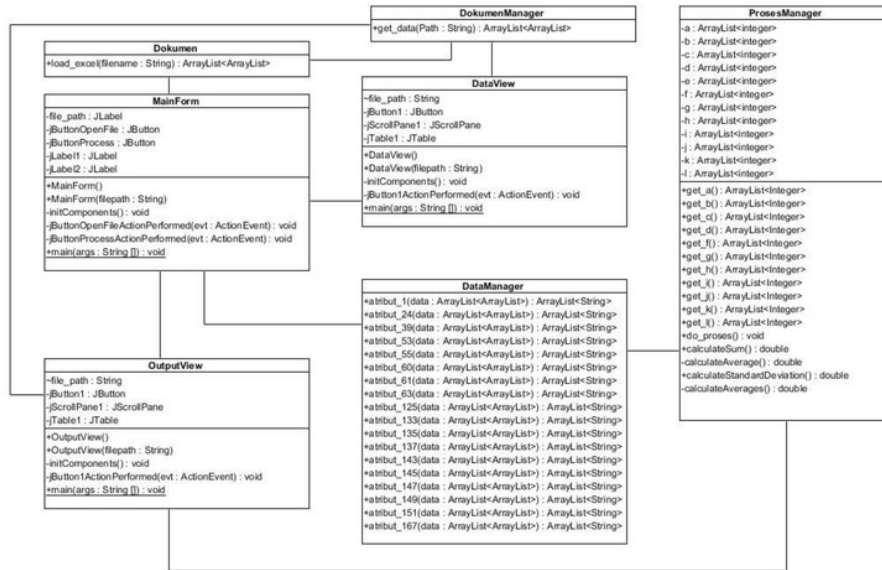
4.4.1 Kebutuhan Sistem

Dalam pengembangan perangkat lunak ini, penulis menggunakan beberapa library. Diantaranya adalah *weka*, dan *Jxl. Library jxl* digunakan untuk memproses data dari file *.xls*, metode klasifikasi *Naïve Bayes* menggunakan *weka*.

4.4.2 Diagram Kelas

Diagram kelas adalah diagram UML yang menggambarkan kelas-kelas dalam sebuah sistem dan hubungannya antara satu dengan yang lain. Terdapat serta dimasukkan pula atribut dan operasi. Terdapat 7 kelas yang terdiri dari 3 kelas *boundary* (kelas *MainForm*, *DataView*

w, OutputView), 3 kelas *control* (kelas DataManager, DokumenManager, ProsesManager), dan 1 kelas *entity* (kelas Dokumen).



Gambar IV-10. Class Diagram

4.4.3 Implementasi

Fase implementasi dalam konstruksi adalah mengembangkan perangkat lunak berdasarkan diagram kelas dan rancangan antarmuka yang telah dibuat dalam fase sebelumnya.

4.4.3.1 Implementasi Kelas

Kelas-

kelas yang telah dirancang pada diagram kelas diimplementasikan dalam bahasa pemrograman Java. Tabel (IV

-

8) menunjukkan implementasi kelas dalam bahasa Java.

Tabel IV-8. Implementasi Kelas

No	Nama Kelas	Nama File	Keterangan
1	MainForm	MainForm.java	Kelas MainForm adalah kelas <i>boundary</i> yang merupakan form utama dan menyediakan akses untuk melakukan masukkan data dan memproses klasifikasi.
2	DataView	DataView.java	Kelas DataView adalah kelas <i>boundary</i> yang bertujuan untuk menampilkan preview dari data (file .xls) yang akan diuji.
3	OutputView	OutputView.java	Kelas DataView adalah kelas <i>boundary</i> yang bertujuan untuk menampilkan <i>output</i> dari nilai atribut yang paling banyak muncul sesuai <i>class</i> yang ditetapkan
4	DokumenManager	DokumenManager.java	Kelas DokumenManager

			merupakan kelas <i>control</i> yang menangani proses pengambilan data dari file excel.
5	Dokumen	Dokumen.java	Kelas Dokumen merupakan kelas <i>entity</i> yang bertujuan untuk menarik data mentah/data training.
6	ProsesManager	ProsesManager.java	Kelas DokumenManager merupakan kelas <i>control</i> yang menangani proses perhitungan klasifikasi.
7	DataManager	DataManager.java	Kelas DataManager merupakan kelas <i>control</i> yang menangani pemisahan data training dan data testing.

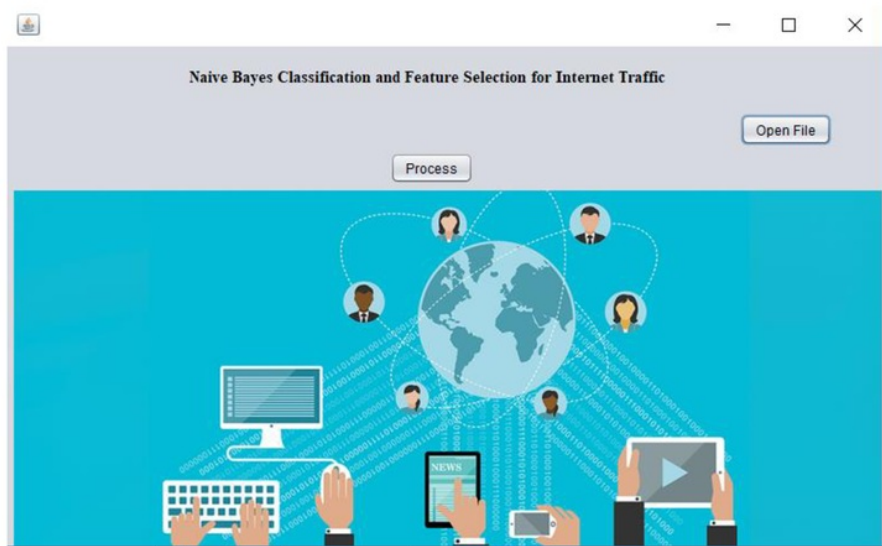
4.4.3.2 Implementasi Antarmuka

Implementasi antarmuka dilakukan berdasarkan perancangan antarmuka pada fase elaborasi. Gambar IV-

11 merupakan antarmuka menu utama utama perangkat lunak, gambar IV-

12 merupakan antarmuka *File Preview*, gambar IV-

13 merupakan antarmuka *OutputView*.



Gambar IV-11. Antarmuka Menu Utama Perangkat Lunak

1	24	39	53	55	60	61	63	125	133	135	137	143	145	147	149	151	167
25.0	677...	20.0	0.0	102...	425.0	0.0	0.0	52.0	0.02...	9.12...	13.0	0.00...	6.47...	0.0	55.0	1.0	722.0
25.0	674...	20.0	0.0	104...	434.0	0.0	0.0	54.0	0.01...	4.13...	13.0	0.00...	9.64...	0.0	57.0	1.0	722.0
25.0	700...	20.0	0.0	101...	441.0	0.0	0.0	46.0	0.02...	7.02...	5.0	0.00...	3.4...	0.0	61.0	1.0	722.0
25.0	699...	20.0	0.0	103...	428.0	0.0	0.0	47.0	0.01...	4.2E...	6.0	0.00...	0.00...	0.0	62.0	1.0	722.0
25.0	715...	18.0	0.0	101...	429.0	0.0	0.0	41.0	0.00...	5.79...	1.0	0.00...	0.0	0.0	67.0	1.0	722.0
25.0	720...	16.0	0.0	101...	440.0	0.0	0.0	41.0	0.01...	0.00...	1.0	8.0E...	0.0	0.0	66.0	1.0	722.0
25.0	707...	20.0	0.0	101...	434.0	0.0	0.0	43.0	0.02...	5.06...	6.0	0.00...	0.00...	0.0	65.0	1.0	722.0
25.0	714...	20.0	0.0	100...	425.0	0.0	0.0	43.0	0.02...	4.8E...	1.0	0.00...	0.0	0.0	63.0	1.0	722.0
25.0	702...	20.0	0.0	101...	445.0	0.0	0.0	46.0	0.02...	4.74...	13.0	0.00...	3.83...	0.0	61.0	1.0	722.0
25.0	688...	20.0	0.0	104...	431.0	0.0	0.0	50.0	0.03...	5.97...	17.0	9.1E...	2.49...	0.0	60.0	1.0	722.0
25.0	718...	17.0	0.0	100...	442.0	0.0	0.0	41.0	0.02...	5.58...	1.0	6.57...	0.0	0.0	65.0	1.0	722.0
25.0	694...	20.0	0.0	102...	428.0	0.0	0.0	48.0	0.02...	4.53...	13.0	8.71...	2.97...	0.0	60.0	1.0	722.0
25.0	710...	20.0	0.0	104...	437.0	0.0	0.0	45.0	0.02...	5.2E...	12.0	0.00...	3.5...	0.0	65.0	1.0	722.0
25.0	650...	20.0	0.0	101...	433.0	0.0	0.0	57.0	0.00...	0.00...	15.0	6.22...	2.16...	0.0	50.0	1.0	722.0
25.0	702...	20.0	0.0	104...	445.0	0.0	0.0	47.0	0.01...	6.11...	8.0	0.00...	3.31...	0.0	64.0	1.0	722.0
25.0	702...	20.0	0.0	100...	431.0	0.0	0.0	44.0	0.01...	3.91...	8.0	9.71...	5.93...	0.0	63.0	1.0	722.0
25.0	721...	15.0	0.0	102...	438.0	0.0	0.0	41.0	0.02...	5.56...	1.0	0.00...	0.0	0.0	68.0	1.0	722.0
25.0	657...	20.0	0.0	100...	428.0	0.0	0.0	57.0	0.01...	0.00...	19.0	0.00...	5.75...	0.0	50.0	1.0	722.0
25.0	696...	20.0	0.0	103...	428.0	0.0	0.0	48.0	0.01...	6.87...	13.0	9.22...	3.67...	0.0	61.0	1.0	722.0
25.0	81.5...	20.0	0.0	109.0	363.0	0.0	0.0	7.0	0.02...	2.7E...	1.0	0.02...	0.0	0.0	0.0	1.0	80.0

Gambar IV-12. Antarmuka *File Preview*

WWW	MAIL	FTP-CO...	FTP-PASV	ATTACK	P2P	DATABA...	FTP-DATA	MULTIM...	SERVIC...	INTERA...	GAMES
16786	50	0	17	2	24	0	0	13	0	0	0
1400	170	24	2	0	1	0	0	21	0	0	0
37	1	10	0	0	0	0	0	0	0	0	0
1	0	0	108	0	0	0	0	0	0	0	0
131	2	0	0	0	0	0	0	1	0	0	0
8	0	0	71	0	12	0	0	3	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	342	0	0	0	21	1	0	0	0
0	0	0	30	0	2	0	0	10	0	0	0
81	1	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

Gambar IV-13. Antarmuka *OutputView*

4.5 Fase Transisi

Pada subbab ini dibahas mengenai pengujian dari perangkat lunak klasifikasi data trafik internet yang telah dibangun. Pengujian dilakukan berdasarkan perangkat lunak hasil pengembangan fase konstruksi.

4.5.1 Permodelan Bisnis

Pengujian perangkat lunak secara *blackbox* dan *whitebox* dengan terlebih dahulu membuat rencana pengujian berdasarkan use case yang dibuat pada fase inisiasi.

4.5.2 Kebutuhan Sistem

Lingkungan pengujian yang digunakan pada fase transisi adalah perangkat keras yang sama saat membangun perangkat

at lunak klasifikasi data trafik internet dengan spesifikasi sebagai berikut:

1. Laptop merk ASUS GL553VE;
2. Processor Intel(R) Core(TM) i7-7700HQ CPU @ 2.80GHz;
3. RAM 8 GB;
4. Hard Disk 1 TB.

Sedangkan perangkat lunak yang digunakan untuk implementasi yaitu:

1. Windows 10 Home Single Language 64-bit;
2. Compiler Netbeans IDE 8.2;
3. Visual Paradigm UML Enterprise Edition V8.0.

4.5.3 Rencana Pengujian

Rencana pengujian pada perangkat lunak klasifikasi data trafik internet digambarkan dalam tabel-tabel. Kolom pada tabel meliputi identifikasi, pengujian, jenis pengujian, serta tingkat pengujian.

4.5.3.1 Rencana Pengujian Use Case Melakukan Proses Data Trafik Internet

Tabel IV-

9 menerangkan rencana pengujian memasukkan data perangkat lunak berdasarkan *Use Case*.

Tabel IV-

9. Rencana Pengujian *Use Case* Melakukan Praproses Data

No	Identifikasi	Pengujian	Jenis Pengujian	Tingkat Pengujian
1	U-1-101	Masukkan folder yang berisi file data trafik Internet berekstensi .xls.	<i>Black Box</i>	Pengujian Unit
2	U-1-102	Melakukan keseluruhan prapengolahan data trafik Internet.	<i>White Box</i>	Pengujian Unit

4.5.3.2 Rencana Pengujian *Use Case* Melakukan Klasifikasi dengan Naïve Bayes

Tabel IV-

10 menerangkan rencana pengujian melakukan klasifikasi dengan Naïve Bayes pada perangkat lunak berdasarkan *Use Case*.

Tabel IV-

10. Rencana Pengujian *Use Case* Melakukan Klasifikasi dengan Naïve Bayes

No	Identifikasi	Pengujian	Jenis Pengujian	Tingkat Pengujian
1	U-2-101	Menekan tombol "Process"	<i>Black Box</i>	Pengujian Unit
2	U-2-102	Melakukan klasifikasi data trafik internet dengan metode Naïve Bayes dan Feature Selection.	<i>White Box</i>	Pengujian Unit

4.5.4 Implementasi

Berikut ini adalah kasus uji yang dilakukan terhadap perangkat lunak yang dibangun. Kasus uji dilakukan berdasarkan rencana uji yang telah dipaparkan sebelumnya.

4.5.4.1 Pengujian Use Case memasukkan Data Trafik Internet

Tabel IV-

11 menerangkan pengujian memasukkan Data Trafik Internet pada perangkat lunak berdasarkan Use Case.

Tabel IV-

11. Pengujian Use Case Melakukan Praproses Data Trafik Internet

Identifikasi	Deskripsi	Prosedur Pengujian	Masukan	Keluaran yang Diharapkan	Hasil yang Didapat	Kesimpulan
U-1-101	Masukkan file data trafik internet berekstensi .xls.	Menekan tombol “ Open File ”	Tidak ada	Data yang sudah diproses oleh sistem.	List data trafik internet dari file .xls.	Terpenuhi
U-1-102	Masukkan file data trafik internet bukan berekstensi .xls.	Menekan tombol “ Open File ”	Tidak ada	Data tidak berhasil dimuat. Kemudian, sistem menampilkan “ File yang dimasukkan harus berekstensi .xls”.	Data tidak berhasil dimuat. Kemudian, sistem menampilkan “File yang dimasukkan harus berekstensi .xls”.	Terpenuhi

**4.5.4.2 Pengujian Use Case Melakukan Klasifikasi dengan
Naïve Bayes**

Tabel IV-

12 menerangkan pengujian klasifikasi Data Trafik Internet
dengan C4.5 pada perangkat lunak berdasarkan Use Case.

Tabel IV-12. Rencana Pengujian Use Case Melakukan Klasifikasi dengan Naïve Bayes dan Feature Selection

Identifikasi	Deskripsi	Prosedur Pengujian	Masukan	Keluaran yang Diharapkan	Hasil yang Didapat	Kesimpulan
U-2-101	Menekan tombol "Process".	Menekan tombol "Process".	Trigger untuk mengaktifkan tombol	Tombol dapat ditekan dan mengaktifkan perhitungan klasifikasi	Tombol dapat ditekan dan mengaktifkan perhitungan klasifikasi	Terpenuhi
U-2-102	Melakukan perhitungan klasifikasi Data Tersebut dengan metode Naïve Bayes dan Feature Selection	Menekan tombol "Process"	Masukan berupa hasil pengolahan	Hasil klasifikasi sebanyak jumlah kelas yang sudah diketahui di awal	Hasil klasifikasi sebanyak jumlah kelas yang sudah diketahui di awal	Terpenuhi

4.6 Kesimpulan

Dari proses pengembangan perangkat lunak, telah jelas diuraikan tahapan pengembangan alat penelitian yang akan membantu proses penelitian klasifikasi data trafik internet. Pengembangan perangkat lunak ini telah disesuaikan dengan kebutuhan penelitian, yaitu pemrosesan data trafik internet, mampu melakukan *replace missing value*, melakukan proses klasifikasi dengan *Naïve Bayes* dan *Feature Selection* dalam pemilihan atribut, menghitung nilai akurasi, *precision*, *recall* klasifikasi. Dijelaskan pula skenario dan alur pengembangan dan pengujian perangkat lunak sehingga dihasilkan perangkat lunak yang sesuai dengan kebutuhan penelitian. Selanjutnya setelah perangkat lunak selesai, dilakukan proses pengujian terhadap data uji serta melakukan analisa dari hasil yang dihasilkan oleh alat penelitian.

BAB V

HASIL DAN ANALISIS PENELITIAN

5.1 Pendahuluan

Pada bab IV telah dilakukan pengembangan perangkat lunak yang menjadi alat penelitian klasifikasi *Network Traffic* menggunakan metode *Naive Bayes* dan *Feature Selection* dalam pemilihan atribut. Untuk membuktikan hasil klasifikasi, maka digunakan perhitungan akurasi, *precision*, dan *recall* untuk mengevaluasi hasil klasifikasi yang diperlukan untuk menyelesaikan proses klasifikasi. Penjelasan dan metode perhitungan akurasi, *precision*, dan *recall* dapat dilihat pada bab II.

5.2 Percobaan Penelitian

Pengujian dilakukan dengan menggunakan data uji berupa dokumen jurnal sebanyak 19384 buah dengan 248 atribut, data diunduh dari situs *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk>) yang kemudian disalin dan disimpan dalam dokumen berekstensi *.xls*. Proses pengujian dilakukan sesuai dengan arsitektur perangkat lunak yang telah dibicarakan pada subbab 4.2.3.1 hingga 4.2.3.3, yakni analisis kebutuhan perangkat lunak prapengolahan, dan desain perangkat lunak.

Setelah didapatkan hasil klasifikasi dari metode klasifikasi *Naive Bayes*, tahapan penelitian dilanjutkan dengan menganalisa hasil klasifikasi yang telah didapatkan.

5.3 Hasil *Feature Extraction*

Bagian ini akan membaca atribut-atribut yang dihasilkan dari data trafik, dimana file data tersebut merupakan *raw data* yang sulit dibaca oleh manusia dikarenakan memiliki susunan yang unik dan juga adanya proses *encapsulated packet data*. Oleh karena itu untuk membaca data dari *raw data* dilakukan disebuah aplikasi yang bernama Weka, dengan tujuan untuk mendapat nilai-nilai dari semua atribut serta hasil *features extraction* diubah dari berekstensi *.arff* ke dalam bentuk file *.xls* yang akan digunakan untuk proses training karena dengan tipe file ini akan lebih mudah dalam proses *training*.

Pada gambar V-

1 menunjukkan tampilan *raw data* saat masih berekstensi *.arff* dan gambar V-2 menunjukkan tampilan data trafik berekstensi *.xls*.

Gambar V-2. Tampilan data trafik berekstensi .xls

5.4 Hasil Feature Selection

Pada bagian ini adalah pekerjaan utama yang dilakukan yaitu menentukan atribut-

atribut yang relevan dari hasil data trafik yang didapat dari situs *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk>) dengan menggunakan metode *feature ranking* yaitu *information Gain* (IG) untuk mencari ranking yang akan divalidasi dengan menggunakan metode klasifikasi *Naive Bayes*. Hasil dari klasifikasi akan dianalisis berdasarkan tingkat akurasi dari setiap atribut.

Bagian pertama yang dilakukan adalah *pre-processing* data yang didapat. Pada tahap ini adalah proses pembersihan data yang digunakan untuk menghilangkan data-data yang *error* sehingga didapat hasil akurasi yang baik. Tahap selanjutnya adalah membagi dataset menjadi dua bagian yaitu *training* sebesar 70% dan *testing* sebesar 30%, hal ini digunakan untuk proses pembelajaran untuk mendapat prediksi atribut-atribut yang kuat dalam menentukan pola sebuah paket data maupun sebuah serangan. Hasil ranking atribut akan dijelaskan pada tabel V-1.

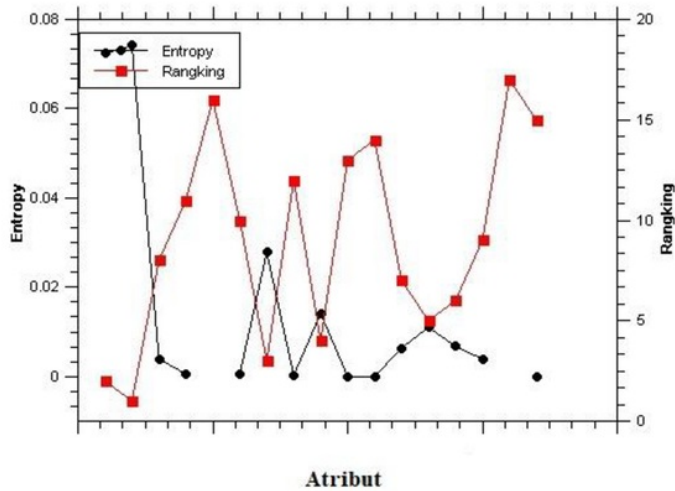
Tabel V-1. Hasil Ranking Atribut dengan menggunakan IG

No	Atribut	Nomor Data	IG	
			Entropy	Ranking
1	4	1	0,869	1

2	35	24	0,496	2
3	180	167	0,470	3
4	137	125	0,387	4
5	145	133	0,361	5
6	149	137	0,329	6
7	73	61	0,316	7
8	75	63	0,293	8
9	65	53	0,283	9
10	67	55	0,283	10
11	163	151	0,280	11
12	51	39	0,275	12
13	157	145	0,246	13
14	155	143	0,245	14
15	161	149	0,244	15
16	72	60	0,241	16
17	159	147	0,237	17
18	147	135	0,231	18

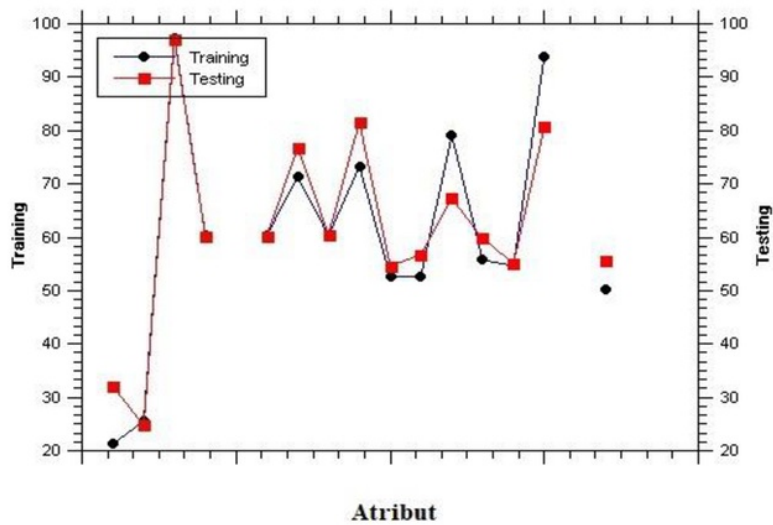
Tabel V-

1 menampilkan hasil *entropy* dan ranking dari setiap atribut dengan menggunakan metode *information Gain* (IG). Hasil ranking didapatkan bahwa atribut 4 merupakan ranking pertama ini artinya atribut ini mempunyai kontribusi yang besar terhadap atribut yang dijadikan *class*, sedangkan tiga atribut yaitu 72, 159, dan 147 merupakan tiga ranking terendah, ini artinya atribut-atribut ini tidak mempunyai kontribusi.



Gambar V-3. Hasil Rangking dan Nilai *Entropy* Atribut

Tahap selanjutnya adalah melakukan proses validasi atribut dengan menggunakan metode klasifikasi *Naïve Bayes*. Proses pengujian ini dengan cara membagi *raw data* menjadi 2 bagian yaitu : 70 % untuk data Training dan 30 % untuk data Testing. Total data yang digunakan sebanyak 19384, maka untuk data *training* didapatkan sebanyak 13568 sedang data untuk testing sebanyak 5816. Tabel V-2 merupakan hasil pengujian *training* dan *testing* akan diuji untuk mendapatkan nilai validasi setiap atribut yang nantinya akan digunakan sebagai pengetahuan baru (*knowledge*) pada proses selanjutnya.



Gambar V-4. Hasil Validasi Atribut dengan *Naive Bayes*

Tabel V-2. Hasil *Training* dan *Testing*

No	Atribut	IG	
		Training	Testing
1	147	0,0000	0,0000
2	161	21.3208	32.0856
3	155	25.4973	24.7892
4	4	97.2942	97.0525
5	73	60.4554	60.265
6	159	0,0000	0,0000
7	75	60.4554	60.265
8	137	71.2501	76.8213
9	149	60.3045	60.5209

10	180	73.2959	81.5143
11	51	52.6728	54.6164
12	67	52.6728	56.7231
13	145	79.0799	67.2691
14	65	55.7368	59.7512
15	163	54.662	55.1056
16	35	93.8767	80.6785
17	72	0,0000	0,0000
18	157	50.2074	55.7061

5.5 Hasil Klasifikasi Dengan Naive Bayes

Tabel V-

3 menunjukkan hasil klasifikasi pada metode Naive Bayes. Proses pengujian dilakukan sebanyak 10 kali yang kemudian didapatkan nilai akurasi, precision, dan recall.

5.6 Analisis Hasil Penelitian

Tabel V-

3 menunjukkan hasil dari klasifikasi *Naive Bayes*, terlihat bahwa *class* yang nilai atributnya banyak muncul adalah WWW. Sedangkan *class* DATABASE, INTERACTIVE, dan GAMES nilai atributnya yang banyak muncul cenderung sedikit dan bahkan tidak ada. Maka dapat dipastikan bahwa pada data trafik yang dipakai oleh penulis dua dari ketiga *class* tersebut tidak ada aktivitas yang terjadi di trafik yang dalam hal ini adalah *Computer Laboratory University of Cambridge* (<http://www.cl.cam.ac.uk>). Untuk penjelasan lebih detail dalam *Confusion Matrix* akan ditampilkan pada gambar V-3.

```
=== Confusion Matrix ===
  a  b  c  d  e  f  g  h  i  j  k  l  <-- classified as
16786 50  0 17  2 24  0  0 13  0  0  0 | a = WWW
1400 170 24  2  0  1  0  0 21  0  0  0 | b = MAIL
 37  1 10  0  0  0  0  0  0  0  0  0 | c = FTP-CONTROL
 1  0  0 108  0  0  0  0  0  0  0  0 | d = FTP-PASV
131  2  0  0  0  0  0  0  1  0  0  0 | e = ATTACK
 8  0  0 71  0 12  0  0  3  0  0  0 | f = P2P
 0  0  0  0  0  0  0  0  0  0  0  0 | g = DATABASE
 0  0  0 342  0  0  0 21  1  0  0  0 | h = FTP-DATA
 0  0  0  30  0  2  0  0 10  0  0  0 | i = MULTIMEDIA
81  1  0  0  0  0  0  0  0  0  0  0 | j = SERVICES
 1  0  0  0  0  0  0  0  0  0  0  0 | k = INTERACTIVE
 0  0  0  0  0  0  0  0  0  0  0  0 | l = GAMES
```

Gambar V-5. Gambar *Confusion Matrix* Hasil Penelitian

5.7 Kesimpulan

Dari hasil pengujian yang telah didapatkan dan dilakukan tahapan analisis hasil klasifikasi, telah jelas diuraikan hasil masing-masing pengujian dari proses klasifikasi dan pemilihan atribut dengan *Naive Bayes* dan *Feature Selection*. Disimpulkan, banyaknya data sangat berpengaruh dalam besarnya nilai akurasi pengklasifikasian, semakin banyak data maka tingkat akurasi akan lebih baik saat pengujian.

KESIMPULAN DAN SARAN

6.1 Pendahuluan

Pada bab ini penulis memberikan kesimpulan dan saran yang diharapkan dapat menjadi acuan bagi penelitian selanjutnya di bidang yang sama.

6.2 Kesimpulan

Berdasarkan analisis dan pembahasan maka dapat ditarik beberapa kesimpulan yaitu sebagai berikut:

1. Di dalam dataset trafik internet yang digunakan oleh penulis nilai atribut yang banyak muncul pada *class* SERVICES, DATABASE, INTERACTIVE, dan GAMES cenderung sedikit dan bahkan tidak ada.
2. Hasil *feature selection* dengan menggunakan metode *feature ranking* didapatkan nilai entropy tertinggi adalah 0.86914856 untuk atribut 4. Untuk validasi atribut nilai tertinggi pada attribute 35 sebesar 97.2942 %.
3. Hasil klasifikasi *Naive Bayes* berdasarkan *class* dihasilkan nilai akurasi, *precision* dan *recall* tertinggi pada *class* WWW (0.994, 0.91, 0.994).

6.3 Saran

Saran yang diajukan oleh penulis, adalah:

1. Pada riset selanjutnya diharapkan dapat melakukan proses visualisasi pemrosesan data secara real-time agar lebih banyak lagi fitur-fitur yang bisa dideteksi.
2. Dapat menemukan solusi untuk meningkatkan kualitas klasifikasi dengan *Naive Bayes* dan *Feature Selection*.
3. Melakukan percobaan dengan metode klasifikasi dan teknik-teknik lain untuk data trafik internet.

Klasifikasi Traffic Network Dengan Menggunakan Naive Bayes dan Feature Selection Dalam Pemilihan Atribut

ORIGINALITY REPORT

15%

SIMILARITY INDEX

8%

INTERNET SOURCES

0%

PUBLICATIONS

16%

STUDENT PAPERS

PRIMARY SOURCES

1

Submitted to Sriwijaya University

Student Paper

14%

2

chemnitzer.linux-tage.de

Internet Source

1%

Exclude quotes On

Exclude matches < 1%

Exclude bibliography On