

**ANALISIS KECENDERUNGAN INFORMASI DENGAN
MENGUNAKAN METODE *TEXT MINING***

(Studi Kasus: Akun *twitter* @detikcom)



SKRIPSI

Oleh:

SYAIFUDIN KARYADI

NIM. 24010212130030

**DEPARTEMEN STATISTIKA
FAKULTAS SAINS DAN MATEMATIKA
UNIVERSITAS DIPONEGORO
SEMARANG**

2016

**ANALISIS KECENDERUNGAN INFORMASI DENGAN
MENGUNAKAN METODE *TEXT MINING*
(Studi Kasus: Akun *twitter @detikcom*)**

**Oleh:
Syaifudin Karyadi
24010212130030**

Tugas Akhir sebagai salah satu syarat untuk memperoleh
gelar Sarjana Sains pada Departemen Statistika

**DEPARTEMEN STATISTIKA
FAKULTAS SAINS DAN MATEMATIKA
UNIVERSITAS DIPONEGORO
SEMARANG
2016**

HALAMAN PENGESAHAN I

Judul Skripsi : Analisis Kecenderungan Informasi dengan Menggunakan
Metode *Text Mining*
(Studi Kasus: Akun *twitter* @detikcom)

Nama : Syaifudin Karyadi

NIM : 24010212130030

Departemen : Statistika

Telah diujikan pada sidang Tugas Akhir dan dinyatakan lulus pada tanggal 16
Agustus 2016

Semarang, 16 Agustus 2016

Mengetahui,
Ketua Departemen Statistika
Fakultas Sains dan Matematika Undip

Panitia Penguji Ujian Tugas Akhir
Ketua,

Dra. Dwi Ispriyanti, M.Si.
NIP. 195709141986032001

Dra. Suparti, M.Si
NIP. 196509131990032001

HALAMAN PENGESAHAN II

Judul Skripsi : Analisis Kecenderungan Informasi dengan Menggunakan
Metode *Text Mining*

(Studi Kasus: Akun *twitter @detikcom*)

Nama : Syaifudin Karyadi

NIM : 24010212130030

Departemen : Statistika

Telah diujikan pada sidang Tugas Akhir dan dinyatakan lulus pada tanggal 16
Agustus 2016

Semarang, 16 Agustus 2016

Dosen Pembimbing I

Dosen Pembimbing II

Hasbi Yasin, S.Si, M.Si
NIP. 198212172006041003

Moch. Abdul Mukid, S.Si, M.Si
NIP. 197808172005011001

KATA PENGANTAR

Puji Syukur penulis ucapkan kehadirat Allah SWT yang telah memberikan rahmat dan karunia-Nya sehingga penulis dapat menyelesaikan penulisan Tugas Akhir dengan judul “Analisis Kecenderungan Informasi dengan Menggunakan Metode *Text Mining*”.

Penulis menyadari bahwa dalam penulisan Tugas Akhir ini tidak lepas dari bimbingan dan dukungan yang diberikan beberapa pihak. Oleh karena itu, penulis ingin menyampaikan terima kasih kepada:

1. Ibu Dra. Dwi Ispriyanti, M.Si. sebagai Ketua Departemen Statistika Fakultas Sains dan Matematika Universitas Diponegoro.
2. Bapak Hasbi Yasin, S.Si., M.Si. selaku dosen pembimbing I dan Bapak Moch. Abdul Mukid, S.Si, M.Si. selaku dosen pembimbing II.
3. Bapak dan Ibu dosen Departemen Statistika Fakultas Sains dan Matematika Universitas Diponegoro
4. Semua pihak yang tidak dapat disebutkan satu per satu yang telah membantu penulis dalam penulisan Tugas Akhir ini.

Penulis menyadari bahwa penulisan Tugas Akhir ini masih jauh dari sempurna. Oleh karena itu, penulis mengharapkan kritik dan saran yang membangun demi kesempurnaan penulisan selanjutnya.

Semarang, 16 Agustus 2016

Penulis

ABSTRAK

Internet merupakan suatu fenomena yang luar biasa. Berawal dari sebuah eksperimen militer di Amerika Serikat, internet telah berkembang menjadi 'kebutuhan' bagi lebih dari puluhan juta orang di seluruh dunia. Jumlah pengguna internet yang besar dan semakin berkembang, telah mewujudkan budaya internet. Salah satu yang berkembang pesat yaitu media sosial *twitter*. *Twitter* merupakan layanan *microblogging* yang menyimpan *text database* yang disebut *tweet*. Untuk memudahkan memperoleh informasi yang dominan dibicarakan, maka dicarilah topik dari *tweet twitter* dengan menggunakan *clustering*. Pada penelitian ini, dilakukan pengelompokan 500 *tweet* dari akun *twitter @detikcom* menggunakan *k-means clustering*. Hasil dari penelitian ini menunjukkan bahwa *Dunn index* yang maksimum, pengelompokan terbaik *k-means Clustering* untuk memperoleh topik yang dominan yaitu sebanyak tiga *cluster*, yaitu mengenai pemerintah, Jakarta, dan politik.

Kata Kunci: *text mining, clustering, k-means , dunn index, dan twitter*

ABSTRACT

The Internet is an extraordinary phenomenon. Starting from a military experiment in the United States, the Internet has evolved into a 'need' for more than tens of millions of people worldwide. The number of internet users is large and growing, has been creating internet culture. One of the fast growing social media twitter. Twitter is a microblogging service that stores text database called tweets. To make it easier to obtain information that is dominant discussed, then sought the topic of twitter tweet using clustering. In this research, grouping 500 tweets from twitter account @detikcom using k-means clustering. The results of this study indicate that the maximum index Dunn, the best grouping K-means clustering to obtain the dominant topic as many as three clusters, namely the government, Jakarta, and politics.

Keywords: text mining, clustering,, k-means , dunn index, and twitter.

DAFTAR ISI

	Halaman
HALAMAN JUDUL	i
HALAMAN PENGESAHAN I	ii
HALAMAN PENGESAHAN II	iii
KATA PENGANTAR	iv
ABSTRAK	v
ABSTRACT	vi
DAFTAR ISI	vii
DAFTAR TABEL	x
DAFTAR GAMBAR	xi
DAFTAR LAMPIRAN	xii
 BAB I PENDAHULUAN	
1.1 Latar Belakang	1
1.2 Rumusan Masalah	3
1.3 Batasan Masalah	4
1.4 Tujuan Penelitian	4
 BAB II TINJAUAN PUSTAKA	
2.1 <i>Twitter</i>	5
2.2 <i>Data Mining dan Text Mining</i>	6
2.2.1 <i>Term-document Matrix</i>	9
2.2.2 Pembobotan	10
2.3 Fitur	11
2.3.1 Tipe Fitur	11
2.3.2 Konsep Kedekatan	12
2.3.3 Ukuran Kedekatan Kontinyu	12
2.4 <i>Clustering</i>	13

2.5	Validasi <i>Cluster</i>	16
-----	-------------------------------	----

BAB III METODOLOGI PENELITIAN

3.1	Sumber Data	19
3.2	Metode Pengumpulan Data	19
3.3	Metode Analisis	19
3.4	Diagram Alir Analisis	21

BAB IV HASIL DAN PEMBAHASAN

4.1	Profil Akun @detikcom	22
4.2	<i>Application Programming Interface</i> (API)	23
4.3	<i>Term-document Matrix</i> dari 5 <i>tweet</i> @detikcom	24
4.3.1	<i>Text Pre-Process</i>	25
4.3.1.1	<i>To Lower Case</i>	25
4.3.1.2	<i>Tokenizing</i>	25
4.3.1.3	<i>Remove Number</i>	26
4.3.1.4	<i>Remove URL</i>	27
4.3.1.5	<i>Remove Punctuation</i>	27
4.3.2	<i>Feature Selection</i>	28
4.3.3	<i>Frequent Terms</i> dari 5 <i>tweet</i> @detikcom	31
4.3.4	<i>Wordcloud</i> dari 5 <i>tweet</i> @detikcom	33
4.3.5	Validasi <i>Cluster</i> dari 5 <i>tweet</i> @detikcom	34
4.4	<i>Term-document Matrix</i> dari 500 <i>tweet</i> @detikcom	37
4.4.1	<i>Frequent Terms</i> dari 500 <i>tweet</i> @detikcom	38
4.4.2	<i>Wordcloud</i> dari 500 <i>tweet</i> @detikcom	39
4.5	<i>K-Means Clustering</i>	40
4.6	Validasi <i>Cluster</i> dari 500 <i>tweet</i> @detikcom	43

BAB V PENUTUP

5.1	Kesimpulan	45
5.2	Saran	45

DAFTAR PUSTAKA	47
LAMPIRAN	50

DAFTAR TABEL

	Halaman
Tabel 1 <i>Term-document Matrix</i>	10
Tabel 2 Tipe Fitur	11
Tabel 3 <i>Term-document Matrix</i> dengan pembobotan tf untuk 5 <i>tweet</i>	29
Tabel 4 <i>Term-document Matrix</i> dengan pembobotan TF-IDF untuk 5 <i>tweet</i> ..	30
Tabel 5 Jumlah kemunculan seluruh <i>terms</i> pada masing-masing dokumen...	31
Tabel 6 <i>Output Dunn index K-means Clustering 5 Tweet</i> dari akun @detikcom	34
Tabel 7 Keanggotaan 2 <i>cluster 5 tweet</i> dari akun @detikcom dengan <i>K-means Clustering</i>	35
Tabel 8 Perhitungan jarak antar data untuk <i>cluster 1</i>	35
Tabel 9 Perhitungan jarak <i>cluster 1</i> dengan <i>cluster 2</i>	36
Tabel 10 <i>Term-document Matrix</i> dengan pembobotan tf untuk 500 <i>tweet</i>	37
Tabel 11 <i>Term-document Matrix</i> dengan pembobotan TF-IDF untuk 500 <i>tweet</i>	38
Tabel 12 Keanggotaan 3 <i>cluster 500 tweet</i> dari akun @detikcom dengan <i>K-means Clustering</i>	43
Tabel 13 <i>Output Dunn index K-means Clustering 500 Tweet</i> dari akun @detikcom	43

DAFTAR GAMBAR

	Halaman
Gambar 1 Diagram Alir Analisis	21
Gambar 2 Tampilan akun <i>twitter @detikcom</i>	22
Gambar 3 Tampilan API.....	24
Gambar 4 Ilustrasi dari <i>tokenizing</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	26
Gambar 5 Ilustrasi dari <i>remove number</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	26
Gambar 6 Ilustrasi dari <i>remove url</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	27
Gambar 7 Ilustrasi dari <i>remove punctuation</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	27
Gambar 8 Ilustrasi dari <i>stopword</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	28
Gambar 9 Diagram Batang Kemunculan <i>Term</i> untuk 5 <i>tweet</i> dari akun <i>twitter @detikcom</i>	32
Gambar 10 Wordcloud 5 <i>tweet</i> dari akun <i>@detikcom</i>	33
Gambar 11 Diagram Batang Kemunculan <i>Term</i> dari 500 <i>tweet @detikcom</i> (Frekuensi ≥ 6)	38
Gambar 12 Wordcloud 500 <i>tweet</i> dari akun <i>@detikcom</i>	40
Gambar 13 <i>Network Graph K-Means Clustering</i> untuk <i>Cluster 1</i>	42

DAFTAR LAMPIRAN

	Halaman
Lampiran 1 <i>Syntax software R untuk retrieve data 5 tweet media sosial twitter dari akun twitter @detikcom</i>	50
Lampiran 2 <i>Syntax software R untuk membuat Term-document Matrix dengan Pembobotan TF-IDF, wordcloud, validasi cluster, dan k-means clustering dari 5 tweet media sosial twitter dari akun twitter @detikcom</i>	55
Lampiran 3 <i>Syntax software R untuk retrieve data 500 tweet media sosial twitter dari akun twitter @detikcom</i>	57
Lampiran 4 <i>Syntax software R untuk membuat Term-document Matrix dengan Pembobotan TF-IDF, wordcloud, validasi cluster, dan k-means clustering dari 500 tweet media sosial twitter dari akun twitter @detikcom</i>	61

BAB I

PENDAHULUAN

1.1 Latar Belakang

Internet merupakan suatu fenomena yang luar biasa. Berawal dari sebuah eksperimen militer di Amerika Serikat, internet telah berkembang menjadi 'kebutuhan' bagi lebih dari puluhan juta orang di seluruh dunia. Jumlah pengguna internet yang besar dan semakin berkembang, telah mewujudkan budaya internet.

Menurut Asosiasi Penyelenggara Jasa Internet Indonesia (APJII) (2014), pengguna internet di Indonesia selalu bertambah dari tahun ke tahun. Jumlah pengguna internet di Indonesia mencapai 88 juta orang hingga akhir tahun 2014 atau mengalami kenaikan sebesar 34,9% jika dibandingkan dengan tahun 2013.

Hasil riset tahun 2014 secara signifikan menunjukkan pengguna jejaring sosial (sosial media) menduduki peringkat tertinggi yang dimanfaatkan, mengalahkan pencarian informasi (*browsing/searching*) di posisi kedua. Posisi ke-3 *chatting* (*messaging*), pencarian berita (ke-4), video (ke-5), *email* (ke-6). Pencarian berita dan penggunaan *email* saat ini anjlok tak populer (APJII, 2014).

Menurut Kominfo (2013), di era globalisasi perkembangan telekomunikasi dan informatika (IT) sudah begitu pesat. Teknologi membuat jarak tak lagi jadi masalah dalam berkomunikasi. Internet tentu saja menjadi salah satu medianya. Situs jejaring sosial yang paling banyak diakses adalah *facebook* dan *twitter*. Indonesia menempati peringkat 4 pengguna *facebook* terbesar setelah USA, Brazil, dan India.

Sedangkan, untuk *twitter* Indonesia menempati peringkat 5 pengguna *twitter* terbesar di dunia. Posisi Indonesia hanya kalah dari USA, Brazil, Jepang dan Inggris. Pengguna *twitter* di Indonesia berdasarkan data PT Bakrie Telecom sebesar 19,5 juta pengguna dari total 500 juta pengguna global. *Twitter* menjadi salah satu jejaring sosial paling besar di dunia sehingga mampu meraup keuntungan mencapai USD 145 juta.

Menurut Francis dan Flynn (2010), *text mining* adalah teknologi baru yang digunakan untuk data perusahaan yang selalu bertambah sehingga data teks yang tidak terstruktur tersebut dapat dianalisis. Salah satu inovasi *software* yang dapat meringankan biaya bagi penambang teks adalah *software* yang bersifat *open source*. Dua jenis *software open source* yang sangat populer dan diunggulkan adalah R dan Perl. R adalah bahasa pemrograman yang mendukung hal-hal yang berkaitan dengan statistik dan digunakan pada hal-hal yang berhubungan dengan ilmu pasti, matematis.

Menurut Zhao (2012), metode *text mining* telah digunakan untuk menganalisa data pada *twitter*. Metode ini dimulai dengan mengambil *text* yang ada pada *twitter*, *text* yang sudah diambil kemudian diubah menjadi *document-term matrix*. Setelah itu, *frequent words* dan *assosiation* yang diperoleh dari *matrix*. *Wordcloud* digunakan untuk menunjukkan kata-kata penting yang ada pada dokumen. Terakhir untuk mendapatkan topik dari *tweet*, kata-kata dalam *tweet* atau biasa disebut *term* akan dikelompokkan dengan metode *k-means cluster*.

R adalah salah satu *software open source* untuk komputasi statistik dan grafik. R menyediakan berbagai varian metode statistik dan grafik. R dapat dipermudah dengan adanya *packages*. Berdasarkan pada CRAN (2016) terdapat 8042 *packages*

yang tersedia pada CRAN *packages repository* per 5 Maret 2016. Untuk melakukan analisa dengan metode *text mining* pada sebuah akun *twitter* dibutuhkan beberapa packages, seperti *packages twitter* dan *tm* diperlukan untuk membantu mendapatkan data pada akun tersebut serta menjelmakan teks. Ada juga *packages word cloud* yang digunakan untuk merepresentasikan visual untuk data teks, biasanya untuk menggambarkan metadata kata kunci (*tag*) di situs *web*. *Tags* biasanya satu kata, dan pentingnya setiap *tag* ditunjukkan dengan ukuran *font* atau warna (Zhao, 2012).

Beberapa informasi penting yang dapat diperoleh dari *twitter* antara lain seperti melihat sejarah perkembangan manusia, sejarah obama terpilih menjadi presiden, dll. Tersedia dalam *tweet-tweet* yang bisa dirunut di *twitter*. Penelitian ini dilakukan pengelompokkan 500 *tweet* dari akun *twitter* @detikcom menggunakan metode *k-means clustering* yang bertujuan untuk untuk mengetahui kecenderungan topik pemberitaan dan mengetahui topik yang paling sering muncul. Hasil analisis pada akun *twitter* berita tersebut akan memberikan gambaran pemberitaan akhir-akhir ini. Penelitian ini menjadi penting mengingat akun @detikcom merupakan akun berita *online* dengan *followers* terbanyak, sehingga berita yang disampaikan juga akan mempengaruhi pengetahuan dan persepsi publik terhadap suatu masalah.

Berdasarkan uraian diatas maka peneliti tertarik untuk menganalisa kecenderungan topik informasi pemberitaan yang disampaikan melalui akun *twitter* @detikcom dengan menggunakan metode *text mining*.

1.2 Rumusan Masalah

Berdasarkan uraian latar belakang dapat dirumuskan permasalahan sebagai berikut:

1. Bagaimana kecenderungan topik informasi yang disampaikan melalui akun *twitter* @detikcom?
2. *Cluster tweet* apa saja yang terbentuk dari akun *twitter* @detikcom?

1.3 Batasan Masalah

Dalam penelitian ini, masalah dibatasi hanya pada 500 *tweets* teratas yang diambil dari *timeline* akun *twitter* @detikcom pada hari Jum'at, 3 Juni 2016 jam 18.30 WIB.

1.4 Tujuan Penelitian

Berdasarkan rumusan masalah, maka tujuan yang ingin dicapai dalam penelitian ini adalah sebagai berikut :

1. Untuk mengetahui kecenderungan topik informasi yang disampaikan melalui akun *twitter* @detikcom?
2. Untuk mengetahui *cluster tweet* apa saja yang terbentuk dari akun *twitter* @detikcom.