



ANALISIS SENTIMEN TERHADAP WACANA POLITIK PADA MEDIA MASA ONLINE MENGGUNAKAN ALGORITMA SUPPORT VECTOR MACHINE DAN NAIVE BAYES

Andi Nurul Hidayat
Fakultas Ilmu Komputer, Stmik Bina Mulia
Email: andinurulhidayat@yahoo.co.id

ABSTRAK

Analisis sentimen merupakan salah satu domain Text Mining ataupenggalian data berupateks, yang di antaranya terdapat proses mengolah dan mengekstrak data teks tual secara otomatis untuk mendapatkan informasi. Manfaat analisis sentimen dalam dunia politik antara lain

untuk membantu dalam menganalisis kejadian publik pemerintah serta memberikan efisiensi waktu dan efisiensi kerja bagi para penyedia berita dalam mengklasifikasi klasikan berita dan membantu para pencari berita untuk mendapatkan wacana berita politik harian yang mereka inginkan. Proses pada analisis sentimen ini awal dengan preprocessing, dilanjutkan dengan pembobotan kata, kemudian penghitungan cosine similarity, dan klasifikasi. Preprocessing terdiri dari beberapa tahap yaitu cleansing, tokenizing, stopword removal, dan stemming. Metode pembobotan kata yang digunakan adalah Term Frequency Inverse Document Frequency (TF-IDF) dan menggunakan Support Vector Machine (SVM). Naive Bayes Classifier (NBC), sebagai metode klasifikasinya. Adalah suatu metode pengklasifikasian berdasarkan mayoritas dari polarity document subjectivity yang di hasilkan dari crawling. Metode ini bertujuan untuk mengklasifikasi objek baru berdasarkan atribut dan training sample. Pengujian akurasi dari Analisis Sentimen Terhadap Wacana Politik Pada Media Masa Online Berbahasa Inggris dengan metode NB menghasilkan rata-rata akurasi sebesar 59,98 % dan nilai tertinggi akurasi sebesar SVM 90,50%.

Kata Kunci: Analisis Sentimen Wacana Politik Pada Media Masa Online, Text Mining, SVM (Support Vector Machine), NBC (Naive Bayes Classifier).

1. Pendahuluan

Analisis sentimen atau *opinion mining* mulai populer pada tahun 2002 mempublikasikan ide di balik penelitiannya yang di lakukan analisis sentimen adalah proses menyajikan informasi dengan membangun sebuah sistem yang dapat mengklasifikasikan dokumen teks ke dalam dua kategori, yakni nilai positif dan negatif yang sesuai dengan keseluruhan sentimen yang dinyatakan di dalam setiap dokumen tersebut. Dalam sebuah metode klasifikasi dalam analisis sentimen menggunakan metode-metode klasifikasi yang biasa digunakan untuk kategorisasi teks antara lain metode *supervised learning machine* (SVM) Based Approach maupun Maximum Entropy.

Website adalah tempat yang baik bagi orang-orang untuk mengekspresikan pendapat mereka pada berbagai topik salah satunya adalah memanfaatkan situs jejaring sosial misalnya facebook, twitter, bahkan pemberi opini secara profesional, seperti reviewer berita politik dan film, pemilik blog dimana publik dapat mengomentari dan merespon apa yang mereka pikirkan. Kemampuan untuk merangkak dari website serta mengekstrak pendapat dari baris-baris teks dapat menjadi sangat berguna bidang ini adalah area studi yang banyak dikaji karena kemungkinan nilai komersialnya.

Kebanyakan informasi disimpan sebagai teks, sehingga *text mining* memiliki potensi nilai komersial. *Sentiment analysis* atau *opinion mining*

adalah studi komputisional dari opini-opini orang, sentimen dan emosi melalui entitas dan atribut yang dimiliki yang diekspresikan dalam bentuk teks. Analisis sentimen 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.

Text mining adalah proses mengambil informasi berkualitas tinggi dari teks informasi berkualitas tinggi biasanya diperoleh melalui peramalan pola dan kecenderungan melalui analisis seperti pembelajaran pola statistik. Secara umum proses *text mining* dapat meliputi kategorisasi teks *text clustering*, ekstraksi konsep/entitas, produksi taksonomi granular, *sentiment analysis*, penyimpulan dokumen, dan pemodelan relasi entitas. Salah satu metode klasifikasi yang dapat digunakan adalah metode Naive Bayes yang sering disebut dengan *Naive Bayes Classifier* (NBC). Kelebihan NBC adalah sederhana tetapi memiliki akurasi tinggi. Terbukti dapat digunakan secara efektif untuk mengklasifikasi kan berita secara otomatis. Algoritma NBC yang sederhana dan kecepatannya yang tinggi dalam proses pelatihan dan klasifikasi membuat algoritma ini menjadi cukup digunakan sebagaimana salah satu metode klasifikasi. Kedua algoritma klasifikasi tersebut banyak digunakan dalam klasifikasi teks. Padahal eksperimen untuk klasifikasi teks berbahasa Inggris didapatkan bahwa SVM yang menunjukkan performansi sedikit lebih baik dibandingkan metode NBC. Metode NBC adalah metode yang jauh lebih sederhana dan mudah diaplikasikan. Sehingga pada penelitian ini menggunakan metode yang manakah memiliki performansi yang

lebih baik untuk diimplementasikan dalam *sentiment analysis* terhadap wacana politik media massa Indonesia berbahasa Inggris. Dalam pengklasifikasi sentimen analisis daratahun 2004 sampai hari ini dengan menggunakan google trend hal ini dapat diamatibahkan dengan pertumbuhan yang tersedia data di internet, kebutuhan untuk analisis sentimen juga meningkat, berbagai statistik dan linear metode quiistik telah dikembangkan untuk analisis sentimen teknologi wacana politik berbahasa Inggris.

2. Penelitian Terkait

Penelitian mengenai klasifikasi sentimen analisis pada wacana politik salah satu permasalahannya adalah tingkat pengklasifikasian sebuah teks kalimat berbahasa Inggris yang dimana dimensi sebuah teks opini wacana politik ini yang di analisa terdapat ambigu dalam penggunaan kata, tidak adanya intonasi dalam sebuah teks sehingga menyebabkan banyak atribut yang kurang spesifik dan relevan sehingga menurunkan kinerja dan performa klasifikasi teks opini sentimen analisis opini wacana politik. Untuk mendapatkan accuracy yang baik, atribut yang ada harus dipilih dengan algoritma yang tepat. Bagian yang penting untuk mengoptimalkan sebuah klasifikasi dokumen teks adalah menggunakan *Feature Selection* salah satunya, yaitu unigram, unigram + bigram, unigram + Part of Speech (POS), adjective, dan nigram dan dikombinasikan dengan unigram.

Berdasarkan beberapa penelitian sentimen analisis men kombinasikan beberapa algoritma *Feature Selection* dan algoritma untuk mendapatkan hasil yang baik dan *performance* yang baik. Penelitian *sentiment analysis* yang dilakukan oleh Ahmed Abbasi, Hsinchun Chen, & Arab Salem berjudul *Sentiment Analysis in Multiple Languages: Feature Selection for Opinion Classification in Web Forums* di gabungkan metode hybrid sisial algoritma genetika EWGA mendapatkan hasil yang lebih baik. Sedangkan *Bayesian Opinion Mining* dilakukan oleh Ian Barber pada data review film berbahasa Inggris dan diujikan untuk 5000 record opini negative dan 5000 record opini positif sebagai data latih dan 333 record opini negatif sebagai data

ujisertamenghasilkanakurasisebesar 80% menggunakan metode naive bayes classifier sedangkan Klasifikasi Berita Berbahasa Indonesia Menggunakan *Naive Bayes Classifier*, dilakukan oleh yudi Wibisono. Sedangkan penelitian yang berjudul *text mining* dengan metode *naive bayes classifier support vector machine* untuk analisis oleh ni wayan sumartini saraswati dalam penelitiannya menguji data sebanyak 3000 record positif dan negatif menggunakan metode *Naive Bayes Classifier* 80.18% dan *support vector machine* 80.15%. sedangkan pada penelitian yang dilakukan Fatimah Wulandin dan Anto Satriyo Nugroho dengan judul *Text Classification Using Support Vector Machine for Webmining Based Spatio Temporal Analysis of the Spread of Tropical Diseases*, menggunakan 4 metode klasifikasi dan mengkomparasikan dengan metode klasifikasi yang digunakan yakni algoritma SVM, NBC, KNN dan C45 hasilnya pada data 3713 feature dan 360 instance. 360 instance sebagai data latih dan 120 instance. sedangkan pada penelitian Fabrice Colas & Pavel Brazdil yang berjudul *comparison of SVM and Some Older Classification Algorithms in Text Classification Tasks*, melakukan pengujian 3 metode dengan menggunakan algoritma SVM, KKN dan NBC.

Sedangkan pada penelitian Blitzer, J., Dredze, M. & Pereira yang berjudul *Biographies, Bollywood, Boom-boxes and Blenders: Domain Adaptation for Sentiment Classification*, menguji dokumen bahwasannya menggunakan metode Structural Correspondence learning (SCL), b aseline, SCLMI SCLMI menunjukkan performansi yang lebih baik untuk adaptasi domain. hasil eksperimen menunjukkan bahwa menguji data sentimen analisis wacana politik media masa online menggunakan algoritma *naive bayes* hasilnya mendapatkan akurasi sebesar 59,98 % dan membandingkan dengan algoritma *support vector machine* mendapatkan hasil yang lebih baik mendapatkan akurasi 90,50%.

Tabel 1 Hasil Akurasi algoritma klasifikasi

Algoritma	Accuracy
NB	59,98 %
SVM	90,50%

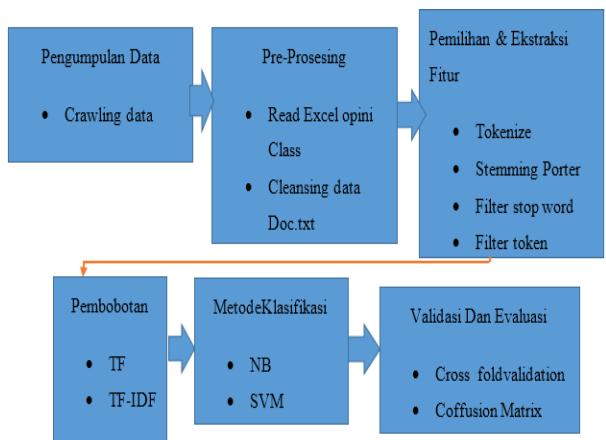
3. METODE YANG DI USULKAN

Peneliti mengusulkan untuk mengkomparasikan 2 algoritma klasifikasi

yakni (naive bayes classifier dan support vector machine) dan menggunakan fitur ekstraksi dokumen (tokenize, stemming porter, filter token dan stop word). Peneliti menguji data tiap-tiap dokumen berita yang di kumpulkan secara online melalui teknik web mining yakni tools crawling data melalui situs portal berita. Data yang sudah di bersihkan yang berisi tag-tag, doc, html, dan untuk menghitung pembobotan teks / term adalah menggunakan *space vector* pembobotan TF-IDF serta menganalisa hasil performa suatu metode dalam klasifikasi dokumen harian politik dan meghasilkan masing-masing tingkat akurasi algoritma *Naive Bayes* dan *support Vector Machine*. Dalam pemilihan fitur ekstrasi dokumen teks wacana harian politik media masa online pada situs portal berita ini. Sebelum melakukan komparasi/kombinasi dataset di lakukan *text processing* terlebih dahulu *text processing* bertujuan untuk mempersiapkan dokumen teks yang tidak terstruktur menjadi data terstruktur yang siap di gunakan untuk proses data selanjutnya.

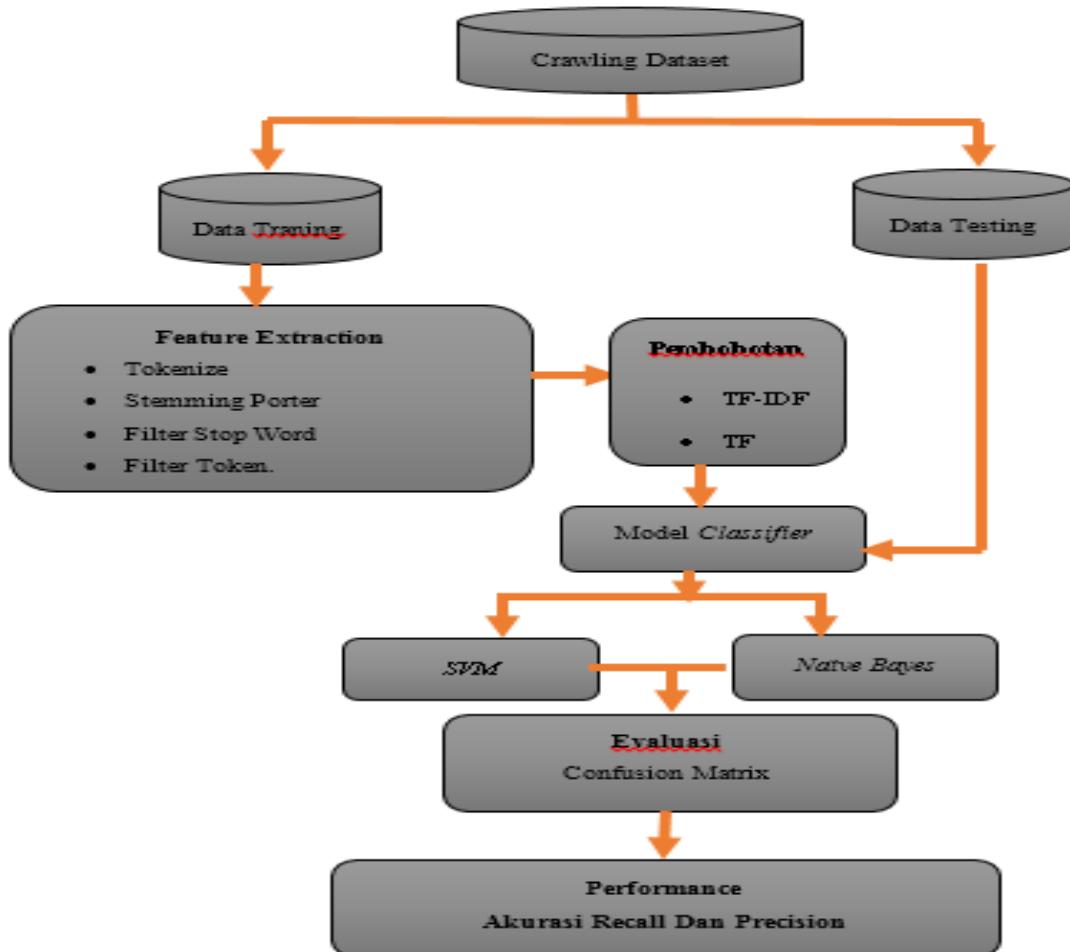
Adapun beberapa tahap-tahap dan implementasi *text processing* yang meliputi:

1. Tokenize merupakan proses untuk memisah-misahkan kata. Pemotongan kata tersebut yang sering disebut token term.
2. Filter Token merupakan pengambilan/menyaring sebuah kata yang berkarakter misal di input nilai karakter 3 maka panjang dalam sebuah karakter kata akan di filter menjadi panjang 3 karakter sesuai panjang karakter yang diinputkan.
3. Stemming yaitu proses menghilangkan



kata-kata yang tidak penting dalam teks namun sering muncul yang tidak memiliki pengaruh apapun dalam proses ekstraksi

sentimen suatu preview. Misalnya kata yang termasuk kata penunjuk waktu dan kata



tanya.

Gambar. 2 Diagram Alur Klasifikasi

Gambar.1 Desain Eksperimen

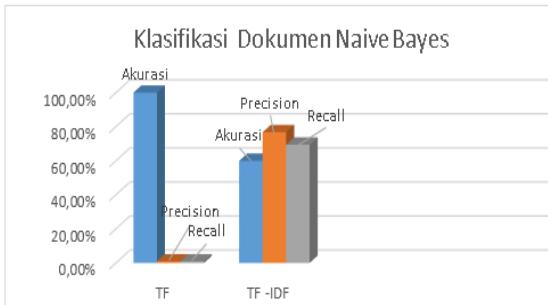
4. Coffusion Matrix

Tabel2.Cuffusion Matrix NB

Dokumen aktual	true negatif	true positif	class precision
aktual. negatif	54	214	20.15%
aktual. positif	147	487	76.81%
class recall	26.87%	69.47%	

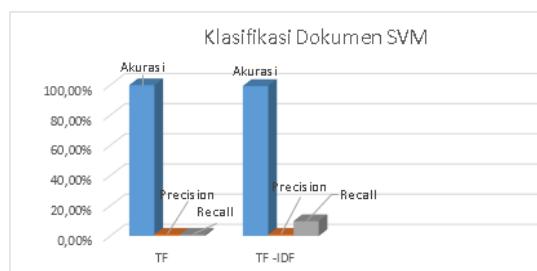
Total Akurasi: 59.98% +/-4.42% (mikro: 59.98%)

Precision dan *Recall* adalah dua perhitungan yang banyak digunakan untuk mengukur kinerja dari sistem/ metode yang digunakan. *Precision* adalah tingkat ketepatan antara informasi yang diminta oleh pengguna dengan jawaban yang diberikan oleh sistem. Sedangkan *Recall* adalah tingkat keberhasilan sistem dalam menemukan kembali sebuah informasi. Sedangkan di dunia lain seperti di dunia statistik dikenal juga istilah *accuracy*. Akurasi di definisikan sebagai tingkat kedekatan antara nilai target



diksidengannilaiaktual.Untuk proses evaluasiini dimanaakan memberikanakurasiataup erformance teks sentimen yang telah dilakukan. Untuk proses evaluasiini menggunakan *confusion matrix* yang dimanaakan memudahkan penelitiuntuk mendapatkan tingkat akurasi dari klasifikasi dokumen teks sentimen analisis berita politik.

Gambar2. Klasifikasi dokumen NB



Gambar3. Klasifikasi dokumen SVM

Untuk mencari akurasi, digunakan alat ukur confusion matrix precision dan recall dapat di gunakan rumus sebagai berikut :

$$\text{Accuracy} = \left(\frac{a + d}{a + b + c + d} \right) = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

Dokumen Aktual	true negatif	true positif	Class precision
aktual. Negative	100	19	84.0%
aktual. Positif	0	81	100.00%
class recall	100.00%	81.00%	

$$\text{Recall} = \frac{TP}{TP + FN}$$

Keterangan :

TP = True positif yang positif

FN = False positif yang negatif

FP = False negatif yang positif
TN = True negatif yang negatif

Tabel 3. Coffusion Matrix SVM

Total Akurasi: 90.50% +/- 6.87%(mikro: 90.50%)

Data yang di gunakan dalam penelitian ini juga berasal website thejakartapost.com tahun 2014 dengan menggunakan metode crawling dan menggunakan pendekatan text mining. Data uji coba adalah dokumen 700 positif dan dokumen 700 negatif hasil penambahan dari data yang sebelumnya 200 example menjadi 700 data di antaranya negatif dan positif total data 1400 data example data berita.

5. Hasil Penelitian

Hasil penelitian yang dilakukan menggunakan spesifikasi komputer AMD E- 450 APU dan sistem operasi windows 7-Ultimate 32-bit. Aplikasi yang di gunakan adalah rapidminer 5.3. Data penelitian ini di ambil dari situs portal di khususkan pada dokumen teks berbahasa inggris situs berita yang digunakan pada penelitian ini dari www.thejakartapost.com/channel/opinion. yaitu kumpulan teks opini-opini berita politik. Pada pengolahan data dari beberapa halaman tag web dan dilakukan crawling data pada website www.thejakartapost.com/channel/opinion dengan menggunakan bantuan tools dokumen subjectivity collection yaitu data kumpulan hasil crawling website dalam satu folder dan di kumpulkan dalam bentuk ekstensi.txt sebanyak 700 record data positif negatif yang terdiri dari beberapa politik harian tersebut dan di tentukan beberapa analisis sentimen di antaranya positif, netral, dan negatif dan melakukan proses klasifikasi dan ekstraksi sebuah data. Analisis sentimen wacana politik harian tersebut uji coba mengklasifikasikan opini wacana politik harian berbahasa Inggris dengan data opini yang di dapatkan secara online portal berita tersebut kemudian di proses data yang berupa teks opiniberita berbahasa Inggris melakukan preprocessing teks antara lain (pemecah kata), (menghilangkan kata tidak penting dalam teks), (menyaring kata karakter panjang), (mengurangi kata-kata dasar atau induknya)

dilanjutkan mencari performa akurasi dengan mengg

unakanmetodealgoritma Naive Bayes dan Support Vector Machine hasilnyaakanberupaefektifitaskalimatopinipositif, netral, negatifterhadapwacanapolitiktersebut. Berdasarkanpengetahuanseorangahliuntukmenent ukann class/label positifnetraldannegatif di lakukansecara manual berdasarkanorangahliwartawanmisbah (redaksikoranharianumummercusuar), yaitupenentuansecara manual kalimatsentimenanalisisdenganmenggunakan 3 klasifikasi class/ label berikut:

- pos: Politik yang membawa sentimen positif terhadap topik.
- neg: Politik yang membawa sentimen negatif terhadap topik.
- campuran: Politik yang membawa kedua positif dan negatif sentimen terhadap topik.

6. Kesimpulan

Dari

hasilpenelitianinimenujkanbahwapenggunaan metode algoritma SVM dapatmengimplementasikanhasildaripengujian data beritawacanapolitikdandaripenggunaan Naive Bayes dapat di hasilkanperbedaan denganalgoritma SVM padamasing - masing data yang telah di analisisdandarihasilpengukuranmenggunakan teknik pengukuran Precision, Recall dan di dapattingkatakurasiklasifikasiberitapolitikmasing - masing. Sehinggapadahasilanalisatersebutdapat di ambilkesimpulanbahwapenerapanalgoritma SVM menunjukkanakurasilebihbesar 90,50% SedangkanNaive Bayes di uji data sebanyak 700 data 59,98% dalamhasilklasifikasiopiniwacanapolitiksertadeng anadanyahasilanalisistersebutdapatmemperolehkla sifikasiopini yang baik.

Reference

- [1] Barber, I. 2009. *Support Vector Machine In PHP*.
- [2] Brige, C.2011. Unstructured Data and the 80 Percent Rule.
- [3] Liu B. Sentiment analysis and opinion mining, *Synth Lect Human Lang Technol* (2012).
- [4] Pang, B. & Lee, L. 2008. *Subjectivity Detection and Opinion Identification. Opinion Mining and Sentiment Analysis*.
- [5] Bo Pang and Lillian Lee *Foundations and Trends in Information Retrieval* 2(1-2), pp.1–135, 2008. Also available as a book or e-book.
- [6] Wibisono,Y. 2005. Klasifikasi Berita Berbahasa Indonesia menggunakan Naïve Bayes Classifier. diases 29 September 2012).
- [7] Saraswati, N.W.S., 2011, *Text Mining dengan Metode Naive Bayes Classifier dan Support Vector Machine untuk Sentimen Analysis*.
- [8] Wulandini, F. & Nugroho, A. N. 2009. *Text Classification Using Support Vector Machine for Webmining Based Spation Temporal Analysis of the Spread of Tropical Diseases*. International
- [9] Tan, P. N., Steinbach, M. & Kumar, V. 2006. *Introduction to Data Mining*. Boston : Pearson Addison Wesley.
- [10] Pang, B. & Lee, L. 2004. A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts. *Proceedings of the Association for Computational Linguistics (ACL)*.pp.271–278.
- [11] Parikh, R., & Movassate, M. (2009). Sentiment Analysis of User Generated Twitter Updates using Various Classification Retrieved November 12, 2011, from CS224N Final Projects 2008-9.
- [12] Prasad,S.(n.d.).Microblogging Sentiment Analysis Using Bayesian Classification Methods. Retrieved November 29, 2011, from The Stanford NLP (Natural Language Processing).

- [13] Read, J. (2005). Using Emoticons to Reduce Dependency in MachineLearning Techniques for Sentiment Classification. *ACLstudent '05 Proceedings of the ACL Student Research Workshop*(pp. 43-48). Stroudsburg: Association for Computational Linguistic.
- [14] Naradipha, A. R., & Purwarianti, A. (2011). Sentiment Classification for Indonesian Message in Social Media. *International Conference on Electrical Engineering and Informatics*, (pp.14).
- [15] Bo Pang and Lilian Lee. 2008. Opinion Mining and Sentiment Analysis, Foundations and Trends in Information Retrieval, vol. V olume 2, no. Issue 1-2, pp. 1-135.
- [16] Blitzer, J., Dredze, M. & Pereira, F. 2006. *Biographies, Bollywood, Boom boxes and Blenders: Domain Adaptation for Sentiment Classification.*
- [17] Feldman, R & Sanger, J. 2007. The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data. Cambridge University Press: New York.
- [18] Berry, M.W. & Kogan, J. 2010. Text Mining Application and theory. WILEY: United Kingdom
- [19] Indah Tri, R. 2010. Pembuatan Judul Otomatis Dokumen Berita Berbahasa Indonesia Menggunakan Metode K-Nearest Neighbor. Prodi Ilmu Komputer, Universitas Brawijaya.
- [20] Dehaff, M. 2010. Sentiment Analysis, Hard But Worth It.
- [21] Pang, B. & Lee, L. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *Proceedings of the Association for Computational Linguistic.*
- [22] Snyder B. & Barzilay R. 2007 Multiple Aspect Ranking using the Good Grief Algorithm. *Proceedings of the Joint Human Language Technology/North American Chapter of the ACL.*
- [23] Zhang, H. 2004. The Optimality of Naive Bayes. *FLAIRS2004 conference.*
- [24] Caruana, R. & Niculescu-Mizil, A. 2006. An empirical comparison of supervised learning algorithms. *Proceedings of the 23rd international conference on Machine learning,*