

KLASIFIKASI POSTING TWITTER CUACA PROVINSI DIY MENGUNAKAN ALGORITMA C4.5 UNTUK INFORMASI PADA WEB PARIWISATA

Yuli Astuti

Manajemen Informatika STMIK AMIKOM Yogyakarta
Jl Ring Road Utara, Condongcatur, Sleman, Yogyakarta 55281
yuli@amikom.ac.id

ABSTRAK

Setiap hari server Twitter menerima data tweet dengan jumlah yang sangat besar. Dengan demikian, perlu dilakukan pemanfaatan data dengan dikembangkan sistem untuk melakukan data mining dari tumpukan data tersebut yang akan digunakan untuk kepentingan tertentu, salah satunya adalah untuk visualisasi kondisi daerah tertentu berdasarkan cuaca. Teknik klasifikasi akan diterapkan untuk mengklasifikasikan data tweet yang menginformasikan cuaca di kota Yogyakarta. Sebelum dilakukan klasifikasi, data tweet melalui preprocessing dan pembobotan term frequency dan tf-idf. Kemudian dari pembobotan ini menghasilkan data training dan data testing yang akan dilakukan klasifikasi menggunakan pohon keputusan dengan algoritma C4.5 untuk memprediksi cuaca di daerah kota Yogyakarta untuk kemudian diinformasikan pada web pariwisata dengan pemanfaatan web service. Hasil pengujian dengan perangkat lunak Rapid Miner 5.3 diperoleh nilai akurasi terkecil 71% dengan sampel sebanyak 100 dan nilai akurasi tertinggi 95,58% dengan sampel 15106 dengan algoritma C4.5.

Kata Kunci : Twitter, cuaca, klasifikasi, algoritma C4.5

ABSTRACT

Every day the Twitter server receives data tweet with a very large number. Thus, it is necessary to use the data to develop a system to perform data mining on a pile of data that will be used for specific purposes, one of which is for the visualization of certain areas based on the weather conditions. Classification techniques will be applied to classify the data tweet informing the weather in the city of Yogyakarta. Before the classification, the data tweets through preprocessing and weighting term frequency and tf-idf. Then from this weighting generates training data and data testing to be performed classification using C4.5 decision tree algorithm to predict the weather in the city of Yogyakarta to then be informed on the tourism web with the use of the web service. Hasil pengujian with Rapid Miner 5.3 software obtained value of the smallest 71% accuracy with a sample of 100 and the highest value of 95.58% accuracy dengan sampel 15106 with C4.5 algorithm.

Keywords: Twitter, weather, classification, algorithm C4.5

1. Pendahuluan

1.1 Latar Belakang

Saat ini dunia internet sedang berada pada fase *user generated content*, yang berarti seluruh konten yang berada di internet adalah buatan pengguna secara umum. Dengan demikian, internet seperti gudang super besar yang diisi oleh pengguna dan pengguna juga dapat menggunakan isi gudang tersebut. Salah satu aplikasi internet yang mendukung *user generated content* adalah *microblogging*. Saat ini, *microblogging* menjadi populer sebagai alat

komunikasi antara pengguna internet. Jutaan pesan setiap hari dikirim oleh pengguna melalui situs *microblogging* ini. Pengguna biasanya menulis tentang hal-hal yang terjadi dalam kehidupan sehari-harinya, berdiskusi tentang berbagai topik yang sedang banyak diperbincangkan dan lain-lain. Salah satu aplikasi internet yang memungkinkan pengguna dapat berbagi konten secara bebas adalah aplikasi Twitter. Twitter adalah salah satu situs *microblogging* paling populer yang dapat membagikan update status atau pesan yang

tidak lebih dari 140 karakter kepada pengguna dalam satu jaringan dalam minat yang sama dikenal dengan istilah *follower*, atau pengguna yang mengikuti karena memiliki kesamaan minat, ketertarikan pada isu tertentu.

[1] menyatakan bahwa pada pertengahan tahun 2010 Twitter memiliki pengguna lebih dari 106 juta pengguna diseluruh dunia dan terus meningkat setiap harinya sebanyak 300.000 pengguna dan Twitter setiap harinya mendapatkan lebih dari 3 juta request. Dari angka tersebut Indonesia menjadi negara yang menduduki peringkat 8 dalam mengakses situs Twitter. Twitter menerima tweet dari pengguna sebanyak 55 juta pesan setiap harinya.

Berdasarkan data tersebut, Twitter memiliki sumber data yang besar. Data dalam hal ini adalah tweet dari pengguna yang berjumlah sangat banyak. Hal ini merupakan sebuah sumber daya yang bisa kita manfaatkan untuk kepentingan tertentu misalkan untuk mengetahui cuaca di suatu kota berdasarkan tweet yang dikirim oleh pengguna yang berisi informasi cuaca.

Saat ini pengguna Twitter, khususnya yang berada di kawasan kota Yogyakarta sering menginformasikan melalui Twitter mengenai informasi yang berhubungan dengan kota Yogyakarta. Informasi tersebut antara lain, informasi kegiatan (event), iklan (advertising), lowongan kerja di kota Yogyakarta dan sekitarnya, informasi cuaca dan lain-lain. Kumpulan informasi di Twitter yang berhubungan dengan kota Yogyakarta biasanya ditandai dengan hastag #eventyk untuk informasi kegiatan, #infoyk untuk informasi umum, #cuacayk untuk informasi cuaca, dan #diskonyk untuk informasi diskon.

Kota Yogyakarta merupakan kota yang sering mengalami cuaca tidak menentu disetiap daerahnya, misalnya daerah condong catur cerah namun daerah jalan kaliurang hujan. Terkait dengan ini, sebagian besar pengguna Twitter di kota Yogyakarta sering menginformasikan cuaca disekitarnya melalui media Twitter dengan hastag #cuacayk. Dengan demikian, apabila data Twitter tersebut dikumpulkan dan kemudian dilakukan knowledge discovery terhadap data tersebut, diharapkan dapat memberikan informasi yang akurat mengenai cuaca di kota Yogyakarta yang akan dijadikan sebagai bahan pertimbangan dalam pengambilan keputusan oleh wisatawan. Hal ini, seperti yang diungkapkan [2] bahwa

sekarang ini melalui data twitter, dapat menentukan gaya hidup, sikap dan perilaku seseorang serta keinginan-keinginannya.

Data mining merupakan sebuah proses dari knowledge discovery (penemuan pengetahuan) dari data yang sangat besar [3]. Sementara itu [4] berpendapat bahwa data mining adalah proses secara otomatis untuk menemukan informasi yang berharga dari repositori data yang sangat besar. Dengan demikian, dari tumpukan data tersebut akan didapat beragam informasi yang berharga dan penting yang sebelumnya tidak diketahui.

Salah satu cabang data mining yang mengkhususkan pada penggalian informasi dari data yang berupa data teks adalah text mining. Text mining merupakan bidang data mining yang bertujuan untuk mengumpulkan informasi yang berguna dari data teks dalam bahasa alami atau proses analisis data teks kemudian mengekstrak informasi yang berguna untuk tujuan tertentu [5]. Dengan demikian, text mining merupakan teknik yang cocok untuk melakukan ekstraksi informasi dari data tweet yang banyak tersebut sehingga dapat menghasilkan informasi yang akurat yang akan digunakan dalam proses pengambilan keputusan.

Ada banyak teknik yang bisa dilakukan untuk melakukan klasifikasi data diantaranya adalah decision tree, bayesian classifiers, bayesian belief network, *Nearest Neighbor* dan rule based classifiers [3]. Berdasarkan penelitian-penelitian sebelumnya, teknik-teknik ini mempunyai kelemahan dan kekurangan. Berdasarkan penelitian yang dilakukan oleh [6], kinerja algoritma decision tree lebih baik jika dibandingkan dengan Multiple Discriminant Analysis (MDA) dalam memprediksi kebangkrutan perusahaan. Algoritma decision tree juga merupakan algoritma paling populer dalam teknik klasifikasi.

Dalam Penelitian ini diusulkan untuk memanfaatkan data twitter untuk memprediksi cuaca kota Yogyakarta secara *realtime* untuk pemanfaatannya pada web pariwisata. Hasil penelitian ini akan menjadi sebuah model bagi daerah lain yang akan dan sedang mengembangkan sistem informasi prakiraan cuaca secara akurat dan *realtime*.

1.2 Rumusan Masalah

Dari latar belakang yang dijelaskan di atas, maka rumusan masalah yang diangkat

adalah Bagaimana mengklasifikasikan data *tweet* yang mengandung *hashtag* #cuacayk yang merupakan informasi cuaca di kota Yogyakarta?

1.3 Batasan Masalah

Proses klasifikasi data twitter cuaca Yogyakarta, dengan batasan masalah sebagai berikut :

1. Data tweet yang digunakan adalah tweet yang diposting dari tanggal 1 Juni 2015 sampai tanggal 24 September 2015.
2. Data tweet yang di download adalah data tweet yang mengandung hashtag #cuacayk
3. Tidak ada keterkaitan satu kata dengan kata yang lain dalam sebuah tweet.
4. Sesuai dengan konsep user generated content sistem ini menganggap pengguna Twitter yang menginformasikan cuaca sebagai agen informasi yang kemudian informasi ini dijadikan sumber data untuk sistem ini.
5. Cuaca yang akan dianalisis dalam sistem ini adalah cuaca di beberapa daerah kabupaten sleman, kota Yogyakarta dan Bantul

1.4 Metode Penelitian

Penelitian ini dilakukan dengan tahap-tahap sebagai berikut :

1. Pengumpulan data

Melakukan pengumpulan data tweet secara realtime yang mengandung hashtag #cuacayk. Pengumpulan data ini dilakukan dengan memanfaatkan API search yang disediakan oleh Twitter. Kemudian hasil daripencarian ini disimpan di database. Data tweet yang dikumpulkan antara lain username, isi tweet, URI tweet, tanggal dan waktu tweet dan URI profil pengguna.

2. Preprocessing

Preprocessing dilakukan dengan langkah-langkah sebagai berikut :

- a) Setelah data terkumpul, kemudian dilakukan preprocessing data sehingga data menjadi lebih mudah untuk dilakukan proses selanjutnya. Pada tahapan ini dilakukan :
 1. Penghapusan kata tertentu yang tidak dipakai dalam proses klasifikasi.
 2. Konversi menjadi huruf kecil.
 3. Menghapus url (<http://bit.ly/mHibqV>).
 4. Melakukan perbaikan data apabila ada kata yang diperlukan tetapi data tidak

sesuai, misalkan ada kata “panaaaaasssss” kemudian di ubah menjadi “panas”.

5. Menghapus mention (@xxx).
 6. Menghapus karakter selain a-z.
 7. Mengganti sinonim (yg=yang, gerah=panas, gerimis=hujan, sumuk=panas, terik=panas).
 8. Menghilangkan stopword.
 9. Menghapus kata dengan satu karakter.
- b) Menghitung bobot tiap *term* dari data tweet yang terkumpul dengan menggunakan teknik *tf-idf*. Ditunjukkan pada Persamaan 1.

$$tfidf(d, w) = tf(d, w) \times \log N/dfw \quad \dots(1)$$

Dimana : $tf(d,w)$ adalah frekuensi kemunculan term w pada dokumen d , n adalah jumlah keseluruhan dokumen dan dfw adalah jumlah dokumen yang mengandung term w .

3. Klasifikasi

Klasifikasi dilakukan dengan langkah-langkah sebagai berikut :

- a) Pada proses *learned* model digambarkan dalam bentuk classification rule atau formula matematika yang biasa dikenal dengan algoritma. Algoritma C4.5 merupakan algoritma yang digunakan untuk membentuk pohon keputusan. Menurut [7] langkah-langkahnya sebagai berikut :
 1. Pilih variabel sebagai akar
 2. Buat cabang untuk tiap-tiap nilai
 3. Bagi kasus dalam cabang
 4. Ulangi proses untuk setiap cabang sampai semua kasus pada cabang memiliki kelas yang sama

Untuk memilih variabel sebagai akar, didasarkan pada nilai *gain* tertinggi dari variabel-variabel yang ada. Untuk menghitung *gain* digunakan rumus pada Persamaan 2.

$$Gain(S, A) = Entropy(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} \times Entropy(S_i) \dots(2)$$

Keterangan :

- S : Himpunan kasus
A : Variabel
n : jumlah partisi variabel A
 $|S_i|$: jumlah kasus pada partisi ke-i
 $|S|$: jumlah kasus dalam S

sementara itu untuk menghitung nilai entropi dapat dilihat pada Persamaan 3.

$$Entropy(S, A) = \sum_{i=1}^n -p_i \log_2 p_i \quad \dots\dots \quad (3)$$

Keterangan :

S : Himpunan Kasus
A : Fitur
n : jumlah partisi S
 p_i : proporsi dari S_i terhadap S

- b) Kelas dalam klasifikasi ini adalah kelas panas untuk tweet yang mengandung informasi cuaca panas (cerah) dan kelas hujan untuk tweet yang mengandung informasi cuaca hujan di suatu jalan atau tempat tertentu
- c) Memisahkan data untuk digunakan sebagai data training dan data testing.
- d) Dilakukan pengujian akurasi menggunakan metode *K-Fold Cross Validation* pada model klasifikasi algoritma C4.5 dengan menggunakan Rapid Miner 5.3.

4. Pengujian hasil akurasi

Analisa hasil akurasi model klasifikasi algoritma C4.5 dengan menggunakan metode 10 fold cross validation. Cara kerja metode ini yaitu Pada pengujian ini, sebanyak 10% dari jumlah posting secara bergantian dijadikan data uji sebanyak 10 kali terhadap 90% posting lainnya yang dijadikan data training. Nilai akurasi diperoleh dari rata-rata nilai akurasi dari 10 kali pengujian tersebut. Dengan demikian, setiap posting tweet akan menjadi data training dan data testing secara bergantian. Hal ini bertujuan untuk meminimalkan nilai akurasi yang dihasilkan oleh faktor kebetulan. Untuk menghitung nilai akurasinya digunakan persamaan (4).

$$\text{akurasi} = \frac{\text{jumlah klasifikasi benar}}{\text{jumlah data uji}} \times 100\% \quad (4)$$

2. Pembahasan

2.1 Persiapan Data

Data yang digunakan dalam penelitian ini adalah data tweet yang diambil secara realtime. Setiap satu menit sekali sistem memberikan request ke server Twitter untuk mengambil dokumen XML yang berisi tweet yang kemudian disimpan di database. Sesuai

dengan batasan masalah, tweet yang dikumpulkan

adalah tweet yang mengandung hastag #cuacayk. Pengambilan data tweet yang akan digunakan dalam penelitian ini dimulai sejak tanggal 1 Juni 2015 sampai tanggal 24 September 2015 yang menghasilkan data tweet sebanyak 15401 record.

Adapun struktur data tweet yang dikumpulkan dalam sistem adalah sebagai berikut :

1. User adalah data pengguna twitter yang memposting tweet.
2. Tweet adalah data konten tweet yang berisi informasi cuaca di kota Yogyakarta.
3. uriTweet adalah link yang menuju halaman tweet.
4. Date adalah tanggal tweet tersebut diposting.
5. Time adalah waktu tweet tersebut diposting.
6. Uri adalah link yang menuju halaman profil pengguna Twitter yang memposting tweet.

Dalam sistem ini, data diklasifikasi menjadi dua kelas/label yaitu kelas 'cerah' untuk tweet yang menginformasikan keadaan cuaca yang cerah atau panas di suatu jalan atau tempat, kelas 'hujan' untuk tweet yang menginformasikan keadaan cuaca yang hujan di suatu jalan atau tempat.

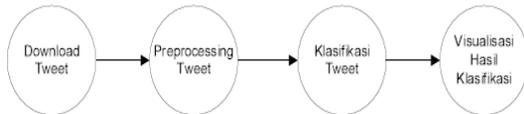
2.2 Analisis dan Pengembangan Sistem

Berdasarkan permasalahan yang telah dirumuskan dan pemahaman terhadap teori-teori dan konsep-konsep sebagai dasar pemikiran, maka pada bagian ini merupakan konsep penyelesaian masalah dengan melakukan analisis dan perancangan sistem. Sistem baru yang dirancang merupakan sistem yang dapat melakukan proses pengambilan data tweet dari API search Twitter secara otomatis dan realtime. Kemudian dari tumpukan data teks tweet tersebut dilakukan pembobotan setiap kata (term) yang kemudian hasil pembobotan tersebut diolah dengan algoritma C4.5 untuk menentukan model yang akan digunakan untuk menentukan prediksi kelas pada data tweet yang baru. Inti dari permasalahan yang diangkat adalah bagaimana melakukan ekstraksi data tweet untuk digunakan untuk training dan testing model klasifikasi algoritma C4.5 yang digunakan untuk prediksi kelas untuk data tweet yang baru. Kemudian hasil klasifikasi tersebut divisualisasikan ke Web Pariwisata yang

terkoneksi dengan web service dengan menggunakan Google Map.

Secara umum sistem ini terdiri dari empat bagian diantaranya adalah download (pengambilan data), preprocessing tweet, klasifikasi tweet dan visualisasi hasil klasifikasi.

Adapun desain arsitektur sistem dapat terlihat pada Gambar 1.

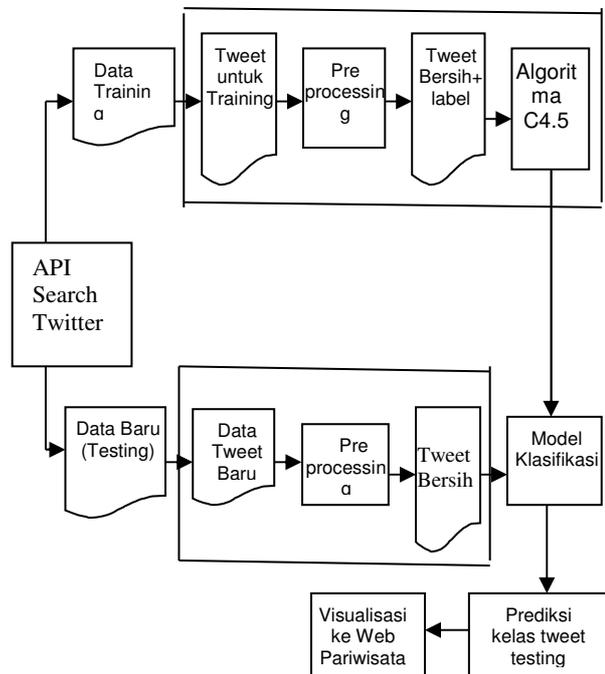


Gambar 1 Rancangan Arsitektur Sistem

2.2.1 Perancangan Sistem

Penelitian ini dimulai dengan mengambil data tweet yang tersimpan di server Twitter dengan menggunakan bantuan API Search Twitter. Pada proses pengambilan data ini juga dilakukan pe-label-an data. Kemudian data tersebut diproses untuk tahap persiapan (preprocessing) agar data siap digunakan untuk proses klasifikasi. Hal ini dilakukan karena tidak semua data tweet tersebut dapat digunakan. Pada data preprocessing ini dilakukan pembersihan data tweet dari tweet dan kata-kata yang tidak digunakan, penggantian kata-kata tertentu dengan daftar sinonim yang sudah ada, melakukan perbaikan data penting tetapi tidak sesuai (tidak lengkap). Kemudian hasil pembersihan data ini disimpan di tempat yang berbeda dari “data kotor”-nya.

Kemudian data tweet bersih yang sudah diberi label/kelas atau data training diolah oleh algoritma C4.5 untuk menghasilkan model pohon keputusan. Model keputusan tersebut terbagi menjadi dua bagian, yaitu: probabilitas kelas dan probabilitas semua kata (term) pada suatu kelas yang dalam sistem ini disimpan di database. Setelah itu, model probabilitas ini digunakan untuk memprediksi data tweet yang baru yang sudah dibersihkan pada tahap preprocessing. Kemudian hasil prediksi data tweet yang baru di visualisasikan ke Web Pariwisata yang terkoneksi dengan Web Service dengan menggunakan Google Map. Gambaran umum sistem dapat dilihat pada Gambar 2.



Gambar 2 Gambaran Umum Sistem

2.2.2 Perancangan Download Tweet

Proses download tweet diawali dengan mengakses API Search Twitter dengan query yang diberikan adalah hashtag cuaca. Hasil query yang diberikan API Search Twitter adalah file XML yang kemudian isi file XML tersebut di parsing ke variabel tertentu yang digunakan untuk proses penyimpanan data. Data yang dihasilkan oleh XML tersebut tidak sepenuhnya langsung disimpan, melainkan di cek terlebih dahulu apakah tweet yang akan dimasukkan ke database tersebut sudah ada sebelumnya dengan mengecek kesamaan antara user dan tweet dengan yang sudah tersimpan di database. Apabila sudah ada, maka data tidak akan disimpan tetapi apabila data belum ada maka selanjutnya data disimpan di database.

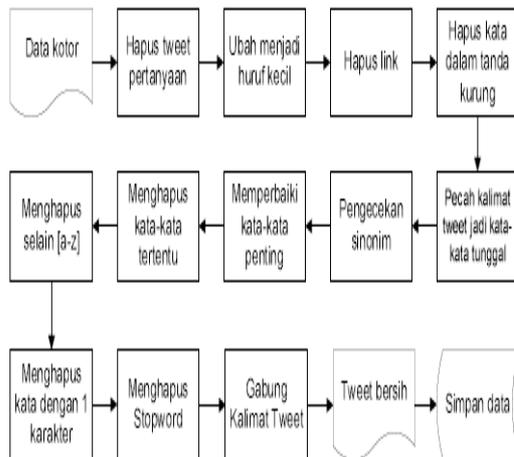
2.2.3 Perancangan Preprocessing

Secara umum preprocessing ini melewati beberapa tahap antara lain :

1. Pembersihan tweet pertanyaan. Dalam sistem ini tweet pertanyaan diasumsikan dengan tweet tersebut mengandung karakter “?”. Tweet pertanyaan dihapus diasumsikan bahwa tweet pertanyaan tidak mengandung informasi cuaca.

2. Pembersihan tweet kata per kata.
3. Penyimpanan data tweet bersih ke database.

Pada tahap preprocessing dilakukan pembersihan semua kata yang terkandung pada satu buah tweet. Proses yang dilewati pada preprocessing inidapat dilihat pada gambar 3. Adapun tahap pembersihan kata dalam sistem iniantara lain adalah mengkonversi menjadi huruf kecil untuk menstandarkan data,melakukan perbaikan kata-kata yang diperlukan dalam dalam klasifikasi,menghapus url/link yang merupakan data yang tidak diperlukan, menghapus mention (@userTwitter) karena data ini tidak diperlukan, menghapus karakter selain a-z dan 0-9, mengganti sinonim, menghilangkan stopwords dan menghapus kata yang hanya satu karakter.



Gambar 3 Flow Chart Preprocessing

2.2.4 Perancangan Bobot Term Teknik TF-IDF

Sebelum tweet diklasifikasi menggunakan algoritma C4.5,Pada bagian ini akan dijelaskan terlebih dahulu mengenai pembentukan bobot tiap term dari data tweet yang terkumpul dengan menggunakan teknik tf-idf. Misalkan, terdapat 4 buah tweet yang sudah bersih yang terdapat pada Gambar 4.

1. imogiri hujan angin kencang
2. panas malioboro
3. concat panas tenanan
4. panas imogiri

Gambar 4 Contoh Data Tweet

Dari data tweet pada Gambar 4 tersebut kemudian dipilih kata-kata unik dari semua tweet tersebut sehingga menjadi matrik yang tertuang pada Tabel 1.

Tabel 1 Matriks Term Data Tweet Gambar 4

Tweet	Imogiri	hujan	Angin	kencang	Panas	Malioboro	concat	tenanan
1	1	1	1	1				
2					1	1		
3					1		1	1
4	1				1			

Dari matrik tersebut kemudian dihitung dengan menggunakan persamaan (1) sehingga setiap kemunculan kata dalam tweet tersebut diubah menjadi nilai tf-idf. Berikut adalah contoh perhitungan tf-idf untuk kata “imogiri” pada tweet1 :

Jumlah dokumen = 4, tf (imogiri pada tweet1) = 1, df (imogiri) = 2

$$tf-idf (\text{Tweet 1, "imogiri"}) = 1 \times \log \frac{4}{2} = 0,301$$

Apabila proses perhitungan tersebut dilakukan untuk semua dokumen dan semua kata maka akan dihasilkan matrik perhitungan tf-idf yang terdapat pada tabel 2.

Tabel 2 Matriks tf-idf Data Tweet Gambar 4

Tweet	Imogiri	hujan	Angin	kencang	Panas	Malioboro	concat	tenanan
1	0,301	0,602	0,602	0,602	0	0	0	0
2	0	0	0	0	0,124	1	0	0
3	0	0	0	0	0,124	0	0,602	0,602
4	0,301	0	0	0	0,124	0	0	0

2.2.5. Perancangan Klasifikasi Algoritma C4.5

Setelah menemukan hasil perhitungan tf-idf tweet, selanjutnya hasil perhitungan tf-idf ini akan diolah dengan menggunakan algoritma C4.5 sehingga menghasilkan model pohon keputusan. Misalkan tweet 1 merupakan tweet dengan kelas hujan, sedangkan tweet 2,tweet 3 dan tweet 4 merupakan tweet dengan kelas panas. Sehingga matrik tf-idf tersebut akan berubah menjadi seperti pada Tabel 3.

Tabel 3 Matriks tf-idf Berdasarkan Kelas

Tweet	Imogiri	hujan	Angin	kencang	Panas	Malioboro	concat	tenanan	Kelas
1	0,301	0,602	0,602	0,602	0	0	0	0	Hujan
2	0	0	0	0	0,124	1	0	0	Panas
3	0	0	0	0	0,124	0	0,602	0,602	Panas
4	0,301	0	0	0	0,124	0	0	0	Panas

Node		Jumlah kasus (S)	Panas (S ₁)	Hujan (S ₂)	Entropy	Gain
1	Total	14	10	4	0,8631	
Daerah	imogiri	4	4	0		0,2585
	malioboro	5	4	1	0,7219	
	concat	5	2	3	0,9709	

Baris **Total** kolom Entropy pada Tabel 4

$$\text{Entropy (Total)} = \left(-\frac{4}{14} \times \log_2\left(\frac{4}{14}\right)\right) + \left(-\frac{10}{14} \times \log_2\left(\frac{10}{14}\right)\right)$$

$$\text{Entropy (Total)} = 0,8631$$

Perhitungan tersebut juga dilakukan untuk mencari entropy pada tiap daerah. Sementara itu, nilai Gain pada baris imogiri dihitung dengan menggunakan persamaan 2, sebagai berikut :

$$\text{Gain (Total, Daerah)} = \text{Entropy (Total)} - \sum_{i=1}^n \frac{|Daerah_i|}{|Total|} \times \text{Entropy (Daerah}_i)$$

$$\text{Gain (Total, Daerah)} = 0,8631 - \left(\frac{4}{14} \times 0\right) + \left(\frac{5}{14} \times 0,7219\right) + \left(\frac{5}{14} \times 0,9709\right)$$

$$\text{Gain (Total, Daerah)} = 0,2585$$

2.3 Hasil Pengujian Akurasi Klasifikasi

Sesuai dengan metode penelitian yang disampaikan, maka pengujian yang akan dilakukan pada sistem ini adalah dengan menggunakan metode k fold cross validation. Nilai k dalam penelitian ini adalah 10. Dengan demikian, akan dilakukan 10 kali pengujian akurasi klasifikasi dari data bersih yang dijadikan sampel. Adapun aspek yang diuji adalah akurasi model terhadap data penelitian yang dikumpulkan dengan berbagai porsi jumlah data penelitian yang dijadikan sampel. Secara lebih rinci, berikut merupakan mekanisme pengujian sistem yang dilakukan pada sistem ini :

1. Membagi sampel menjadi 10 bagian yang sama rata.
2. Sebanyak 10% dari jumlah sampel tersebut secara bergantian dijadikan sebagai data testing dan 90% lainnya dijadikan sebagai data training.
3. Dari 10% yang dijadikan data testing tersebut kemudian dibandingkan hasil klasifikasi oleh sistem dengan kelas yang sudah ditentukan sebelumnya.
4. Dihitung nilai akurasinya menggunakan persamaan (4) untuk masing k.

5. Dihitung nilai rata-rata seluruh nilai akurasi untuk semua k untuk memperoleh nilai akurasi keseluruhan.

Sistem yang sudah dibangun diuji dengan menggunakan data penelitian yang sudah terkumpul sebelumnya. Pada proses pengujian ini, dilakukan pengklasifikasian data dengan berbagai porsi jumlah sampel yang dijadikan sebagai data input dari sistem ini. Adapun jumlah porsi sampel adalah 100, 1000, 5000, 10000 dan 15106 sampel data yang digunakan.

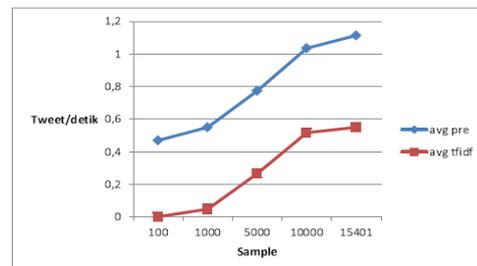
2.3.1 Pengujian Running Time Preprocessing

Pada pengujian ini dilakukan perhitungan running time preprocessing dari berbagai porsi jumlah sampel yang telah ditentukan sebelumnya. Tabel 5 tertuang hasil pengujian running time preprocessing.

Tabel 5 Running Time Preprocessing

No	Jumlah Sampel	Preprocessing (detik)	Rata-rata Preprocessing tweet/detik	Tf-idf (detik)	Rata-rata Tf-idf tweet/detik
1	100	41	0,410	0	0,000
2	1000	530	0,530	42	0,042
3	5000	3810	0,762	1319	0,263
4	10000	10337	1,033	5110	0,511
5	15106	16585	1,097	8385	0,555

Dari Table 5 dapat dilihat bahwa semakin besar jumlah sampel yang digunakan, semakin besar juga rata-rata tweet yang dapat diolah pada proses preprocessing dan tf-idf dalam satu detik. Gambar 5 merupakan grafik perbandingan running time dengan berbagai porsi jumlah sampel.



Gambar 5 grafik perbandingan running time

Dengan memperhatikan gambar 5 dapat dilihat bahwa semakin besar jumlah sampel yang diolah, maka semakin besar pula rata-rata tweet yang melalui preprocessing dalam satu

detik. Sementara itu, tidak terjadi kenaikan yang signifikan pada nilai rata-rata tweet yang yang dapat diproses tf-idf dalam satu detik.

2.3.2 Pengujian Akurasi Model

Pada pengujian ini dilakukan pengujian akurasi model pohon keputusan terhadap data bersih yang terbentuk dengan berbagai porsi jumlah. Seperti pada preprocessing, porsi jumlah data pada pengujian akurasi model pohon keputusan terdiri dari 100, 1000, 5000, 10000 dan 15106 sampel data yang digunakan yang dihasilkan dari preprocessing. Berikut merupakan rincian hasil pengujian dengan berbagai porsi jumlah data yang ditunjukkan pada Tabel 6.

Tabel 6. Rata-rata Nilai Akurasi Model

Jumlah Sampel	Rata-rata Akurasi
100	71 %
1000	83,40 %
5000	88,52 %
10000	91,58 %
15106	95,58 %

Tabel 6. merupakan tabel rata-rata nilai akurasi model untuk masing-masing pengujian dengan berbagai porsi data. Dari tabel tersebut dapat diketahui bahwa nilai akurasi tertinggi terdapat pada pengujian dengan menggunakan sampel data sebanyak 15106.

3. KESIMPULAN

Dari penelitian yang telah dilakukan terdapat beberapa kesimpulan yaitu :

1. Metode algoritma C4.5 mengasumsikan bahwa didalam setiap kelas dikelompokkan dengan pasti dan model pohon keputusan lebih mengarah pada perhitungan probabilitas tiap-tiap record, dalam hal ini kata. setiap kata unik bebas bersyarat satu sama lain. Asumsi ini digunakan agar proses training untuk membuat model klasifikasi dapat dijalankan.
2. Dari hasil pengujian preprocessing yang dilakukan, dapat disimpulkan bahwa semakin banyak jumlah data yang di-preprocessing maka semakin lama waktu yang diperlukan untuk preprocessing. Hal ini terlihat dari semakin besar rata-rata tweet yang dapat di-preprocessing dalam satu detik.
3. Dari hasil pengujian akurasi model dari sistem yang dikembangkan, menghasilkan

nilai akurasi terkecil sebesar 71% pada proses pengujian dengan menggunakan sampel sebanyak 100 dan menghasilkan nilai akurasi tertinggi sebesar 95,58% pada proses pengujian dengan menggunakan sampel sebanyak 15106.

SARAN

penelitian ini masih memiliki keterbatasan yang dapat dijadikan acuan untuk pengembangan dimasa yang akan datang, sehingga dapat disarankan beberapa hal sebagai berikut :

1. Perlunya parameter untuk mengidentifikasi twitter yang di-input-kan adalah user sebenarnya (manusia) atau mesin
2. Perlu dilakukan analisis hubungan antar kata pada setiap klasifikasi
3. Perlu dilakukan perbandingan dengan teknik lain pada proses *preprocessing*
4. Perlu dilakukan validasi apakah hasil klasifikasi ini sesuai dengan keadaan di tempat sebenarnya.

DAFTAR PUSTAKA

- Hepburn, A., 2010, Infographic: Twitter Statistics, Facts & Figures, <http://www.digitalbuzzblog.com/infographic-twitterstatistics-factsfigures/> Di akses tanggal 9 Mei 2011.
- McCormick, H. T., Lee, H., Cesare, N., Shojaie, A., 2013, Using Twitter for Demographic and Social Science Research: Tools for Data Collection, International Joint Conference for Statistics and the Social Sciences University of Washington, Paper no.127.
- Han, J., & Kamber, M., *Data Mining Concept and Technique*, San Fransisco: Morgan Kaufman Publisher, 2006
- Tan, P. N., Steinbach, M., & Kumar, V., 2006, *Introduction to Data Mining*, Pearson Education, Boston.
- Witten, I. H., 2005, Text mining, Dalam *Practical Handbook of Internet Computing Florida*: Chapman & Hall/CRC Press, Boca Raton, 14-22.
- Lonneke Mous. Predicting bankruptcy with discriminant analysis and decision tree using financial ratios, 2005. Faculty of Economics at Erasmus University Rotterdam.

Kusrini & Luthfi, T. E., 2009, Algoritma Data Mining, Andi Offset, Yogyakarta.