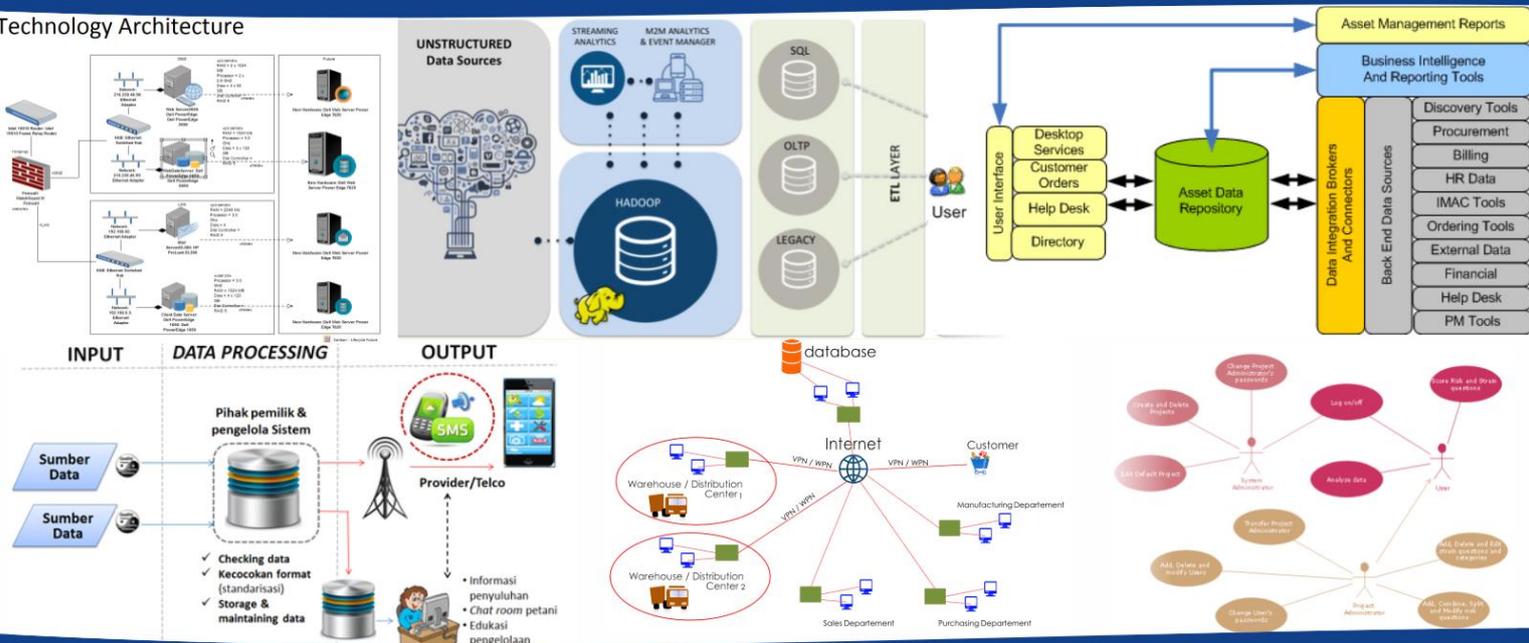




MISI

JURNAL MANAJEMEN INFORMATIKA & SISTEM INFORMASI

Technology Architecture



Diterbitkan Oleh LPPM STMIK Lombok
Jln. Basuki Rahmat No.105 Praya, Lombok Tengah - NTB
Telp dan Fax (0370) 654310 - e-journal.stmiklombok.ac.id/jsi
email. lppm@stmiklombok.ac.id



DEWAN REDAKSI

JURNAL MISI (JURNAL MANAJEMEN INFORMATIKA DAN SISTEM INFORMASI)

Jurnal Manager

Wire Bagye, S.Kom.,M.Kom - STMIK Lombok, SINTA ID : 5992010

Reviewer :

Resad Setyadi.,S.T.,S.Si.,MMSI,Ph,D (cand)- Institut Teknologi Telkom Purwokerto

SCOPUS ID 57204172534, SINTA ID : 6113570

Yesaya Tommy Paulus, S.Kom., MT., Ph.D. - STMIK Dipanegara Makassar

SCOPUS ID 57202829909, SINTA ID : 6002004

Lalu Mutawalli, S.Kom.,M.I.Kom.,M.Kom - STMIK Lombok

SCOPUS ID : 57205057118, SINTA ID : 6659709

Saruni Dwiasnati, ST.,MM.,M.Kom - Universitas Mercu Buana

SCOPUS ID : 57210968603, SINTA ID : 6150854

Ida Bagus Ary Indra Iswara, S.Kom.,M.Kom - STMIK STIKOM Indonesia

SCOPUS ID 57203711945, SINTA ID : 183498

Erlin Windia Ambarsari - Universitas Indraprasta PGRI

SCOPUS ID : 56242503900, SINTA ID : 5998887

Fachrudin Pakaja, S.Kom, M.T - Universitas Gajayana

SINTA ID : 6164357

Ahmad Jufri, S.Kom., M.T - Sekolah Tinggi Teknologi Stikma Internasional

SINTA ID : 172241

Mohammad Taufan Asri Zaen, ST.,MT - STMIK Lombok

SINTA ID : 5992087

Hairul Fahmi, S.Kom., M.Kom - STMIK Lombok

SINTA ID : 5983160

I Ketut Putu Suniantara, S.Si., M.Si - ITB STIKOM Bali

SINTA ID : 6086221

Nawassyarif S. Kom., M.Pd. - Universitas Teknologi Sumbawa

SINTA ID : 6722660

Muhamad Malik Mutoffar, ST., MM., CNSS - Sekolah Tinggi Teknologi Bandung

SINTA ID : 6013819

Editor :

Ahmad Susan Pardiansyah S.Kom.,M.Kom - STMIK Lombok

Wire Bagye, S.Kom.,M.Kom - STMIK Lombok

Vrestanti Novalia Santosa, M.Pd. - Universitas Tribuana Kalabahi

Desain Grafis & Web Maintenance

Jihadul Akbar, S.Kom - STMIK Lombok

Secretariat

Maulana Ashari, M.Kom - STMIK Lombok

DAFTAR ISI

1	AUDIT TEKNOLOGI INFORMASI PADA SISTEM PERKREDITAN ONLINE TERPADU BANK XYZ CABANG PERAWANG MENGGUNAKAN ITIL V3	90 -99
	<i>M. Khairul Anam, Ade Riyanda Putra, Sofiansyah Fadli, Muhammad Bambang Firdaus, Fadli Suandi, Lathifah</i>	
2	SISTEM PENJADWALAN EVENT ORGANIZER DENGAN METODE ROUND ROBIN (RR)	100-107
	<i>Sofiansyah Fadli, Maulana Ashari, Khairul Imtihan</i>	
3	APLIKASI PENDAFTARAN SISWA BARU MENGGUNAKAN ALGORITMA <i>BEST FIRST SEARCH</i> PADA SMP NEGERI 1 MEDAN	108-115
	<i>Maulana Ikhsan, Muhammad Irwan Padli Nasution, Ali Ikhwan</i>	
4	IMPLEMENTASI SCRUM DALAM PENGEMBANGAN SISTEM INFORMASI JASA DESAIN GRAFIS	116-122
	<i>Lalu Mutawali, Buyung Kurnia Fathoni, Hasyim Asyari</i>	
5	RANCANG BANGUN APLIKASI E VOTING BERBASIS ANDROID MENGGUNAKAN FRAMEWORK 7 STUDI KASUS DI PIMPINAN CABANG IPNU IPPNU KABUPATEN JOMBANG	123-130
	<i>Hudan Aminulloh, Ivan Dwi Fibrian, Mukhammad Masrur</i>	
6	SISTEM INFORMASI GEOGRAFIS LOKASI PRAKTEK DOKTER DI KOTA PALEMBANG BERBASIS MOBILE WEB	131-137
	<i>Ari Muzakir, Alfian Egi Erlangga</i>	
7	DATA MINING KETERKAITAN ANTARA KEBERADAAN TAMBAK MENURUT JENIS IKAN PADA KABUPATEN ATAU KOTA DI PROVINSI JAWA TENGAH DENGAN ALGORITMA A PRIORI	138-145
	<i>Tohirin, Widhy Al Mauludyansah, Sanjaya Endra Setyawan, Ronny Regawa Budiman Djatisara</i>	
8	APLIKASI PREDIKSI PENJUALAN AC MENGGUNAKAN DECISION TREE DENGAN ALGORITMA C4.5	146-156
	<i>Ade Izyuddin, Setyawan Wibisono</i>	
9	RANCANG BANGUN SISTEM PENGARSIPAN SURAT KEDINASAN BERBASIS WEB MENGGUNAKAN FRAMEWORK CODEIGNITER	157-165
	<i>Puja Irawan, Dimas Aulia Pudjie Prasetya, Petrus Sokibi</i>	
10	KLASIFIKASI KOMENTAR PUBLIK TERHADAP KEBIJAKAN PEMERINTAH PADA FACEBOOK FRONTPAGE KOMPAS MENGGUNAKAN NAIVE BAYES	166-173
	<i>I Wayan Dikse Pancane, I Wayan Suriana</i>	

KLASIFIKASI KOMENTAR PUBLIK TERHADAP KEBIJAKAN PEMERINTAH PADA FACEBOOK FRONTPAGE KOMPAS MENGUNAKAN CLUSTERING K-MEANS, FURTHEST FIRST

I Wayan Dikse Pancane¹, I Wayan Suriana²

Program Studi Teknik Elektro, Fakultas Teknik dan Informatika, Universitas Pendidikan Nasional
Jl. Bedugul No. 39 Sidakarya, Denpasar, Bali 80224
diksapancane@undiknas.ac.id, wayansuriana@undiknas.ac.id

ABSTRACT

The number of Facebook users in Indonesia is very high. This condition is used as an opportunity by the online media like news portal by making a page on facebook. The content of their wall is a preview or short description about news / articles which is posted on their official website. By using k-Means and Farthest-First, the wall can be automatically grouped by similarity topics. Preprocess in this grouping or clustering is done by Porter Stemmer Stemmer and Naizef Stemmer. Amount of data to be tested in this study were 466 wall. The best clustering quality produced by the k-Means $k = 2$ without pre-stemming process. Accuracy achieved for the labeling of "national" and "non-national" is 92.92%. Clustering wall labeled "national" into the label of "corruption" and "non-corruption", produced 77.78% accuracy. This result was achieved by the k-means clustering with $k = 2$ pre-process Stemming Nazief.

Keywords: Clustering, K-Means, Farthest-First, Porter Stemmer, Nazief Stemmer

ABSTRAK

Begitu tingginya jumlah pengguna Facebook di Indonesia membuat media cetak nasional di Indonesia juga membuat Facebook *page* dengan isi *wall* adalah cuplikan dari berita yang ada di websitenya. Dengan menggunakan k-Means dan Farthest-First, *wall* tersebut dapat dikelompokkan secara otomatis berdasarkan kesamaan topik bahasannya. Pre-proses dari pengelompokkan ini menggunakan Porter Stemmer dan Nazief Stemmer. Dari hasil uji coba 466 data *wall* facebook, cluster terbaik didapatkan dengan k-Means $k=2$ tanpa pre-proses *stemming*. Akurasi yang dicapai untuk pelabelan "nasional" dan "non-nasional" adalah 92.92%. Clustering *wall* "nasional" terbaik juga dihasilkan menggunakan k-Means $k=2$ dengan menerapkan pre-proses Nazief. Akurasi yang didapatkan dengan label "korupsi" dan "non-korupsi" adalah 77.78%.

Kata kunci: Pengelompokan, K-Means, Farthest-First, Porter Stemmer, Nazief Stemmer

1. PENDAHULUAN

Disaat internet sudah merupakan kebutuhan sehari-hari di negara kita seperti saat ini maka kebutuhan informasi terkini sangat tinggi. Pusat penyedia informasi seperti media TV, Radio dan portal berita di internet menjadi point-point media yang paling dicari oleh masyarakat. Media informasi seperti media cetak harian, mingguan, bulanan apalagi tiga bulan sudah menjadi sesuatu yang tidak relevan lagi mengingat keterkinian beritanya sangat tertinggal.

Media TV dan Radio walaupun mampu menyajikan informasi terupdate masih menjadi media informasi yang mahal. Ketersediaan

perangkat untuk melakukan pemancaran siaran serta daya jangkauan pemancar masih menjadi kendala. Dengan mahalnya dana yang dibutuhkan maka stasiun TV masih berpikir ulang untuk melakukan siaran langsung. Media Internet adalah satu-satunya yang paling tepat untuk mengatasi berbagai kesulitan tersebut, Internet mampu membawa informasi langsung ke orang yang membutuhkan dengan cepat dan memungkinkan real time. Apalagi saat ini semua provider seluler telah membuka layanan internet dengan murah dan cepat.

Hal ini membuat media-media informasi cetak akhirnya perlahan mulai melakukan manuver dengan menduplikasi konten dari media cetaknya ke

dalam bentuk digital. Hampir semua koran nasional saat ini mempunyai koran versi onlinenya. Contoh Harian Kompas (kompas.com), Harian JawaPos (jawapos.co.id), Majalah Tempo (majalah.tempointeraktif.com) dan masih banyak lagi. Belum lagi majalah teknologi, kesehatan dan yang lainnya. Kebutuhan pembaca akan informasi yang terbatas pada versi cetak diakomodir pada versi onlinenya. Selain menjadi lebih update biaya yang dikeluarkan masyarakat untuk mendapatkannya pun jauh lebih murah. Hanya bermodalkan HP, laptop, komputer atau tablet serta koneksi internet baik via modem HP, telpon rumah bahkan wifi gratis, masyarakat sudah langsung bisa mendapatkan informasi.

Walaupun mampu mengatasi segala faktor keterkinian, teknologi juga datang tanpa terlepas dari eksese dan efek negatifnya. Di Internet, siapa saja dapat menjadi siapa saja dan apa saja. Jika pengguna internet tidak waspada dan berhati-hati tentunya selain mendapatkan keuntungan bisa jadi juga mendapatkan kerugian. Demikian juga dengan berjamurnya keberadaan media-media online di internet. Tidak jarang situs-situs informasi yang ada di internet menyajikan informasi yang tidak akurat. Beberapa situs pribadi seperti blog lebih sering mengungkapkan opini dibandingkan fakta. Tentunya hal ini menjadi sangat berbahaya jika pembaca akhirnya mengikuti bentuk-bentuk opini yang kemungkinan jauh melenceng dari fakta.

Keuntungan lain dari media-media informasi online ini adalah adanya interaktif antara pembaca dan media. Hal ini sebenarnya terjadi juga dengan media cetak. Surat pembaca, opini, dan tajuk adalah contoh-contoh cara menjembatani pembaca dengan media. Tetapi sifatnya sangat kaku jika dibandingkan dengan media internet. Pada media informasi internet, pembaca bisa langsung mengomentari setiap berita yang muncul. Hampir di seluruh media informasi internet menyediakan fitur komentar yang memungkinkan pembaca memberikan komentar dengan mudah dan tanpa banyak aturan. Langsung tulis komentar menggunakan ID di internet apakah itu menggunakan akun media social seperti facebook atau twitter atau ada media yang mengharuskan pembaca untuk menjadi anggota gratis dulu untuk bisa menuiskan komentar.

Hal ini menjadi sangat menarik karena memperlihatkan bagaimana masyarakat kita selain haus akan informasi kekinian juga ternyata sangat concern masalah-masalah sosial yang timbul. Masyarakat bahkan berlomba-lomba memberika solusi yang tertepat dalam berbagai

kasus. Masyarakat juga tidak ragu-ragu mengancam, menghujat atau justru mendukung terhadap sebuah kejadian. Untuk media-media internet yang berstatus koran nasional dan sudah terpercaya di masyarakat kita, hal ini sebenarnya bisa dimanfaatkan juga oleh pemerintah dalam hal mendapatkan reaksi yang cepat dan tepat dari masyarakat.

Dengan latar belakang inilah penulis tertarik untuk mengeksplorasi secara lebih detail dari semua komentar pembaca yang ada pada portal berita kompas yang berhubungan dengan pemerintahan dan situasi politik di Indonesia. Komentar akan dikoleksi menjadi sebuah data untuk dianalisa dan didapatkan cluster-cluster komentar dari masyarakat tersebut dengan harapan dapat membantu pemerintah didalam mengetahui secara cepat tingkat keberpihakan masyarakat terhadap program-program pemerintah atau kasus-kasus tertentu yang sedang terjadi.

2. Tinjauan Pustaka

2.1 Text Mining

Text mining dapat diartikan sebagai penemuan informasi yang baru dan tidak diketahui sebelumnya oleh komputer, dengan secara otomatis mengekstrak informasi dari sumber-sumber yang berbeda. Kunci dari proses ini adalah menggabungkan informasi yang berhasil diekstraksi dari berbagai sumber (Hearst, 2003). Sedangkan menurut (Milkha Harlian) *text mining* memiliki definisi menambang data yang berupa teks dimana sumber data biasanya didapatkan dari dokumen, dan tujuannya adalah mencari kata-kata yang dapat mewakili isi dari dokumen sehingga dapat dilakukan analisa keterhubungan antar dokumen. Dengan *text mining* tugas-tugas yang berhubungan dengan penganalisaan teks dengan jumlah yang besar, penemuan pola serta penggalian informasi yang mungkin berguna dari suatu teks dapat dilakukan.

Sebagai bentuk aplikasi dari *text mining*, sistem klasifikasi berita menggunakan berita sebagai sumber informasi dan informasi klasifikasi sebagai informasi yang akan diekstrak dari sumber informasi. Informasi klasifikasi dapat berbentuk angka, angka probabilitas, set aturan atau bentuk lainnya.

Walaupun inti dari suatu sistem klasifikasi adalah tahap penemuan pola (*pattern discovery*) namun secara lengkap proses *text mining* dibagi menjadi 3 tahap utama, yaitu proses awal terhadap teks (*text preprocessing*), transformasi teks ke dalam

bentuk antara (*text transformation/feature generation*), dan penemuan pola (*pattern discovery*). (Even dan Zohar, 2002). Masukan awal dari proses ini adalah suatu data teks dan menghasilkan keluaran berupa pola sebagai hasil interpretasi.

2.2 Text Preprocessing

Tahapan awal dari *text mining* adalah *text preprocessing* yang bertujuan untuk mempersiapkan teks menjadi data yang akan mengalami pengolahan pada tahapan berikutnya. Beberapa contoh tindakan yang dapat dilakukan pada tahap ini, mulai dari tindakan yang bersifat kompleks seperti *partofspeech (pos) tagging*, *parse tree*, hingga tindakan yang bersifat sederhana seperti proses parsing sederhana terhadap teks, yaitu memecah suatu kalimat menjadi sekumpulan kata. Selain itu pada tahapan ini biasanya juga dilakukan *casefolding*, yaitu pengubahan karakter huruf menjadi huruf kecil. Proses *partofspeech* melakukan parsing terhadap seluruh kalimat dalam teks kemudian memberikan peran kepada setiap kata, misalnya : petani (subyek) pergi (predikat) ke (kata hub) sawah (keterangan). Hasil dari *partofspeech tagging* dapat digunakan untuk *parse tree*, di mana masing masing kalimat berdiri sebagai sebuah pohon mandiri. Untuk proses parsing sederhana tidak dibangun *parse tree* seperti carasebelumnya. Pada proses parsing sederhana sistem akan memecah teks menjadi sekumpulan kata kata, yang kemudian akan dibawa sebagai input untuk tahap berikutnya pada proses *text mining*.

2.3 Text Transformation (feature generation)

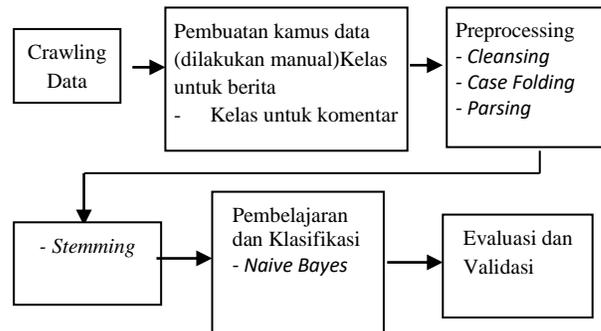
Pada tahap ini hasil yang diperoleh dari tahap *text preprocessing* akan melalui proses transformasi. Adapun proses transformasi ini dilakukan dengan mengurangi jumlah kata kata yang ada dengan penghilangan *stopword* dan juga dengan mengubah kata kata ke dalam bentuk dasarnya (*stemming*). *Stopword* adalah kata kata yang bukan merupakan ciri (kata unik) dari suatu dokumen seperti kata sambung, kata kepunyaan. Memperhitungkan *stopword* pada transformasi teks akan membuat keseluruhan sistem *text mining* bergantung kepada faktor bahasa. Hal ini menjadi kelemahan dari proses penghilangan *stopword*. Namun proses penghilangan *stopword* tetap digunakan karena proses ini akan sangat mengurangi beban kerja system. Dengan menghilangkan *stopword* dari suatu teks maka sistem hanya akan memperhitungkan kata kata yang dianggap penting.

Stemming adalah contoh tindakan lain yang dapat dilakukan pada tahap transformasi teks. *Stemming* adalah proses untuk mereduksi kata ke bentuk dasarnya. Kata yang memiliki bentuk dasar samawalaupun imbuhan berbeda seharusnya memiliki kedekatan arti. Disamping itu juga, proses *stemming* akan sangat mengurangi jumlah dan beban database. Jika setiap kata disimpan tanpa melalui proses *stemming*, maka satu macam kata dasar saja akan disimpan dengan berbagai macam bentuk yang berbeda sesuai dengan imbuhan yang mungkin melekatinya. Hal ini sangat berbeda jika kita menerapkan proses *stemming* pada tahap ini, satu kata dasar hanya akan disimpan sekali walaupun mungkin kata dasar tersebut pada sumber data sudah berubah dari bentuk aslinya dan mendapatkan berbagai macam imbuhan.

Proses *stemming* dan penghilangan *stopword* dapat digunakan secara mandiri atau tergabung, dimana dilakukan proses penghilangan *stopword* terlebih dahulu yang diikuti dengan proses *stemming*. Hal ini dilakukan untuk menemukan pola dari teks dalam berita tersebut. Karena pada Tugas Akhir ini menggunakan teks berita dalam bahasa Indonesia sebagai objek penelitian maka akan dibahas proses *stemming* pada teks/kata berbahasa Indonesia pada bahasan selanjutnya.

3. Metodologi Penelitian

3.1 Diagram alir Penelitian



Gambar 1. Blok diagram aliran kerja penelitian

3.2 Pemodelan Sistem

1. Model Probabilistic Naive Bayes

Model probabilitas untuk classifier adalah model kondisional Lihat hasil pada table 1 dan 2

$$p(C|F_1, \dots, F_n) \quad (2.1)$$

terhadap variabel kelas dependen C dengan sejumlah ke cil hasil atau kelas, tergantung pada beberapa

variabel fitur F1 sampai Fn. Masalahnya adalah bahwa jika jumlah fitur n besar atau bila fitur bisa mengambil sejumlah besar nilai, maka membuat sebuah model pada tabel probabilitas adalah tidak mungkin. Oleh karena itu kita mereformulasi model untuk membuatnya lebih fleksibel.

Menggunakan teorema Bayes, kita menulis

$$p(C|F_1, \dots, F_n) = \frac{p(C)p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)} \quad (2.2)$$

Dalam bahasa Inggris persamaan di atas dapat ditulis sebagai

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}} \quad (2.3)$$

Dalam praktiknya kita hanya tertarik pada pembilang dari persamaan tersebut, karena penyebut tidak tergantung pada C dan nilai-nilai fitur F_i diberikan, sehingga penyebut secara efektif konstan. Pembilang ini setara dengan model probabilitas gabungan $p(C|F_1, \dots, F_n)$ yang dapat ditulis ulang sebagai berikut, menggunakan penggunaan berulang dari definisi probabilitas bersyarat:

$$\begin{aligned} p(C|F_1, \dots, F_n) &= p(C)p(F_1, \dots, F_n|C) \\ &= p(C)p(F_1|C)p(F_2, \dots, F_n|C, F_1) \\ &= p(C)p(F_1|C)p(F_2|C, F_1)p(F_3, \dots, F_n|C, F_1, F_2) \\ &= p(C)p(F_1|C)p(F_2|C, F_1)p(F_3|C, F_1, F_2)p(F_4, \dots, F_n|C, F_1, F_2, F_3, \dots, F_{n-1}) \\ &= p(C)p(F_1|C)p(F_2|C, F_1)p(F_3|C, F_1, F_2) \dots p(F_n|C, F_1, F_2, \dots, F_{n-1}) \end{aligned} \quad (2.4)$$

Sekarang asumsikan mandiri bersyarat yang "naif" memegang peranan.

Menganggap bahwa setiap fitur F_i adalah secara kondisi independen terhadap setiap fitur lainnya F_j untuk j ≠ i. Ini berarti bahwa

$$p(F_i|C, F_j) = p(F_i|C) \quad (2.5)$$

untuk i ≠ j, sehingga joint model dapat dinyatakan sebagai

$$\begin{aligned} p(C|F_1, \dots, F_n) &= p(C)p(F_1|C)p(F_2|C)p(F_3|C) \dots \\ &= p(C) \prod_{i=1}^n p(F_i|C) \end{aligned} \quad (2.6)$$

Ini berarti bahwa dibawah asumsi independen di atas, distribusi bersyarat dari variabel kelas C dapat dinyatakan seperti ini:

$$p(C|F_1, \dots, F_n) = \frac{1}{Z} p(C) \prod_{i=1}^n p(F_i|C) \quad (2.7)$$

dimana Z (bukti) adalah faktor skala tergantung hanya pada F₁, ..., F_n, yaitu, sebuah konstanta jika nilai dari variabel fitur diketahui.

Model dari bentuk ini jauh lebih mudah dikelola, karena mereka memecah menjadi class prior p(C) dan distribusi probabilitas independen

p(F_i|C). Jika ada k kelas dan jika model untuk masing-masing p(F_i|C) = c dapat dinyatakan dalam bentuk parameter, maka model naif Bayes yang sesuai memiliki (k-1) + nrk parameter. Dalam prakteknya, sering k = 2 (klasifikasi biner) dan r = 1 (variabel Bernoulli sebagai fitur) yang umum, sehingga jumlah parameter model Bayes naif adalah 2n+1, dimana n adalah jumlah fitur biner yang digunakan untuk klasifikasi dan prediksi.

2. Estimasi Parameter

Semua model parameter (yaitu, prior kelas dan distribusi probabilitas fitur) dapat didekatkan dengan frekuensi relatif dari himpunan pelatihan. Ini merupakan perkiraan kemungkinan maksimum dari probabilitas. Sebuah prior class dapat dihitung dengan asumsi kelas *equiprobable* (yaitu, prior = 1 / (jumlah kelas)), atau dengan menghitung perkiraan probabilitas kelas dari himpunan pelatihan (yaitu, (prior untuk kelas tertentu) = (jumlah sampel di kelas) / (jumlah sampel)). Untuk memperkirakan parameter untuk distribusi fitur ini, seseorang harus mengasumsikan distribusi atau menghasilkan model

untuk parameter untuk fitur.

Jika seseorang berhadapan

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad (2.8)$$

dengan data kontinu, asumsi khas adalah distribusi Gaussian, dengan parameter model dari mean dan varians.

Mean, μ, dihitung dengan dimana N adalah jumlah sampel dan x_i adalah nilai dari suatu contoh yang diberikan.

$$\text{Varian dihitung dengan } \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \quad (2.9)$$

Jika sebuah kelas tertentu dan nilai fitur tidak pernah terjadi bersama-sama dalam himpunan pelatihan maka estimasi probabilitas berbasis frekuensi akan menjadi nol. Hal ini bermasalah karena akan menghapus seluruh informasi dalam probabilitas lain ketika mereka dikalikan. Oleh karena itu sering diinginkan untuk memasukkan koreksi sampel kecil dalam semua perkiraan probabilitas bahwa tidak ada probabilitas untuk menjadi persis nol.

3.3 Penyusunan Algoritma

3.3.1. Crawling Data

Data komentar facebook Kompas ditarik memanfaatkan GraphAPI dari Facebook. Berikut adalah hasil XMLnya.

```
ID: 398904863461485_5074958
From ID: 100000202236281
From Name: Mutia Nugraha
Message: makan pakis enak ke timbang pks hehhehehehe
Timestamp: 2012-04-07T12:50:49+0000

ID: 398904863461485_5074959
From ID: 100001648513699
From Name: Buyut Jaga Makam
Message: 2014 golput lebih baik.
Timestamp: 2012-04-07T12:50:51+0000

ID: 398904863461485_5074960
From ID: 100001594936036
From Name: Mhail Suarez
Message: pak tifatul,,knp mesti takut,,kalian punya power koq,,,
Timestamp: 2012-04-07T12:50:52+0000

ID: 398904863461485_5074961
From ID: 100002163565224
From Name: Een Egi
Message: gtu dong.
Timestamp: 2012-04-07T12:50:56+0000

ID: 398904863461485_5074963
From ID: 100002244511563
From Name: Harri Mimi
Message: hahahaha
Timestamp: 2012-04-07T12:51:02+0000
```

Gambar 3.2 Contoh hasil retrieve XML komentar Facebook

Data hasil retrieve ini dimasukkan ke database MySQL sebagai data training.

autoid	iduser	name	comment	date
440	100001744211876	Fatimah Wulandari	Hmmm. PKS bgai telur diujung tanduk. ywdh met mnun...	2012-04-
441	100000362892573	Aziz Adzy	Koreksilah diri sbum mengkoreksi org laen...	2012-04-
442	100001095828969	H Roedi Effendi Nst	me juah2 ya	2012-04-
443	10000307274325	Mahd Ajhun Piliang	ssuka hi low lah...	2012-04-
444	100000130773054	Dian Utami	mohon dukungannya teman2 dalam pemilihan ICON FK U...	2012-04-
445	100002614022440	Ahmad Fahri Ryanfariansyah	waduh gawat brgt...	2012-04-
446	100001284420843	Wily Kusnadi	Mo di ujung tanduk, mo di ujung dunia terserah. yg ...	2012-04-
447	100003177745104	Muhdesrianto Kh	Mau bela koalisi? Boleh Mau suka PKS? Silahkan M...	2012-04-
448	1850458182	Hary Agustia	setuju bgt	2012-04-
449	1671095582	Mas Adi Winamo	pada kenyataannya pemerintah tdk jadi menaikan buk...	2012-04-
450	100002908309940	Arip Gunawan	Berdayung sampan ke pulau feri, menatap bintang sam...	2012-04-
451	10000011582674	Opa Tjpek Aryana	Hwaratah...kono...f	2012-04-
452	100002913184367	Ardan Dhana	PKS??	2012-04-
453	100000062859377	Widarto Plat H	Ini yang namanya jual suara rakyat untuk cari jaba...	2012-04-
454	100001436973462	Afizal Kasim Caniogo	>>Salah tuh Sembiring. Reshuffle itu perintah Sby...	2012-04-

Gambar 3.3 Contoh komentar sudah dimasukkan ke MySQL

3.3.2. Pembuatan Kamus Kelas

Tahapan yang dilakukan dari pembuatan kamus kelas ini adalah mengambil data hasil crawling dan melakukan klasifikasi secara manual dan disimpan menjadi sebuah kamus kelas. Ada 2 (dua) kamus kelas yaitu untuk kelas kategori berita (kebijakan) dan kelas kategori pendapat (pro, kontra).

3.3.3. Pre-Processing

Tahapan yang dilakukan dari dokumen *pre-processing* adalah sebagai berikut:

1. *Cleansing*, yaitu proses membersihkan dokumen dari kata yang tidak diperlukan untuk mengurangi noise. Kata yang dihilangkan adalah karakter HTML dan simbol.
2. *Case folding*, yaitu penyeragaman bentuk huruf serta penghapusan angka dan tanda baca. Dalam hal ini yang digunakan hanya huruf latin antara a sampai dengan z.
3. *Parsing*, yaitu proses memecah dokumen menjadi sebuah kata. Hal ini sesuai dengan fitur digunakan yaitu *unigram*.

SOURCE

Ini yang kalian bilang Demokrasi ? YANG BENER DEMOCRAZY !!!

Cleaning and Folding

ini yang kalian bilang demokrasi yang bener democrazy

3.3.4. Pemilihan dan Ekstraksi Fitur

Berikut adalah proses pemilihan dan ekstraksi fitur yang akan digunakan sebagai dasar proses klasifikasi.

1. *Stemming*, bertujuan mengurangi variasi kata yang memiliki kata dasar sama. Proses *stemming* dilakukan dengan menggunakan bantuan KBBI.

Penghapusan partikel = menggelapkan
 Penghapusan milik = menggelapkan
 Penghapusan prefiks 1 = gelapkan
 Penghapusan prefiks 2 = gelapkan
 Penghapusan sufiks = gelap

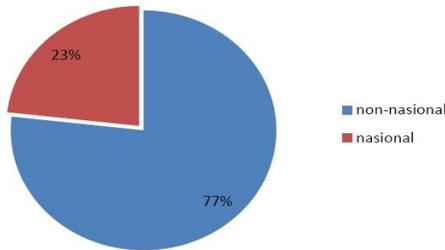
Gambar 3.5 Contoh hasil *stemming*

3.4 DATA

Data yang digunakan pada penelitian ini adalah data wall Facebook dari Facebook *page* KompasCom yang diposting selama bulan Januari sampai dengan Juli 2012. Pelabelan manual untuk melakukan interpretasi hasil *cluster*, memanfaatkan caption yang sudah ada pada setiap wall KompasCom. Jumlah data yang diolah sebesar 466 wall. Penelitian ini menggunakan 3 jenis pre-proses yaitu *wall* tanpa stem, *wall* dengan melalui Porter *stemmer* dan *wall* dengan melalui Nazief *stemmer*. Stopwords yang digunakan adalah stopwords dari penelitian Tala 2003 yaitu sebanyak 758 kata. Total wall yang diproses adalah 466 wall dengan jumlah atribut/fitur 3.665.

4. HASIL DAN PEMBAHASAN

Pengelompokan yang dilakukan oleh KompasCom pada setiap wall adalah untuk label “nasional” sebesar 108 wall dan untuk label “non-nasional” adalah 358 wall.



Gambar 2 Grafik Prosentase Jumlah Wall berdasarkan 2 label manual

Tabel 1 Purity k-Means k=2

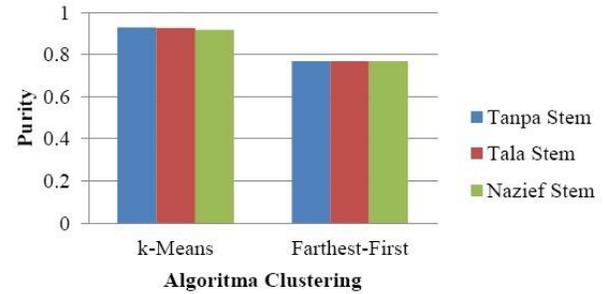
SEED	TANPA STEM	TALA STEM	NAZIEF STEM
10	0,7682	0,7682	0,8906
20	0,8863	0,8864	0,8864
30	0,9292	0,9270	0,9163
40	0,7682	0,7682	0,7682
50	0,7682	0,7682	0,7682
60	0,7682	0,7682	0,7682
70	0,7682	0,7682	0,7682
80	0,7682	0,7682	0,7682
90	0,7682	0,7682	0,7682
100	0,7682	0,7682	0,7682

Tabel 2 Purity Farthest-first k=2

SEED	TANPA STEM	TALA STEM	NAZIEF STEM
10	0,7682	0,7682	0,7682
20	0,7682	0,7682	0,7682
30	0,7682	0,7682	0,7682
40	0,7682	0,7682	0,7682
50	0,7682	0,7682	0,7682
60	0,7682	0,7682	0,7682
70	0,7682	0,7682	0,7682
80	0,7682	0,7682	0,7682
90	0,7682	0,7682	0,7682
100	0,7682	0,7682	0,7682

Percobaan dilakukan dengan menggunakan seed yang berbeda-beda untuk mendapatkan inisialisasi centroid yang paling baik. Setelah dilakukan percobaan menggunakan k-Means dengan membandingkan pre-proses tanpa stem, dengan Tala Stem dan Nazief Stem didapatkan seperti pada Tabel 1. Purity tertinggi dihasilkan oleh percobaan tanpa pre-proses stem dengan nilai 0.9292. Tabel 2 menunjukkan bahwa hasil purity yang dihasilkan Farthest- First adalah 0.7682. Nilai ini lebih kecil dari yang dihasilkan oleh k-Means.

Perbandingan k-Means dan Farthest-First ditunjukkan pada Gambar 3.



Gambar 3 Grafik Perbandingan Purity antara k-Means dan Farthest-First

Purity didapatkan dengan cara membandingkan nilai presisi maksimal setiap kelasnya. Nilai purity 0.9292 ini hampir mencapai nilai maksimal dari purity yaitu 1. Ini menandakan cluster yang dihasilkan oleh k-Means dengan pre-proses tanpa stem memiliki kualitas yang sangat baik atau bisa diartikan bahwa pembeda antar cluster sangat jelas. Tabel 1 menunjukkan bahwa pada seed 10, 20 dan 30, stemming memang menghasilkan perbaikan purity walaupun sangat kecil. Tapi kemudian di seed berikutnya dan pada Tabel 2, stemming sama sekali tidak membawa perbaikan. Hal ini menunjukkan bahwa stemming hanya membawa pengaruh yang kecil baik terhadap kualitas cluster yang dihasilkan.

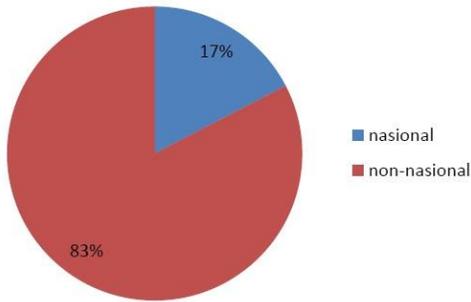
Tabel 3 Matching Matrix k-Means k=2

C0	C1	Label
30	78	Nasional
355	3	Non-Nasional

Tabel 3 menunjukkan komposisi jumlah label di masing-masing cluster. Matching Matrix yang diperlihatkan ini adalah pada purity tertinggi yaitu seperti yang dihasilkan oleh k-Means di atas. Dengan melihat jumlah label terbanyak di setiap clusternya, dapat diinterpretasikan bahwa Cluster 0 adalah kelompok yang terdiri dari wall yang memunyai topik “non-nasional”. Sedangkan Cluster 1 adalah kelompok wall yang bertopik “nasional”. Akurasi dari interpretasi ini adalah sebesar 92.92%. Tabel 4 menunjukkan hasil interpretasi dan nilai precision recall masing-masing cluster.

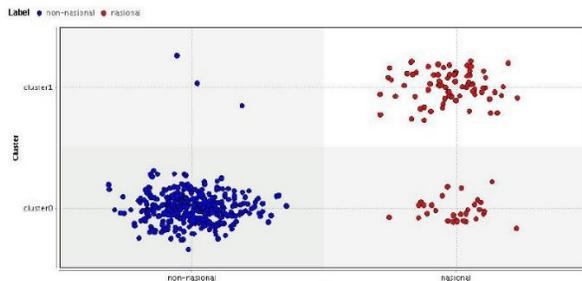
Tabel 4 Hasil Interpretasi Cluster

CLUSTER	LABEL	PRECISION	RECALL
0	Non-Nasional	0,9221	0,9916
1	Nasional	0,9630	0,7222



Gambar 4 Grafik Perbandingan Label “nasional” dan “non-nasional”

Hasil *Clustering* Kondisi *purity* yang tinggi yang dihasilkan oleh k-Means ditunjukkan pada Gambar 4.5. Walaupun sangat tinggi, namun masih belum mencapai nilai sempurna sehingga terlihat masih ada label yang menyeberang cluster.



Gambar 5 Grafik Hubungan Label Terhadap Cluster Pada k-Means k=2

Tabel 5 Centroid PerCluster

CLUSTER	CENTROID	LABEL	FEATURE CENTORID
0	0,0196	Non-nasional	Ekonomi, emirates, enam, esdm, fpi, garuda, jokowi, juli, jurnal, kalinya, kanada, kandidat, merusak, meter, mikro, nomor, nigroho, obat oktober, olimpiade, opsi, organisasi, oscar, otak, pabrik, pakai, pakar, palu, paul, pegawai, pecan, pelatihan, telekomunikasi
1	0,0526	nasional	Ode, partainya, pdip, pks, pramono

Cluster 0 adalah kumpulan wall berlabel “non-nasional” dengan centroid cluster 0.0196. Jarak terdekat dengan centroid tersebut dicapai oleh fitur-fitur seperti yang ditunjukkan oleh Tabel 5. Sedangkan Cluster 1 adalah kumpulan wall yang berlabel “nasional” dengan centroid 0.0536. Cluster “nasional” pada page KompasCom lebih banyak menuliskan kata-kata yang berhubungan dengan partai. Sedangkan cluster “non-nasional” lebih banyak menggunakan kata-kata yang bersifat umum.

5. KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan hasil percobaan dapat disimpulkan bahwa :

1. k-Means dan Farthest-First dapat digunakan untuk melakukan pengelompokan otomatis topik pada wall Facebook ber-Bahasa Indonesia.
2. Pre-proses *stemming* dalam percobaan ini mampu memberikan pengaruh perbaikan kualitas *cluster* sebesar 5%.
3. Untuk label “nasional” dan “non-nasional”, hasil percobaan terbaik didapatkan dari hasil k-Means dengan nilai *Purity* = 0.0.9262 dan akurasi sebesar 92.92%. Hasil ini didapatkan dari wall yang tidak melewati *pre-proses stemming*. *Purity* dan akurasi yang didapatkan dari percobaan menggunakan Farthest-First masih lebih kecil dari k-Means dengan rata-rata akurasi sebesar 75.69%.

5.2 Saran

Agar diperoleh hasil yang lebih maksimal, maka beberapa saran untuk pengembangan lebih lanjut adalah sebagai berikut:

1. Perlu diuji coba menggunakan teknik *clustering* yang lain sebagai perbandingan teknik yang menghasilkan *cluster* lebih baik lagi.
2. Komentar dari wall perlu diujicoba untuk *dicluster* juga untuk menemukan sentiment terhadap wall yang dikomentari.

6. Ucapan Terima Kasih

Ucapan Terima Kasih Terima kasih disampaikan kepada Universitas Pendidikan Nasional Denpasar dan Fakultas Teknik Dan Informatika Undiknas Juga Program Profesi Insinyur Undiknas yang telah memberikan masukan juga mendanai keberlangsungan jurnal ini.

DAFTAR PUSTAKA

- [1] J. Informatika, W. Mega, and P. Duhita, "CLUSTERING MENGGUNAKAN METODE K-MEANS UNTUK," vol. 15, no. 2, 2015.
- [2] D. P. Langgeni, Z. K. A. Baizal, and Y. F. A. W, "CLUSTERING ARTIKEL BERITA BERBAHASA INDONESIA," vol. 2010, no. semnasIF, pp. 1-10, 2010.
- [3] P. Ke, "K-MEANS CLUSTERING," pp. 1-16.
- [4] M. Kumar, "An Optimized Farthest First Clustering Algorithm," no. November 2013, 2019, doi: 10.1109/NUiCONE.2013.6780070.
- [5] S. Godara, "A Comparative Performance Analysis of Clustering Algorithms," vol. 1, no. 3, pp. 441-445.
- [6] M. H. Adiya and Y. Desnelita, "Jurnal Nasional Teknologi dan Sistem Informasi Penerapan Algoritma K-Means Untuk Clustering Data Obat-Obatan Pada RSUD Pekanbaru," vol. 01, pp. 17-24, 2019.
- [7] N. Valarmathy and S. Krishnaveni, "Performance Evaluation and Comparison of Clustering Algorithms used in Educational Data Mining," no. 6, pp. 103-113, 2019.
- [8] G. Sehgal and K. Garg, "Comparison of Various Clustering Algorithms," vol. 5, no. 3, pp. 3074-3076, 2014.
- [9] A. V. D. Sano and H. Nindito, "APPLICATION OF K-MEANS ALGORITHM FOR CLUSTER ANALYSIS ON POVERTY OF PROVINCES IN INDONESIA," no. 6, pp. 141-150, 2011.
- [10] S. Gnanapriya, "Evaluation of Clustering Capability Using Weka Tool," vol. 8, no. 1, pp. 181-187.
- [11] B. K. Teknomo, "K-Means Clustering Tutorial," pp. 1-12, 2007.
- [12] D. A. Vadeyar and H. K. Yogish, "Farthest First Clustering in Links Reorganization," vol. 5, no. 3, pp. 17-24, 2014.
- [13] J. O. Ong, "IMPLEMENTASI ALGORITMA K-MEANS CLUSTERING UNTUK MENENTUKAN STRATEGI MARKETING," no. April, pp. 10-20, 2013.
- [14] M. M. C. K-means, "VOL. 9 NO. 1 April 2016," vol. 9, no. 1, pp. 94-101, 2016.
- [15] M. H. Zafar and M. Ilyas, "A Clustering Based Study of Classification Algorithms," vol. 8, no. 1, pp. 11-22, 2015.
- [16] M. M. Hassan, "ENHANCING CLUSTERING-BASED CLASSIFICATION ALGORITHMS IN E-COMMERCE APPLICATIONS," vol. 96, no. 18, pp. 6095-6105, 2018.
- [17] S. Dasgupta and P. M. Long, "Performance guarantees for hierarchical clustering," no. July 2010.
- [18] R. Pebria, B. H. I, and I. Sugihartono, "GFK-12 : IDENTIFIKASI PENYEBARAN GEMPA DI INDONESIA DENGAN METODE CLUSTERING," pp. 366-370.
- [19] L. Mutawalli, M. T. A. Zaen, and W. Bagye, "KLASIFIKASI TEKS SOSIAL MEDIA TWITTER MENGGUNAKAN SUPPORT VECTOR MACHINE (Studi Kasus Penusukan Wiranto)," *J. Inform. dan Rekayasa Elektron.*, vol. 2, no. 2, p. 43, 2019, doi: 10.36595/jire.v2i2.117.