

## **Integrasi Algoritma Genetika Dan Information Gaint Untuk Menganalisis Sentimen Review Hotel Menggunakan Algoritma Naive Bayes**

**Ari Abdilah<sup>1</sup>, Elva Mardiyani<sup>2</sup>, Mahmud Safudin<sup>3</sup>**

<sup>1</sup>Program Studi Manajemen Informatika  
Akademi Manajemen Informatika dan Komputer Bina Sarana Informatika (AMIK BSI) Pontianak  
Jl. Abdurrahman Saleh No.18A, Pontianak – Kalimantan Barat 78124  
e-mail: ari.aab@bsi.ac.id

<sup>2</sup>Program Studi Manajemen Informatika  
Akademi Manajemen Informatika dan Komputer Bina Sarana Informatika (AMIK BSI) Pontianak  
Jl. Abdurrahman Saleh No.18A, Pontianak – Kalimantan Barat 78124  
e-mail: [elva.eim@bsi.ac.id](mailto:elva.eim@bsi.ac.id)

<sup>3</sup>Program Studi Manajemen Informatika  
Akademi Manajemen Informatika dan Komputer Bina Sarana Informatika (AMIK BSI) Pontianak  
Jl. Abdurrahman Saleh No.18A, Pontianak – Kalimantan Barat 78124  
e-mail: [mahmud.mud@bsi.ac.id](mailto:mahmud.mud@bsi.ac.id)

**Abstract** – *Input and advice is one important part of the application site, in order to assess and improve a quality and quality, Reading reviews helps consumers choose the best hotels, help companies and developers to monitor user satisfaction to improve the quality and quantity of features and services, read as a whole and in manual can spend quite a long time, if read at a glance, the information is not delivered perfectly. This study analyzes the user sentiment Agoda Hotels by automatically classifying reviews for a positive or negative opinion. To improve the accuracy of Naive Bayes methods Feature Selection, Information Gain and genetic algorithms. This model was evaluated using 10 Fold Cross Validation. Measurements were made with the Confusion Matrix and the ROC curve, comparing accuracy before and after the addition of feature selection methods. The results showed an increase in accuracy, 60.50% to 83.00%.*

*Keywords: Sentiment Analysis, Review, Text Classification, Naive Bayes*

### **I. PENDAHULUAN**

Pada saat ini situs microblogging telah menjadi alat komunikasi yang sangat populer di kalangan pengguna internet. Dimana jutaan pesan yang muncul setiap hari di situs web populer yang menyediakan layanan microblogging seperti Twitter, Tumblr, dan Facebook (Alexa, 2013).

Sentiment analysis atau opinion mining adalah studi komputasional dari opini- opini orang, sentimen dan emosi melalui entitas dan atribut yang dimiliki yang diekspresikan dalam bentuk teks (Liu, 2012). Analisis sentiment akan mengelompokkan polaritas dari teks yang ada dalam kalimat atau dokumen untuk mengetahui pendapat yang dikemukakan dalam kalimat atau dokumen tersebut apakah bersifat positif, negatif atau netral (Pang & Lee, 2008). Text mining adalah salah satu teknik yang dapat digunakan untuk melakukan klasifikasi dokumen dimana text mining merupakan variasi dari data mining yang berusaha menemukan pola yang menarik dari sekumpulan data tekstual yang berjumlah besar (Feldman & Sanger, 2007).

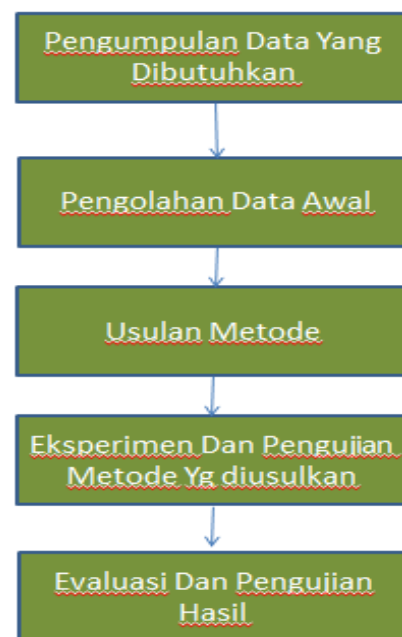
Dunia internet saat ini semakin maju dan semakin banyak diminati dikarenakan media ataupun pasilitas yang dibutuhkan sangatlah mudah untuk didapatkan baik itu melalui *computer*, laptop, maupun smartphone. Seiring dengan perkembangan ini para pengunjung website/konsumen yang menulis opini online terus meningkat. Memilah dan mendata review tersebut secara detail dan keseluruhan dapat mengurangi efisiensi waktu, sehingga dapat menimbulkan ketidak efektifan waktu kerja. Namun, apabila dengan sedikit review yang dibaca evaluasi akan bias. Klasifikasi sentimen bertujuan untuk mengatasi masalah ini dengan secara otomatis mengelompokkan review pengguna menjadi opini positif atau negatif (Z. Zhang, 2011). Terdapat beberapa penelitian yang sudah dilakukan dalam melakukan klasifikasi sentimen terhadap review yang tersedia secara online yaitu, Analisa sentimen pada opini review film menggunakan pengklasifikasi Support Vector Machine dan Particle Swarm Optimization (Basari, 2013). Analisa sentiment pada komentar review film yang ada di jejaring sosial Digg menggunakan pengklasifikasi Naive Bayes, Decision Tree, Maximum-Entropy, dan K-Means

clustering (Yessenov, 2009). Analisa sentimen pada reuters dan teks bahasa China menggunakan pengklasifikasi Naïve Bayes dan dua metric pengevaluasi fitur yaitu Multi-class Odds Ratio (MOR) dan Class Discriminating Measure (CDM) (Chen, 2009). Pengklasifikasian sentimen pada review restoran di internet yang ditulis dalam bahasa Canton menggunakan pengklasifikasi Naïve Bayes dan Support Vector Machine (Z. Zhang, 2011). Klasifikasi sentimen pada review online tempat tujuan perjalanan menggunakan pengklasifikasi Naïve Bayes, Support Vector Machine, dan Character Based N-gram Model (Ye, Zhang, & Law, 2009). Pengklasifikasi Naïve Bayes sangat sederhana dan efisien, (Chen, 2009). Di samping kesederhanaannya, pengklasifikasi Naïve Bayes adalah teknik machine learning yang populer untuk klasifikasi teks, dan memiliki performa yang baik pada banyak domain (Ye, Zhang, & Law, 2009). Namun, Naïve Bayes memiliki kekurangan yaitu sangat sensitif dalam pemilihan fitur (Chen et al., 2009). Tingkatan lain yang umumnya ditemukan dalam pendekatan klasifikasi sentimen adalah pemilihan fitur- fitur yang ada. Pemilihan fitur bisa membuat pengklasifikasi baik lebih efisien/efektif dengan meminimalisasi banyaknya data yang dianalisa, maupun mengidentifikasi fitur yang sesuai untuk dipertimbangkan dalam proses pembelajaran (Moraes, 2013). Metode filter terdiri dari document frequency, mutual information, information gain, dan chi-square. Tidak ada dari keempat metode tersebut yang secara luas diterima sebagai metode penyeleksi fitur terbaik untuk klasifikasi sentimen atau kategorisasi teks, namun, information gain sering lebih unggul dibandingkan yang lain (Moraes, 2013). Menurut Kohavi dalam Yang (Yang, 2010) wrapper mengevaluasi fitur secara berulang dan menghasilkan akurasi klasifikasi yang tinggi. Menurut Gunal (Gunal, 2012) salah satu metode wrapper yang bisa digunakan dalam pemilihan fitur adalah Genetic algorithm (GA). Umumnya metode pemilihan fitur yang lebih disukai adalah filter dikarenakan waktu pemrosesannya yang relatif rendah. Information Gain mengukur berapa banyak informasi kehadiran dan ketidakhadiran dari suatu kata yang berperan untuk membuat keputusan klasifikasi yang benar dalam class apapun. Information Gain adalah salah satu pendekatan filter yang sukses dalam pengklasifikasian teks (Uysal & Gunal, 2012). Untuk mengurangi kerumitan perhitungan dilakukan pemilihan fitur dengan menghitung Information Gain (Z. Zhang, 2011). Pada penelitian ini penulis mencoba pengklasifikasi Naïve Bayes dengan Information Gain dan Genetic algorithm sebagai metode yang akan diterapkan untuk menganalisa dan mengklasifikasikan teks/data *review* pada komentar dari sebuah hotel ternama untuk meningkatkan akurasi penganalisaan.

## II. METODOLOGI PENELITIAN

Ketersediaan data akan sangat menentukan dalam proses pengolahan dan analisa selanjutnya, karenanya, dalam pengumpulan data harus dilakukan teknik yang menjamin bahwa data diperoleh itu benar, akurat dan bisa dipertanggung jawabkan sehingga hasil pengolahan dan analisa data tidak bias.

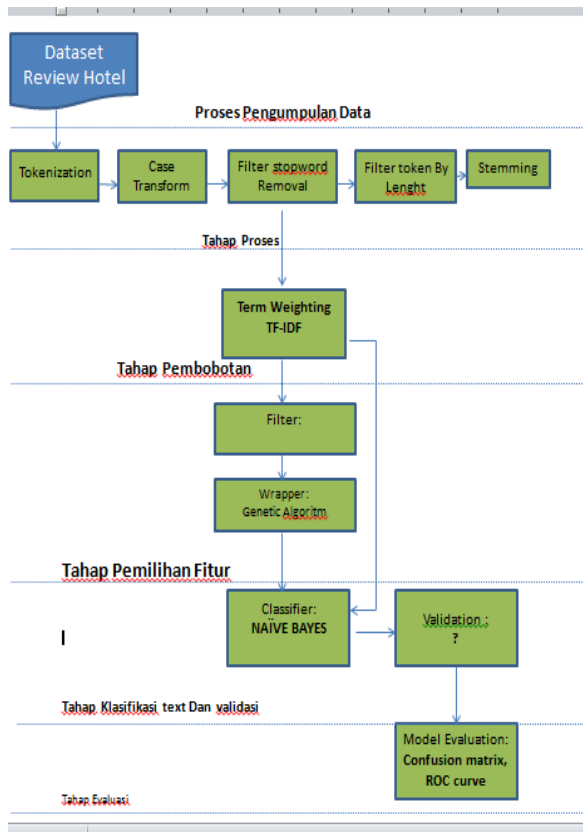
Metode yang digunakan yaitu metode penelitian eksperimen dengan tahapan penelitian yaitu:



Sumber: (Gorunescu, 2011)

Gambar 1. Tahapan Penelitian

Dalam penelitian ini dilakukan penggabungan dua jenis metode pemilihan fitur, yaitu Information Gaint, dan Genetic Algorithm dari jenis wrapper sebagai metode pemilihan fitur untuk meningkatkan akurasi dari hasil analisa sentimen ulasan review Pada Hotel AGODA yang menggunakan pendekatan algoritma Naïve Bayes. Data ulasan yang telah dikumpulkan melalui tahap preprocessing terlebih dahulu agar didapatkan kata-kata relevan untuk diklasifikasi. Proses evaluasi dilakukan menggunakan 10 fold cross validation. Pengukuran akurasi diukur dengan confusion matrix. Hasil yang nantinya akan dibandingkan adalah akurasi Naïve Bayes sebelum dan setelah menggunakan metode pemilihan fitur gabungan, yaitu Information Gaint dan Genetic algorithm. Dimana pada Genetic algorithm, Naïve Bayes diuji di dalam tahap wrapper. Berikut langkah penelitian yang dilakukan:



Sumber: (Gorunescu, 2011)  
Gambar 2. Langkah Penelitian

### III. HASIL DAN PEMBAHASAN

#### 1. Tahap Pengumpulan Data

Data penelitian berupa 200 ulasan positif dan 200 ulasan negatif berbahasa Inggris dari para review yang ada pada website Hotel AGODA, dimana 100 ulasan positif dan 100 ulasan negatif akan digunakan pada tahap training dan data testing pada model klasifikasi, dan sebanyak 100 ulasan positif dan 100 ulasan negatif akan diujikan pada aplikasi yang dirancang untuk implementasi. Data review dikumpulkan dalam bentuk notepad.

#### 2. Tahap Pengolahan Awal Data (Tahap Preprocessing)

Yang termasuk dalam tahapan ini meliputi:

##### a. Tokenization

Semua kata yang ada di dalam tiap dokumen ulasan dikumpulkan lalu tanda baca, dan simbol apapun yang bukan huruf dihilangkan.

Tabel 1. Perbandingan Teks Sebelum dan Sesudah Tokenization

Sebelum Tokenization	Sesudah Tokenization
look less professional hotel services , hotel room condition dirty , moist , slightly smelly and worse no hot water in the shower . for officers menngkontak	look less professional hotel services hotel room condition dirty moist slightly smelly and worse no hot water in the shower for officers

no means the phone and we had to walk back and forth , the column can not be swimming with the boss would be cleaned and children's games much damaged , I suggest that things are all correct for the convenience of hotel guests trims.	menngkontak no means the phone and we had to walk back and forth the column can not be swimming with the boss would be cleaned and children s games much damaged I suggest that things are all correct for the convenience of hotel guests trims
---	--

##### b. Transform Case

Seluruh huruf kapital (uppercase) didalam dokumen ulasan diubah menjadi huruf kecil (lower case).

Tabel 2. Perbandingan Teks Sebelum dan Sesudah Transform Case

Sebelum Transform Case	Sesudah Transform Case
Choosing Hyatt Regency Bandung is one mistake , I arrived at 4 pm is too late 2 hours from the time of check -in , but it turns out the rooms that I have a new booking can be ready at 5 pm without a clear explanation . After the available room turned out some towels still no , when I asked the housekeeping supply , takes up to 4 hours . food was very standard at all to the level of five-star hotel.	choosing hyatt regency bandung is one mistake i arrived at pm is too late hours from the time of check in but it turns out the rooms that i have a new booking can be ready at pm without a clear explanation after the available room turned out some towels still no when i asked the housekeeping supply takes up to hours food was very standard at all to the level of five star hotel

##### c. Filter Stopwords Removal

Kata-kata yang tidak relevan dalam dokumen ulasan akan dihapus, seperti kata the, of, for, with, dan lain-lain yang merupakan kata-kata yang tidak mempunyai makna tersendiri jika dipisahkan dengan kata yang lain dan tidak terkait dengan kata sifat yang berhubungan dengan sentimen.

Tabel 3. Perbandingan Teks Sebelum dan Sesudah Stopword Removal

Sebelum Transform Case	Sesudah Transform Case
Choosing Hyatt Regency Bandung is one mistake , I arrived at 4 pm is too late 2 hours from the time of check -in , but it turns out the rooms that I have a new booking can be ready at 5 pm without a clear explanation . After the available room turned out some towels still no , when I asked the housekeeping supply , takes up to 4 hours . food was very standard at	choosing hyatt regency bandung mistake i arrived pm late hours time check turns rooms i booking ready pm clear explanation available room turned towels i asked housekeeping supply takes hours food standard level star hotel

all to the level of five-star hotel.	
--------------------------------------	--

d. Filter Token By Length

Menyaring kembali kata dalam dokumen ulasan yang kurang dari batas minimal karakter dan batas maksimal karakter, kemudian menghapusnya. Pada penelitian ini batas minimal yang digunakan sebanyak 2 karakter, dan batas maksimal sebanyak 20 karakter.

Tabel 4. Perbandingan Teks Sebelum dan Sesudah Filter Tokens By Length

Sebelum Transform Case	Sesudah Transform Case
choosing hyatt regency bandung mistake i arrived pm late hours time check turns rooms i booking ready pm clear explanation available room turned towels i asked housekeeping supply takes hours food standard level star hotel.	choosing hyatt regency bandung mistake arrived late hours time check turns rooms i booking ready clear explanation available room turned towels asked housekeeping supply takes hours food standard level star hotel

e. Stemming

Mencari root kata dari tiap kata dalam dokumen ulasan, memecahkan atau memotong setiap varian-varian kata untuk menemukan kata dasarnya agar dapat mengelompokkan bentuk-bentuk yang berasal dari kata dasar yang sama, contohnya seperti write, wrote dan written di mana kata dasar dari semuanya adalah kata write.

Tabel 5. Perbandingan Teks Sebelum dan Sesudah Stemming

Sebelum Transform Case	Sesudah Transform Case
choosing hyatt regency bandung mistake arrived pm late hours time check turns rooms i booking ready lear explanation available room turned towels asked housekeeping supply takes hours food standard level star hotel.	choos hyatt regenc bandung mistak arriv pm late hour time check turn room book readi pm clear explan avail room turn towel ask housekeep suppli take hour food standard level star hotel

3. Tahap Transformation (Pembobotan Kata dengan Algoritma TF-IDF)

Pada tahap ini akan dilakukan perhitungan bobot kata yang dihasilkan dari tahap preprocessing menggunakan algoritma TF-IDF.

Gambar 3. Sample Hasil Pembobotan Kata dengan Algoritma TF-IDF

4. Tahap Klasifikasi Teks Sentimen Menggunakan Algoritma Naïve Bayes

Klasifikasi teks dilakukan untuk menentukan apakah sebuah ulasan termasuk sebagai class positif atau class negatif berdasarkan nilai perhitungan probabilitas dari rumus Algoritma Naïve Bayes yang lebih besar.

Kehadiran kata di dalam suatu dokumen ulasan akan diwakili oleh angka 1 dan angka 0 jika kata tersebut tidak muncul di dalam dokumen ulasan.

Tabel 6. Tabel Vector Dokumen Boolean & Label Class Hasil Klasifikasi

Deskripsion	good	bad	nice	dirty	Clean	hot	class
doc_n001	0	1	0	0	0	0	negatif
doc_n002	0	0	1	0	0	0	negatif
doc_n019	0	1	0	0	0	0	negatif
doc_n011	0	1	0	1	0	1	negatif
doc_n038	1	1	0	0	0	0	negatif
doc_p063	1	0	0	0	0	0	positif
doc_p056	1	0	0	0	0	0	positif
doc_p051	1	0	1	0	1	0	positif
doc_p095	1	0	0	0	0	0	positif
doc_p68	1	0	1	0	1	0	?

Berikut adalah perhitungan probabilitas bayes untuk dokumen ulasan “docp68.txt”.

a. Menghitung Probabilitas Bersyarat

Untuk class positif:

$$P(\text{docp68}|\text{positif}) =$$

$$P(\text{good} = 1 | \text{positif}) \times P(\text{bad} = 0 | \text{positif}) \times P(\text{nice} = 1 | \text{positif}) \times P(\text{dirty} = 0 | \text{positif}) \times P(\text{clean} = 1 | \text{positif}) \times P(\text{hot} = 0 | \text{positif})$$

$$= 4/4 \times 4/4 \times 1/4 \times 4/4 \times 1/4 \times 4/4$$

$$= 1 \times 1 \times 0,25 \times 1 \times 0,25 \times 1$$

$$= \underline{0,0625}$$

Untuk class negatif:

$$P(\text{docp68}|\text{negatif}) =$$

$$P(\text{good} = 1 | \text{negatif}) \times P(\text{bad} = 0 | \text{negatif}) \times P(\text{nice} = 1 | \text{negatif}) \times P(\text{dirty} = 0 | \text{negatif}) \times P(\text{clean} = 1 | \text{negatif}) \times P(\text{hot} = 0 | \text{negatif})$$

$$= 1/5 \times 1/5 \times 1/5 \times 4/5 \times 0/5 \times 4/4$$

$$= 0,2 \times 0,2 \times 0,5 \times 0,8 \times 0 \times 1$$

$$= \underline{0}$$

**b. Menghitung Probabilitas Prior**

Perhitungan probabilitas prior dari class positif dan negatif dihitung dengan proporsi dokumen pada tiap class:

$$P(\text{positif}) = 4/9 = \underline{0,44}$$

$$P(\text{negatif}) = 5/9 = \underline{0,55}$$

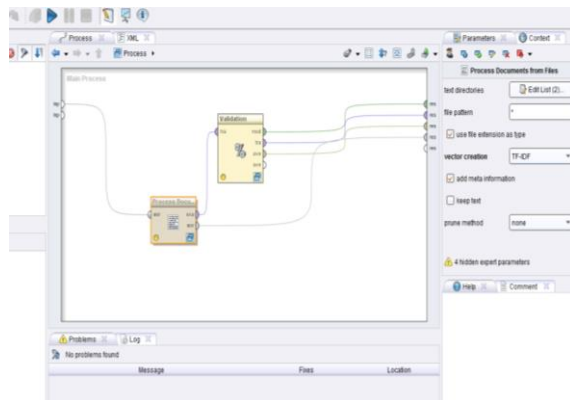
**c. Menghitung Probabilitas Posterior**

Perhitungan probabilitas posterior dengan memasukkan rumus Bayes dan menghilangkan penyebut P(docp68):

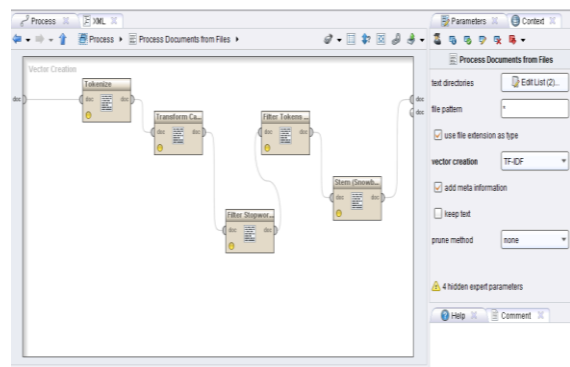
$$P(\text{doc\_p068}|\text{positif}) = \frac{(0,0625)(0,44)}{P(\text{doc\_p068})} = \underline{0,0275}$$

$$P(\text{doc\_p068}|\text{negatif}) = \frac{(0)(0,55)}{P(\text{doc\_p068})} = \underline{0}$$

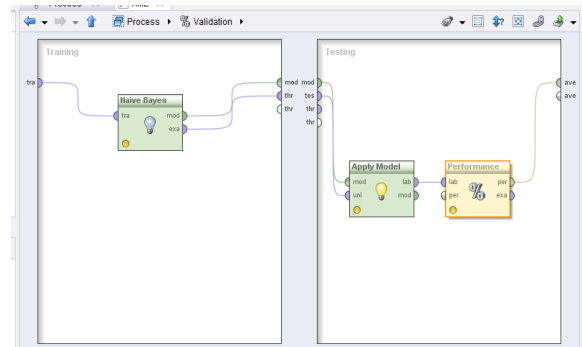
Berdasarkan probabilitas di atas, maka dapat disimpulkan bahwa kata doc\_p068 termasuk dalam class positif, karena  $P(\text{positif}|\text{doc\_p068}) > P(\text{negatif}|\text{doc\_p068})$ . Perhitungan di atas dapat dibuat suatu model dengan menggunakan RapidMiner 5. Desain model dapat dilihat pada berikut:



Gambar 4. Desain Model Klasifikasi Naïve Bayes Menggunakan RapidMiner 5.2.



Gambar 5. Tahap Preprocessing



Gambar 6. Tahap Validation dengan 10 Fold Cross Validation

**5. Tahap Validasi/Pengujian Model dengan 10 Fold Cross Validation**

Pada penelitian ini, penulis melakukan pengujian model dengan menggunakan teknik 10 *cross validation*, di mana proses ini membagi data secara acak ke dalam 10 bagian. Proses pengujian dimulai dengan pembentukan model dengan data pada bagian pertama. Model yang terbentuk akan diujikan pada 9 bagian data sisanya. Setelah itu proses akurasi dihitung dengan melihat seberapa banyak data yang sudah terklasifikasi dengan benar.

**6. Tahap Optimasi Model**

Dengan menggabungkan metode pemilihan fitur filter dan wrapper, di mana dalam penelitian ini metode yang digunakan adalah **Information gain** dari filter dan **Genetic algorithm** dari wrapper. Data yang akan diolah diberikan bobot dari Information gain untuk meningkatkan akurasi pengklasifikasi Naïve Bayes. Penulis menggunakan operator *select by weight* dengan memilih parameter *weight relation*= top k, dengan k= 10. Di mana nanti akan dihasilkan 10 atribut teratas. 10 atribut yang terpilih akan ditampilkan bobotnya masing-masing, untuk lebih jelasnya dapat dilihat pada tabel 7.

Tabel 7. Daftar 10 Atribut/Fitur Teratas

attribute	
bad	1
braga	1
clean	1
enjoy	1
famili	1
friend	1
good	1
locat	1
star	1
walk	1

Bobot di atas adalah bobot yang sudah di-generate oleh operator *select by weight*. atribut yang



ditampilkan bobotnya dari masing-masing dokumen hanya yang mempunyai bobot 1. Di antara 10 atribut di atas, yang berhubungan dengan sentimen hanyalah kata *bad*, *good* dan *clean*. Tabel 8 menunjukkan 3 atribut tersebut di dalam dokumen bentuk vektornya. Berikut adalah perhitungan bobot Information gainnya.

Tabel 8. Tabel Vector Dokumen Boolean & Label Class Hasil Klasifikasi Setelah Diberikan Bobot Dari Information gain

NO	Dokumen	bad	good	clean	class
1	doc_n001	1	0	0	negatif
2	doc_n002	0	0	0	negatif
3	doc_n019	1	0	0	negatif
4	doc_n011	1	0	0	negatif
5	doc_n038	1	1	0	negatif
6	doc_p063	0	1	0	positif
7	doc_p056	0	1	0	positif
8	doc_p051	0	1	1	positif
9	doc_p095	0	1	0	positif
10	doc_p068	0	1	1	positif

### 7. Eksperimen Terhadap Indikator Model

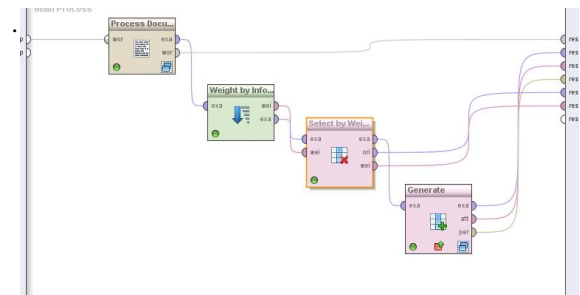
Untuk mendapatkan model yang baik, beberapa indikator disesuaikan nilainya agar didapatkan hasil akurasi yang tinggi. Dalam hal ini, pembobotan menggunakan Information gain akan ditampilkan hanya 10 kata yang sering muncul. Untuk Genetic algorithm, indikator yang disesuaikan adalah *population size=50*, *p initialize=0.8*, *p crossover*, dan *p generate=1.0*. Sedangkan yang diuji coba untuk meningkatkan akurasi adalah nilai *p crossover*. Tabel indikator dan hasil pengujian dapat dilihat pada tabel 9.

Tabel 9. Tabel Indikator dan Hasil Pengujian

nilai P crossover	Hasil Akurasi	Execution time
0.6	82.50%	1:15 detik
0.7	83.00%	1:18 detik
0.8	83.00%	1:26 detik
0.9	83.00%	2:55 detik
1.0	82.50%	1:43 detik

Akurasi yang paling tinggi dengan kombinasi *population size=50*, *p initialize=0.8*, *p crossover=0.7*, dan *p generate=1.0* dengan hasil akurasi mencapai 83.00%.

Desain model Naïve Bayes dengan metode pemilihan fitur Information gain dan Genetic algorithm ini dapat dilihat pada gambar 9.



Gambar 7. Desain Model Naïve Bayes dan Metode Pemilihan Fitur

Dengan memiliki model klasifikasi teks pada review, pembaca dapat dengan mudah mengidentifikasi mana review yang positif maupun yang negatif.

Dari data review yang sudah ada, dipisahkan menjadi kata-kata, lalu diberikan bobot pada masing-masing kata tersebut. Dapat dilihat kata mana saja yang berhubungan dengan sentimen yang sering muncul dan mempunyai bobot paling tinggi. Dengan demikian dapat diketahui review tersebut positif atau negatif.

Dalam penelitian ini, hasil pengujian model akan dibahas melalui *confusion matrix* untuk menunjukkan seberapa baik model yang terbentuk. Tanpa menggunakan metode pemilihan fitur, algoritma Naïve Bayes sendiri sudah menghasilkan akurasi sebesar **60.50%** dan nilai **AUC 0.519**. Akurasi tersebut masih kurang akurat, sehingga perlu ditingkatkan lagi menggunakan metode pemilihan fitur. Setelah menggunakan metode pemilihan fitur dari filter dan wrapper yang digabungkan, akurasi algoritma Naïve Bayes meningkat menjadi **83.00%** dan nilai **AUC 0.872**.

### 8. Tahap Evaluasi (Pengukuran dengan Confusion Matrix dan ROC Curve/AUC)

Hasil dari pengujian model akan dibahas melalui Confusion Matrix untuk menunjukkan seberapa baik model yang terbentuk.

Tabel 10. Confusion Matrix Model Naïve Bayes Tanpa Penambahan Metode Feature Selection

accuracy: 60.50% +/- 11.06% (mikro: 60.50%)			
	true negatif	true positif	class precision
pred. negatif	62	41	60.19%
pred. positif	38	59	60.82%
class recall	62.00%	59.00%	

Tabel 11. Confusion Matrix Model Naïve Bayes Sesudah Penambahan Metode Feature Selection

accuracy: 83.00% +/- 9.00% (mikro: 83.00%)			
	true negatif	true positif	class precision
pred. negatif	92	26	77.97%
pred. positif	8	74	90.24%
class recall	92.00%	74.00%	

Berikut adalah tabel perbandingan hasil pengujian model Algoritma Naïve Bayes sebelum dan sesudah menggunakan metode Feature Selection Information Gain & Genetic Algorithm:

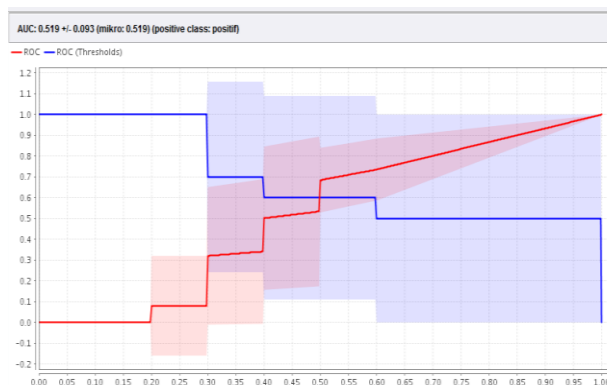
Tabel 12. Perbandingan Model Algoritma Naïve Bayes Sebelum dan Sesudah Penambahan Metode Feature Selection

Perbandingan	Hasil Klasifikasi Ulasan Negatif	Hasil Klasifikasi Ulasan positif	Accuracy Model	AUC
Algoritma Naïve Bayes	62	59	60.50 %	?
Algoritma Naïve Bayes + Information Gain & Genetic Algorithm	92	74	83.00 %	?

Berikut adalah perhitungan akurasinya:

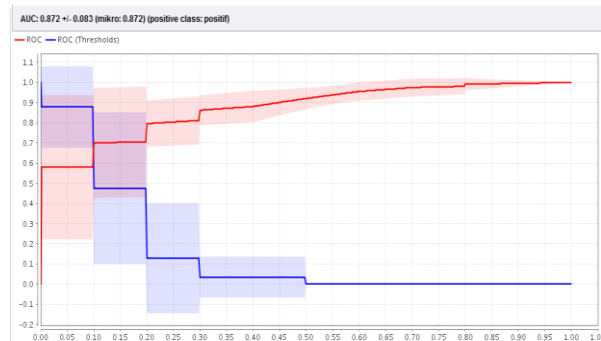
$$\text{Accuracy: } \frac{74 + 92}{74 + 26 + 8 + 92} = \frac{166}{200} = 0.83 \times 100\% = \mathbf{83.00\%}$$

Hasil perhitungan Confusion Matrix diatas digambarkan melalui kurva ROC.



Gambar 8. Kurva ROC Model Naïve Bayes sebelum ditambahkan Metode Feature Selection

Berdasarkan Kurva ROC dari model Naïve Bayes sebelum penambahan metode Feature Selection diatas nilai AUC (Area Under Curve) yang didapatkan yaitu sebesar **0.519** dan termasuk ke dalam kategori *Fair Classification*.



Gambar 9. Kurva ROC Model Naïve Bayes Sesudah ditambahkan Metode Feature Selection

Sedangkan berdasarkan Kurva ROC dari model Naïve Bayes setelah penambahan metode Feature Selection diatas, dapat dilihat bahwa nilai AUC (Area Under Curve) yang didapatkan lebih besar yaitu **0.872** dan termasuk ke dalam kategori *Good Classification*.

#### IV. KESIMPULAN

Untuk mengklasifikasikan teks dengan data berupa review Hotel, salah satu pengklasifikasi yang dapat digunakan adalah pengklasifikasi Naïve Bayes. Hal ini dikarenakan Naïve Bayes sangat sederhana dan efisien. Selain itu Naïve Bayes juga sangat populer digunakan untuk klasifikasi teks dan memiliki performa yang baik.

Berdasarkan hasil dari pengujian model yang dilakukan, didapatkan peningkatan akurasi sebesar **23%**, dan dapat dikatakan bahwa Naïve Bayes yang integrasikan dengan kedua metode seleksi Information Gaint dan Genetic Algorithm memang merupakan metode yang cukup baik dalam mengklasifikasikan teks khususnya pada kasus analisa sentimen seperti pada penelitian ini.

Hasil dari penelitian ini dapat membantu pihak perusahaan, developer, dan juga pengguna aplikasi dalam menganalisis sentimen rievew konsumen dari situs website khususnya **Hotel AGODA** dengan secara otomatis mengklasifikasikan ulasan berbahasa Inggris kedalam dua kategori yaitu positif dan negatif secara otomatis dengan waktu yang singkat.

#### REFERENSI

- Basari, A. S. H., Hussin, B., Ananta, I. G. P., & Zeniarja, J. (2013). Opinion Mining of Movie Review using Hybrid Method of Support Vector Machine and Particle Swarm Optimization. *Procedia Engineering*, 53, 453-462.
- Chen, J., Huang, H., Tian, S., & Qu, Y. (2009). Feature selection for text classification with Naïve Bayes. *Expert Systems with Applications*, 36(3), 5432-5435.
- Feldman, R. (2013). Techniques and applications for sentiment analysis. *Communications of the ACM*, 56(4), 82.
- Gorunescu, Florin. (2011). *Data Mining: Concepts*

- and Techniques. Verlag berlin Heidelberg: Springer
- Gunal, S. (2012). Hybrid feature selection for text classification, 20.
- Han, J., & Kamber, M. (2007). *Data Mining Concepts and Techniques*.
- Moraes, R., Valiati, J. F., & Gavião Neto, W. P. (2013). Document-level sentiment classification: An empirical comparison between SVM and ANN. *Expert Systems with Applications*, 40(2), 621–633.
- Uysal, A. K., & Gunal, S. (2012). A novel probabilistic feature selection method for text classification. *Knowledge-Based Systems*, 36, 226–235.
- Witten, H. I., Frank, E., & Hall, M. A. (2011). *Data Mining Practical Machine Learning Tools and Technique*. Burlington: Elsevier Inc.
- Yessenov, K. (2009). Sentiment Analysis of Movie Review Comments 6.863, 1–17. Yessenov, K. (2009). Sentiment Analysis of Movie Review Comments 6.863, 1–17.
- Ye, Q., Zhang, Z., & Law, R. (2009). *Expert Systems with Applications Sentiment classification of online reviews to travel destinations by supervised machine*
- Zhang, Z., Ye, Q., Zhang, Z., & Li, Y. (2011). Sentiment classification of Internet restaurant reviews written in Cantonese. *Expert Systems with Applications*, 38(6), 7674–7682.