

**SISTEM ANALISIS SENTIMEN PUBLIK TENTANG OPINI
PEMILIHAN KEPALA DAERAH JAWA TIMUR 2018 PADA
DOKUMEN TWITTER MENGGUNAKAN *NAIVE BAYES*
*CLASSIFIER***

SKRIPSI



**OLEH
FAJAR DARWIS DZIKRIL HAKIMI
NIM. H02214001**

**PROGRAM STUDI MATEMATIKA
JURUSAN SAINS
FAKULTAS SAINS DAN TEKNOLOGI
UNIVERSITAS ISLAM NEGERI SUNAN AMPEL SURABAYA
SURABAYA
2018**

PERNYATAAN KEASLIAN

Yang bertandatangan di bawah ini:

Nama : Fajar Darwis Dzikril Hakimi

NIM : H02214001

Program Studi : Matematika

Angkatan : 2014

Menyatakan bahwa saya tidak melakukan plagiat dalam bentuk apapun pada skripsi saya yang berjudul: **Sistem Analisis Sentimen Publik Tentang Opini Pemilihan Kepala Daerah Jawa Timur 2018 Pada Dokumen Twitter Menggunakan *Naive Bayes Classifier***. Apabila suatu saat nanti terbukti saya melakukan tindakan plagiat, maka saya akan menerima sanksi yang telah ditetapkan.

Demikian pernyataan keaslian ini saya buat dengan sebenar-benarnya

Surabaya, 3 Agustus 2018



Fajar Darwis Dzikril Hakimi
NIM. H02214001

LEMBAR PENGESAHAN

**SISTEM ANALISIS SENTIMEN PUBLIK TENTANG OPINI PEMILIHAN
KEPALA DAERAH JAWA TIMUR 2018 PADA DOKUMEN TWITTER
MENGUNAKAN *NAIVE BAYES CLASSIFIER***

Disusun oleh
Fajar Darwis Dzikril Hakimi
NIM. H02214001

Telah dipertahankan di depan Dewan Penguji
Pada tanggal 18 Juli 2018
Dan dinyatakan telah memenuhi syarat
untuk memperoleh gelar
Sarjana Matematika (S.Mat)

Dewan Penguji

Penguji I



Ahmad Hanif, Asyhar, M.Si
NIP. 198601232014031001

Penguji II



Nurissadah Ulinnuha, M. Kom
NIP. 199011022014032004

Penguji III



Aris Fanani, M. Kom
NIP. 198701272014031002

Penguji IV



Wika Dianita Utami, M.Sc
NIP. 1992061002018012003

Mengesahkan
Dekan Fakultas Sains dan Teknologi
UIN Sunan Ampel Surabaya



Dr. Eni Purwati, M.Ag
NIP. 196512211990022001



KEMENTERIAN AGAMA
UNIVERSITAS ISLAM NEGERI SUNAN AMPEL SURABAYA
PERPUSTAKAAN

Jl. Jend. A. Yani 117 Surabaya 60237 Telp. 031-8431972 Fax.031-8413300
E-Mail: perpus@uinsby.ac.id

LEMBAR PERNYATAAN PERSETUJUAN PUBLIKASI
KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademika UIN Sunan Ampel Surabaya, yang bertanda tangan di bawah ini, saya:

Nama : Fajar Darwis Dzikril Hakimi
NIM : H02214001
Fakultas/Jurusan : Saintek / Matematika
E-mail address : fajarddh@gmail.com

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Perpustakaan UIN Sunan Ampel Surabaya, Hak Bebas Royalti Non-Eksklusif atas karya ilmiah :

Skripsi Tesis Desertasi Lain-lain (.....)

yang berjudul :

Sistem Analisis Sentimen Publik tentang Opini Pemilihan Kepala Daerah
Jawa Timur 2018 pada Dokumen Twitter menggunakan Naive
Bayes Classifier

beserta perangkat yang diperlukan (bila ada). Dengan Hak Bebas Royalti Non-Eksklusif ini Perpustakaan UIN Sunan Ampel Surabaya berhak menyimpan, mengalih-media/format-kan, mengelolanya dalam bentuk pangkalan data (database), mendistribusikannya, dan menampilkan/mempublikasikannya di Internet atau media lain secara **fulltext** untuk kepentingan akademis tanpa perlu meminta ijin dari saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan atau penerbit yang bersangkutan.

Saya bersedia untuk menanggung secara pribadi, tanpa melibatkan pihak Perpustakaan UIN Sunan Ampel Surabaya, segala bentuk tuntutan hukum yang timbul atas pelanggaran Hak Cipta dalam karya ilmiah saya ini.

Demikian pernyataan ini yang saya buat dengan sebenarnya.

Surabaya, 1 Agustus 2018

Penulis

(Fajar Darwis D. H.)

dokumen twitter. Opini masyarakat tersebut bisa dijadikan peta kekuatan calon gubernur dan wakil gubernur dalam persiapan menghadapi pemilu. Oleh karena itu, dengan melakukan klasifikasi terhadap dokumen twitter maka peta kekuatan calon gubernur dan wakil gubernur akan diketahui sedini mungkin sehingga dapat dijadikan evaluasi terhadap kinerja kampanye masing-masing calon.

Masalah akan muncul ketika dilakukan analisis dan pengklasifikasian sentimen secara manual. Masalah tersebut yaitu dibutuhkan waktu yang lama untuk melakukan klasifikasi terhadap dokumen twitter tersebut dan hal tersebut tentu tidak efektif mengingat proses evaluasi terhadap kinerja kampanye harus dilakukan secepat dan seefektif mungkin. Untuk memecahkan permasalahan tersebut, akan digunakan solusi berupa sistem analisis sentimen positif dan negatif tentang opini pemilihan Kepala Daerah Jawa timur 2018 pada dokumen twitter.

Untuk mengklasifikasikan dokumen teks, algoritma yang sering digunakan dan mempunyai hasil yang cukup bagus adalah *naive bayes classifier* dan *support vector machine* (Aggarwal & Zhai, 2013). Selain itu, pada penelitian yang dilakukan oleh Ahmad Fathan Hidayatullah dan Azhari SN (Hidayatullah & SN, 2014) yang berjudul “Analisis Sentimen dan Klasifikasi Kategori Terhadap Tokoh Publik pada Twitter” digunakan data *tweet* sebanyak 1329 data yang diperoleh dari hasil *crawling* yang diberi label kategori secara manual untuk kemudian diklasifikasikan dengan metode *naïve bayes*. Proses yang pertama dilakukan adalah tahap pre-proses untuk membersihkan data *tweet* dan mempersiapkannya untuk proses klasifikasi. Proses dilanjutkan dengan menghitung probabilitas masing-masing kata dalam *tweet* berdasarkan data *training* menggunakan metode *term*

yang mereka pasarkan. Dalam pemerintahan, pemerintah tentu akan memperhatikan opini masyarakat tentang kebijakan-kebijakan yang sudah ditetapkan. Opini masyarakat ini tentu berguna untuk memperbaiki dan menyempurnakan produk atau kebijakan yang telah dibuat.

Dengan pertumbuhan media sosial yang pesat saat ini seperti twitter, facebook, blog, beberapa produsen dan penyedia jasa menggunakan opini di media sosial tersebut untuk pertimbangan dalam membuat kebijakan di masa depan (Liu, 2012).

Selain itu, opini atau sentimen di media sosial juga bisa digunakan sebagai pemetaan kekuatan calon kepala daerah. Dengan memperhatikan sentimen positif dan negatif yang ditujukan kepada calon kepala daerah tertentu, maka bisa dilihat kecenderungan masyarakat suatu daerah dalam menentukan pilihan dalam pemilihan kepala daerah selanjutnya.

B. Kategorisasi Sentimen Opini

Dalam penelitian sebelumnya yang dilakukan oleh Faradhillah (Faradhillah, 2016), kategorisasi sentimen dilakukan secara manual. Kategorisasi sentimen ini dilakukan agar data yang dimiliki mempunyai kategori sehingga bisa dilakukan proses pelatihan dan pengujian data. Dalam melakukan kategorisasi secara manual, data yang dimiliki dibagi menjadi dua jenis data yaitu data untuk calon gubernur nomor urut satu dan data untuk calon gubernur nomor urut dua. Pembagian data ini dilakukan agar tidak terjadi bias dalam proses kategorisasi opini. Jadi, ketika data itu sudah terbagi dan masuk dalam data untuk salah satu calon, maka tidak akan

facebook adalah fitur *follow* yang memudahkan seseorang untuk terhubung dengan orang lain dengan jumlah yang tidak terbatas. Hal ini tentu lebih baik daripada facebook yang membatasi pertemanan hanya sejumlah 5000 orang.

Dalam twitter, ada beberapa istilah dasar yang harus dipahami agar seseorang dapat menggunakan situs jejaring sosial twitter (Maulana, 2012) diantaranya:

1. *Tweet*

Tweet adalah apa yang seseorang tulis atau posting di dalam twitter (layaknya status facebook).

2. *Mention* (@)

Mention adalah suatu cara untuk membuat link terhadap suatu akun twitter. Cara ini biasanya digunakan ketika seseorang membalas *tweet* atau ingin menandai suatu *tweet* kepada seseorang. Ketika seseorang membuat *tweet* yang memuat mention akun twitter orang lain, maka *tweet* tersebut akan muncul di beranda orang lain tersebut.

3. *Hashtag* (#)

Cara menuliskan *hashtag* yaitu dengan cara menambahkan tanda # di depan topik yang ditulis. Penulisan topik di dalam hashtag tidak diperkenankan menggunakan spasi, jadi ketika suatu topik terdiri dari dua kata maka penulisannya digabung menjadi satu. *Hashtag* berfungsi agar pengguna lain yang melihat bisa mencari topik sejenis yang ditulis orang lain dengan mudah. Intinya untuk mengelompokkan suatu topik agar mudah dicari. *Hashtag* ini berhubungan dengan *trending topic*. Ketika suatu *hashtag* merupakan *hashtag* yang paling banyak diposting di

twitter, maka *hashtag* tersebut menjadi *trending topic*. Contoh: #pilkadajatim, #pemiluserentak

4. *ReTweet* (RT)

ReTweet dalam bahasa Indonesia adalah mengulang suatu *tweet* seseorang maka *retweet* memiliki makna bahwa seseorang mengulang *tweet* orang lain untuk disampaikan kembali kepada *follower* nya. Dalam penelitian ini, data yang digunakan adalah data dalam twitter. Proses pencarian data dalam penelitian ini menggunakan metode *crawling*. *Crawling* itu prinsipnya seperti *search engine* tetapi mampu mencari data dalam jumlah yang besar dan menyimpannya.

D. Pemilihan Kepala Daerah 2018

Tahun 2018 ini merupakan tahun diadakannya pemilihan kepala daerah secara serentak hampir di seluruh daerah di Indonesia. Pemilihan kepala daerah adalah pemilihan pemimpin baru daerah secara langsung oleh masyarakat daerah tersebut. Sedangkan menurut PP Nomor 6 tahun 2005, pemilihan kepala daerah adalah pemilihan kepala daerah dan wakil kepala daerah sebagai sarana pelaksanaan kedaulatan rakyat di wilayah provinsi atau kabupaten ataupun kota berdasarkan Pancasila dan Undang-Undang Dasar Negara Republik Indonesia tahun 1945 untuk memilih kepala daerah dan wakil kepala daerah.

Joko J. Prihatantoro menyatakan bahwa pemilihan kepala daerah adalah rekrutmen politik yaitu penyeleksian rakyat terhadap tokoh-tokoh yang mencalonkan diri sebagai kepala daerah, baik gubernur/wakil gubernur maupun bupati/wakil bupati atau walikota/wakil walikota. Dalam kehidupan politik di

daerah, pemilihan kepala daerah merupakan salah satu kegiatan yang nilainya ekuivalen dengan pemilihan anggota DPRD. Ekuivalen tersebut ditunjukkan dengan kedudukan yang sejajar antara kepala daerah dan DPRD.

E. Klasifikasi

Klasifikasi merupakan kata serapan dari bahasa Inggris, *classification*. Menurut Kamus Besar Bahasa Indonesia, klasifikasi adalah proses pengelompokan benda berdasarkan ciri-ciri persamaan dan perbedaan. Menurut Saputro, klasifikasi adalah suatu penilaian terhadap sebuah data untuk dimasukkan ke dalam kelas tertentu dari beberapa kelas yang tersedia (Saputro, 2016).

Dalam klasifikasi, ada dua proses utama yang dilakukan. Pertama yaitu pembentukan model klasifikasi berdasarkan data-data yang sudah ada dan dijadikan sebagai memori. Kedua yaitu pengkategorian data baru yang didasarkan pada model klasifikasi yang telah dibentuk tersebut. Ada banyak metode yang bisa digunakan untuk melakukan klasifikasi. Metode yang sering digunakan yaitu *decision tree*, *artificial neural network*, *support vector machine*, dan *naive bayes classifier*.

F. Klasifikasi Dokumen Teks

Klasifikasi merupakan pembentukan suatu model yang akan digunakan untuk memprediksi suatu kategori (Juniawan, 2009). Klasifikasi memiliki dua tahap proses yaitu proses pembelajaran dan proses klasifikasi itu sendiri. Pada tahap pertama, algoritma klasifikasi membentuk suatu model klasifikasi dengan menganalisis data latih. Tahap ini disebut juga sebagai *supervised learning* karena setiap data latih telah memiliki label kelas masing-masing. Tahap ini bisa

- $P(A \cap B)$ adalah peluang hipotesis A yang terjadi bersamaan dengan hipotesis B .
- $P(A)$ adalah peluang terjadinya hipotesis A tanpa melihat kondisi apapun.
- $P(B)$ adalah peluang terjadinya kejadian B tanpa melihat kondisi apapun.

Ide dasar dari teorema bayes adalah hasil dari peluang terjadinya hipotesis A dapat diperkirakan berdasarkan kejadian B yang diamati. Ada beberapa hal penting dari teorema bayes yang merupakan dasar dari *naïve bayes classifier* (Balagatabi, 2012) diantaranya:

1. Sebuah peluang awal atau $P(A)$ adalah peluang dari suatu hipotesis sebelum kejadian lain diamati.
2. Sebuah probabilitas akhir atau $P(A|B)$ adalah peluang dari suatu hipotesis setelah kejadian lain diamati.

H. Naïve Bayes Classifier

Metode *Naïve bayes* merupakan metode yang digunakan untuk memprediksi probabilitas (Rini, Farida, & Puspitasari, 2016). Metode ini memanfaatkan teori peluang yang dikemukakan oleh ilmuwan Inggris Thomas Bayes. Cara kerja metode ini yaitu dengan memprediksi peluang terjadinya kejadian di masa depan berdasarkan data yang ada sebelumnya. *Naïve bayes classifier* adalah konsep probabilitas yang bisa digunakan untuk penentuan kelompok kelas dokumen teks dan dapat mengolah data dalam jumlah besar dengan hasil akurasi yang tinggi (Lestari, Perdana, & Fauzi, 2017). Tingkat performa dari sistem klasifikasi yang dibuat menggunakan *naive bayes classifier* bergantung pada data yang dimiliki dan data yang dipilih sebagai data *training*. Jika data yang dipilih sebagai data *training*

bisa mewakili semua atau sebagian besar data yang dimiliki, maka sistem klasifikasi yang dibuat mempunyai performa yang bagus. Ketika sistem klasifikasi yang dibuat mempunyai performa yang bagus, maka sistem tersebut bisa digunakan untuk melakukan klasifikasi terhadap data yang lebih banyak.

Naïve bayes dituliskan dengan $P(X|Y)$ dimana X merupakan vektor masukan yang berisi fitur-fitur dan Y adalah label kelas. Notasi tersebut berarti probabilitas kelas Y diperoleh setelah fitur-fitur X diamati. Notasi ini biasa disebut probabilitas akhir (*posterior probability*) untuk Y , sedangkan $P(Y)$ merupakan probabilitas awal (*prior probability*) dari Y . Selama proses pelatihan data dilakukan pembelajaran probabilitas akhir $P(X|Y)$ pada model untuk setiap kombinasi dari X dan Y berdasarkan informasi dari data latih. Pelatihan data itu dilakukan untuk mendapatkan model klasifikasi yang nantinya digunakan untuk mengklasifikasikan data uji X dengan mencari nilai Y nya dengan memanfaatkan nilai $P(X|Y)$ yang didapat (Balagatabi, 2012).

Dalam mengklasifikasikan dokumen teks, ada beberapa proses yang harus dilakukan yaitu:

1. Mencari nilai peluang dari setiap kategori dokumen.
2. Mencari nilai peluang kemunculan dari masing-masing kata pada masing-masing kategori dokumen.
3. Menentukan kategori dokumen uji berdasarkan perhitungan dari proses pertama dan kedua.

I. Pre-proses Teks

Tahap pre-proses teks ini merupakan tahap yang dilakukan untuk mempersiapkan dokumen teks agar siap untuk diklasifikasikan. Tahap ini dilakukan karena dokumen teks yang ada biasanya tidak memiliki struktur yang pasti sehingga informasi di dalamnya tidak bisa diproses secara langsung. Alasan lain dilakukannya tahap ini adalah karena tidak semua kata dalam dokumen teks mencerminkan isi yang terkandung dalam dokumen tersebut. Tahap ini dilakukan sebelum dokumen teks diklasifikasikan dengan metode *naïve bayes classifier*. Tahap pre-proses teks ini meliputi *cleansing*, *case folding*, penghapusan *stopword*, tokenisasi, dan pembobotan kata.

1. *Cleansing*

Cleansing termasuk dalam pre-proses teks. *Cleansing* dilakukan untuk menghapus semua karakter html atau web yang tidak mempunyai makna dalam pengklasifikasian dokumen. Hal ini dikarenakan dalam sebuah *tweet* terkadang disertakan suatu alamat web yang jika tidak dihapus akan mengganggu proses klasifikasi.

2. *Case Folding*

Dalam melakukan klasifikasi dokumen, *case folding* merupakan salah satu hal yang harus dilakukan. Cara melakukan *case folding* yaitu dengan mengubah semua huruf menjadi huruf kecil. Hal ini dilakukan agar kata yang berada di awal kalimat mempunyai arti yang sama dengan kata di tengah atau akhir kalimat (Manning, Raghavan, & Schutze, 2009). Selain itu, di dalam

social media seperti twitter. Ada sebagian orang yang menulis *tweet* menggunakan huruf besar semua untuk melakukan penekanan makna. Sehingga sulit jika data tersebut diklasifikasikan nantinya. Setelah itu, dilakukan penghapusan tanda baca dan angka dan menggantinya dengan karakter spasi.

3. Penghapusan *Stopword*

Biasanya, kata-kata yang sering muncul di setiap kategori dokumen merupakan kata-kata yang tak bermakna dalam pengkategorian dokumen (Manning, Raghavan, & Schutze, 2009). Kata-kata tersebut biasanya disebut *stopword*. Kata-kata tersebut juga tak bisa dijadikan pencari suatu dokumen, sehingga kata-kata tersebut seharusnya dibuang. Kata-kata yang dibuang tersebut kemudian disimpan dalam daftar kata yang disebut *stoplist*. Contoh dari *stopword* adalah kata penghubung dan kata ganti.

4. Tokenisasi

Tokenisasi adalah proses untuk membagi dokumen teks yang dapat berupa kalimat atau paragraf menjadi token-token atau bagian-bagian tertentu (Manning, Raghavan, & Schutze, 2009). Tokenisasi ini dilakukan agar kata-kata yang ada pada kalimat atau paragraf bisa diberi bobot untuk setiap kata. Contohnya seperti kalimat “partai politik memainkan peran penting dalam proses pemilihan kepala daerah” setelah dilakukan tokenisasi maka menjadi “partai – politik – memainkan – peran – penting – dalam – proses – pemilihan – kepala - daerah” yang terdiri dari 10 token.

Tabel 2.2 Contoh Pembobotan IDF

<i>Term</i>	DF	IDF
Gubernur	200	2,32
Politik	50	4,32
Rakyat	100	3,32

Tabel 2.2 menunjukkan contoh hasil perhitungan pembobotan kata dengan metode IDF. Metode perhitungan IDF ini selanjutnya dipadukan dengan pembobotan kata TF sehingga menjadi metode TF-IDF. Jenis TF yang biasa digunakan untuk pembobotan adalah TF murni. Namun, pada perhitungan pembobotan TF-IDF, dilakukan normalisasi TF murni. Normalisasi ini dilakukan dengan cara membagi nilai TF murni dengan jumlah *term* yang ada pada suatu dokumen. Oleh karena itu, rumus umum untuk pembobotan TF-IDF adalah penggabungan perhitungan normalisasi TF murni dengan formula IDF dengan cara mengalikan nilai normalisasi TF murni dan IDF nya. Berikut ini adalah Persamaan 2.5 yang digunakan untuk menghitung pembobotan TF-IDF dari suatu term:

$$W_{ij} = \frac{TF_{ij}}{K_j} \cdot IDF_i \quad (2.5)$$

$$W_{ij} = \frac{TF_{ij}}{K_j} \cdot \log_2 \left(\frac{D}{DF_i} \right) \quad (2.6)$$

Dimana:

W_{ij} = bobot term *i* terhadap dokumen *j*

K_j = jumlah semua term yang ada pada dokumen *j*

TF_{ij} = jumlah kemunculan term *i* pada dokumen *j*

false negative merupakan data yang sebenarnya merupakan kelas positif tetapi diklasifikasikan sebagai kelas negatif.

Ada empat jenis evaluasi yang bisa digunakan untuk melakukan evaluasi sistem klasifikasi yaitu akurasi, presisi, *recall*, dan *f-measure*. Perhitungan empat jenis evaluasi ini dilakukan terhadap data uji yang dimiliki. Akurasi merupakan perbandingan antara banyaknya data yang diklasifikasikan secara benar dibagi oleh jumlah semua data. Sedangkan presisi merupakan perbandingan antara banyaknya *true positive* dibagi oleh jumlah semua data yang diklasifikasikan positif. *Recall* merupakan perbandingan antara banyaknya *true positive* dibagi oleh jumlah semua data yang berkategori positif. Sedangkan, *f-measure* merupakan perhitungan performa sistem klasifikasi yang didapat dengan mengkombinasikan perhitungan presisi dan *recall*.

Evaluasi kinerja sistem klasifikasi yang paling sering digunakan adalah akurasi. Akurasi bisa digunakan untuk mengukur kinerja sistem klasifikasi jika data yang digunakan mempunyai perbandingan jumlah kategori setiap data yang seimbang (Prasetyo, 2014). Namun, terkadang nilai akurasi tidak menggambarkan performa sebenarnya dari sebuah sistem klasifikasi. Hal ini bisa terjadi jika perbandingan jumlah kategori setiap data yang ada sangat tidak seimbang (Prasetyo, 2014). Sebagai contoh, jika terdapat 95 data kategori positif dan hanya terdapat 5 kategori negatif dalam sebuah dataset, sebuah sistem klasifikasi mungkin akan mengklasifikasikan semua dataset sebagai data kategori positif. Akurasi dari sistem tersebut memang sebesar 95%, tetapi secara lebih detail sistem klasifikasi tersebut mempunyai *recognition rate (recall)* kelas positif sebesar 100% sedangkan

pengumpulan data dengan cara *crawling* yang digunakan untuk mendapatkan dataset berupa opini pada twitter. Proses yang kedua yaitu pelabelan data menjadi positif, negatif, netral, dan *outlier*. Proses yang ketiga yaitu penghapusan data netral dan *outlier*. Proses yang keempat yaitu pre-proses teks yang digunakan untuk mengubah dataset mentah menjadi data yang siap untuk diklasifikasikan. Proses yang kelima yaitu klasifikasi dengan menggunakan metode *naive bayes classifier*. Proses yang keenam yaitu uji coba sistem sehingga didapatkan hasil klasifikasi sentimen. Proses yang terakhir yaitu visualisasi dan analisis hasil sentimen masyarakat tentang pemilihan Kepala Daerah Jawa Timur pada tahun 2018. Agar lebih jelas, proses-proses tersebut dijelaskan seperti berikut ini:

1. *Crawling* Data Twitter

Untuk mendapatkan data yang dibutuhkan dalam penelitian ini, dilakukan *crawling* data twitter. Langkah pertama yang dilakukan untuk melakukan *crawling* data twitter yaitu memiliki akun twitter. Dengan akun twitter yang sudah dimiliki tadi, langkah selanjutnya yaitu registrasi di <https://dev.twitter.com/apps/new> untuk mendapatkan kode akses API twitter. Setelah mendapatkan kode akses API twitter, selanjutnya digunakan *software* Rstudio untuk melakukan *crawling* data twitter. *Crawling* ini dilakukan dengan menentukan kata kunci apa yang dicari dan menentukan jumlah *tweet* yang diinginkan. Setelah itu, data hasil *crawling* disimpan dalam bentuk dokumen microsoft excel berekstensi .csv.

2. Pelabelan Data

Dari data yang sudah disimpan dalam bentuk dokumen *microsoft excel* tadi, dilakukan pelabelan data secara manual ke dalam empat kategori yaitu positif, negatif, netral, dan *outlier*. Pelabelan data ini dilakukan agar data siap diklasifikasikan.

3. Penghapusan Data Netral dan *Outlier*

Setelah dilakukan pelabelan data, dilakukan penghapusan data netral dan outlier atau yang tidak berhubungan dengan penelitian ini. Penghapusan data netral ini dilakukan karena keberadaan data netral yang sangat sedikit, sehingga keberadaannya akan mengganggu klasifikasi. Selain itu, data yang bisa digunakan sebagai peta kekuatan calon gubernur hanya data berlabel positif dan negatif saja. Salah satu contoh dari data outlier ini adalah data yang berbahasa Inggris.

4. *Cleansing*

Cleansing termasuk dalam pre-proses teks. *Cleansing* dilakukan untuk menghapus semua karakter html atau web yang tidak mempunyai makna dalam pengklasifikasian dokumen. Hal ini dikarenakan dalam sebuah *tweet* terkadang disertakan suatu alamat web yang jika tidak dihapus akan mengganggu proses klasifikasi.

5. *Case Folding*

Case folding ini juga termasuk dalam pre-proses teks. *Case folding* adalah proses penghapusan tanda baca, angka, dan merubah huruf menjadi huruf kecil semua. Hal ini dilakukan agar dalam klasifikasi kata

yang sama yang penulisan huruf kapitalnya berbeda dianggap menjadi satu kata yang sama.

6. Penghapusan *Stopword*

Stopword adalah kata-kata yang tidak mempunyai arti penting dalam pengkategorian dokumen. Hal ini dikarenakan *stopword* cenderung ada pada semua kategori dokumen. Contoh dari *stopword* adalah kata ganti dan kata penghubung.

7. Tokenisasi

Tokenisasi adalah proses untuk membagi dokumen teks yang dapat berupa kalimat atau paragraf menjadi token-token atau bagian-bagian tertentu. Contohnya seperti kalimat “partai politik memainkan peran penting dalam proses pemilihan kepala daerah” setelah dilakukan tokenisasi maka menjadi “partai-politik-memainkan-peran-penting-dalam-proses-pemilihan-kepala-daerah” yang terdiri dari 10 token.

8. Pembobotan Kata

Setelah dilakukan tokenisasi, setiap kata dalam dokumen diberi bobot dengan dua metode yang berbeda, pertama yaitu TF murni dan yang kedua yaitu TF-IDF seperti yang sudah dijelaskan di tinjauan pustaka. Dalam proses ini, semua kata diproses sehingga setiap kata yang ada pada dokumen tersebut mempunyai bobot sendiri-sendiri. Setelah proses ini selesai, maka kumpulan dokumen (dokumen *corpus*) telah siap untuk dilakukan *training* untuk proses klasifikasi.

Tabel 4.4 Contoh Case Folding

No	Sebelum Data Dilakukan <i>Case Folding</i>	Setelah Data Dilakukan <i>Case Folding</i>
1	@PBB2019: Yusril Ihza Mahendra Pimpin Deklarasi Partai Bulan Bintang Dukung Khofifah – Emil	pbb yusril ihza mahendra pimpin deklarasi partai bulan bintang dukung khofifah emil
2	Diusung PPP, Khofifah Raih Dukungan RKH M Syamsul Arifin Pamekasan	diusung ppp khofifah raih dukungan rkh m syamsul arifin pamekasan
3	@liputan6dotcom: Dukung Gus Ipul- Puti, Anak Milenial Sampang Buat Lagu	liputandotcom dukung gus ipul puti anak milenial sampang buat lagu
4	Gus Ipul: Banyak PR yang Harus Digarap	gus ipul banyak pr yang harus digarap

Tahap selanjutnya yaitu penghapusan *stopword*. Pada tahap ini, semua kata yang masuk dalam daftar *stopword* bahasa Indonesia dihapus. Penghapusan *stopword* ini dilakukan karena kata yang masuk dalam kategori *stopword* tidak bisa digunakan sebagai ciri kategori dokumen tertentu. Tabel 4.5 menunjukkan perbedaan data yang belum dilakukan penghapusan *stopword* dan sudah dilakukan penghapusan *stopword*.

Tabel 4.5 Contoh Penghapusan Stopword

No	Sebelum Dilakukan Penghapusan <i>Stopword</i>	Setelah Dilakukan Penghapusan <i>Stopword</i>
1	gus ipul banyak pr yang harus digarap	gus ipul pr digarap
2	kami selalu mendoakan yang terbaik untuk gus ipul dan mbak puti	mendoakan terbaik gus ipul mbak puti
3	kita lihat bu khofifah dengan energi yang seakan tidak ada habisnya terus berkeliling kemarin dari sendang biru	lihat bu khofifah energi seakan tidak habisnya berkeliling kemarin sendang biru
4	relawan gus dur siap sumbang juta suara untuk khofifah	relawan gus dur siap sumbang juta suara khofifah

tersebut pada kumpulan dokumen yang bersangkutan. Semakin jarang sebuah kata muncul pada dokumen, maka nilai IDF nya semakin besar.

Dua metode pembobotan kata yang digunakan untuk pembuatan sistem mempunyai performa yang sedikit berbeda. Dalam penelitian ini, metode pembobotan TF sedikit lebih unggul dibandingkan dengan pembobotan kata TF-IDF. Untuk data calon gubernur nomor urut satu, pembobotan kata TF dan TF-IDF mempunyai performa yang sama. Hal ini terlihat dari nilai akurasi, presisi, *recall*, dan *f-measure* dari kedua pembobotan kata yang mempunyai nilai sama yaitu berturut-turut 98,99%; 97,78%; 93,44%; dan 95,61%. Tetapi, untuk data calon gubernur dengan nomor urut dua, pembobotan kata TF sedikit lebih unggul dari pembobotan kata TF-IDF. Untuk pembobotan kata TF didapatkan akurasi, presisi, *recall*, dan *f-measure* berturut-turut yaitu 98,95%; 98,55%; 97,78%; dan 98,17% . Sedangkan pembobotan kata TF-IDF mempunyai akurasi, presisi, *recall*, dan *f-measure* berturut-turut yaitu 98,25%; 97,32%; 96,57%; dan 96,95%.

F. Pembuatan Sistem Klasifikasi

Setelah semua kata yang ada pada data yang digunakan sebagai data *training* sudah diberi bobot dengan metode TF dan TF-IDF, maka langkah selanjutnya adalah pembuatan sistem klasifikasi menggunakan metode *naive bayes classifier*. Seperti yang sudah dijelaskan pada tinjauan pustaka, cara kerja *naive bayes classifier* ini didasarkan pada peluang kemunculan kata yang sudah diberi bobot tadi.

Pada pembuatan sistem klasifikasi, tahap pertama yang harus dilakukan adalah membagi data *training* dan *testing*. Pembagian data ini tidak bisa dilakukan secara acak karena akan menyebabkan ketidakseimbangan proporsi yang ada pada data *training* dan data *testing*. Selain itu, data *training* yang digunakan juga harus merupakan data yang mewakili sebagian besar data yang ada, sehingga tidak bisa dilakukan secara acak. Proporsi pembagian data yang sering digunakan yaitu 80% sebagai data *training* dan 20% sebagai data *testing*. Selain itu, ada juga pembagian dengan proporsi 75% data *training* dan 25% data *testing*. Pada sistem ini, digunakan pembagian data *training* 80% dan data *testing* 20%.

Pada sistem yang pertama, data yang digunakan berjumlah 2497 data yang dibagi menjadi data *training* dan data *testing*. Data *training* terdiri dari 2000 data yang mempunyai kategori positif berjumlah 1906 data dan kategori negatif berjumlah 94 data. Sedangkan data *testing* berjumlah 497 data yang terdiri dari 469 data berkategori positif dan 28 data berkategori negatif.

Pada sistem yang kedua, data yang digunakan berjumlah 1487 data yang dibagi menjadi data *training* dan data *testing*. Data *training* terdiri dari 1200 data yang mempunyai kategori positif berjumlah 1132 data dan kategori negatif berjumlah 68 data. Sedangkan data *testing* berjumlah 287 yang terdiri dari 238 data berkategori positif dan 49 data berkategori negatif.

Berikut ini adalah contoh perhitungan manual *naive bayes classifier* berdasarkan kata yang sudah diberi pembobotan TF murni. Contoh pertama yaitu *tweet* dengan isi "PKH Plus untuk lansia dan disabilitas, Sebut Khofifah sebagai

Dari gabungan data pertama dan data kedua yang didapatkan dari proses *crawling* twitter, dapat disimpulkan bahwa calon gubernur nomor urut satu lebih unggul daripada calon gubernur nomor urut dua. Kesimpulan ini dapat diambil berdasarkan fakta bahwa:

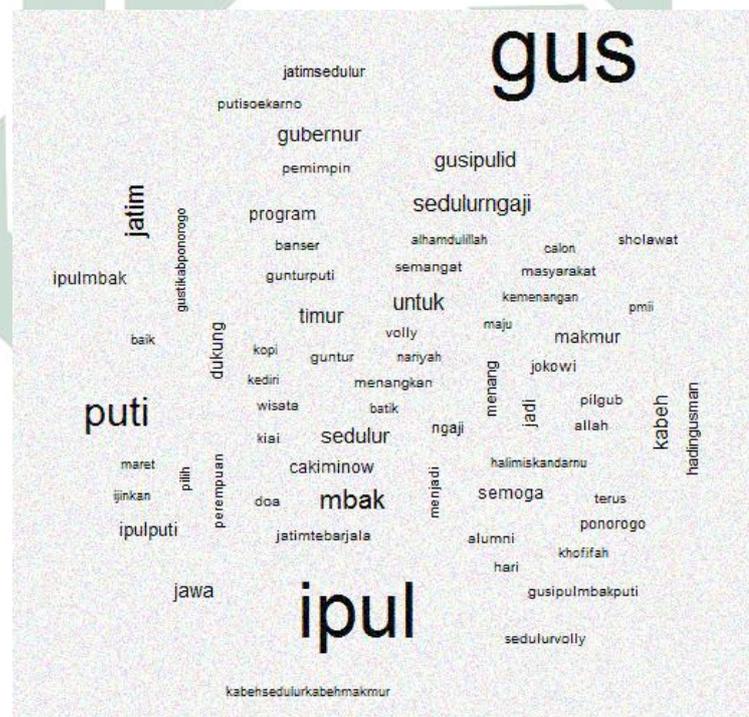
1. Calon gubernur nomor urut satu mendapatkan perhatian yang lebih banyak dari para pengguna twitter dengan bukti lebih banyaknya data yang menyebutkan calon gubernur nomor urut satu.
2. Persentase sentimen positif terhadap calon gubernur nomor urut satu lebih banyak daripada calon gubernur nomor urut dua dan persentase sentimen negatif terhadap calon gubernur nomor urut satu lebih sedikit daripada calon gubernur nomor urut dua.

I. Visualisasi dan Analisis Hasil

Visualisasi dan analisis hasil dalam penelitian ini dibagi menjadi tiga tahap. Tahap yang pertama yaitu menggunakan grafik *wordcloud* yang berasal dari software Rstudio. Grafik *wordcloud* ini merupakan grafik yang menunjukkan frekuensi kata dalam suatu dokumen. Tahap yang kedua yaitu menampilkan analisis hubungan antar kata yang dibuat dengan software VOSviewer. Dalam VOSviewer, visualisasi dari hasil ditampilkan dalam bentuk *network graph*. *Network graph* merupakan sebuah grafik yang menggambarkan hubungan antara kata-kata yang mempunyai frekuensi yang tinggi pada dokumen tersebut. Dari *network graph* tersebut, bisa dianalisis hubungan yang ada pada kata-kata yang ada pada dokumen tersebut. Sedangkan tahap yang ketiga menampilkan perbandingan jumlah *tweet*

pemilihan Kepala Daerah Jawa Timur 2018, tim pemenangan pasangan Khofifah-Emil mengeluarkan fatwa fardhu ain untuk memilih Khofifah-Emil (Ardlyanto, 2018). Hal ini digunakan pendukung pihak lawan untuk melancarkan sentimen negatif yang menyebutkan bahwa tim pemenangan khofifah gunakan segala cara untuk menang, fatwa fardhu ain tak berdasar, dan lain-lain (Ardlyanto, 2018).

Gambar 4.5 adalah tampilan *wordcloud* untuk data sistem kedua yang berkategori positif dengan frekuensi minimal 30 kali:



Gambar 4.5 Wordcloud untuk Data Sistem Kedua yang Berkategori Positif

Gambar 4.5 menunjukkan kata-kata yang mempunyai frekuensi kemunculan yang tinggi pada data sistem kedua yang berkategori positif. Kata-kata yang sering muncul adalah gus ipul, mbak puti, dan slogan-slogan pendukung pasangan calon gubernur nomor urut dua seperti kabehsedulurkabehmakmur, sedulurjatim,

Pada *cluster* pertama yang berwarna merah, ada lima frasa yang berhubungan satu sama lain yaitu khofifah emil unggul, gus ipul putih, bagussholah, khofifah emil, dan tribun jatim. Hubungan yang ada dari frasa-frasa tersebut adalah khofifah emil unggul atas gus ipul putih yang diberitakan oleh tribun jatim. Selain itu, adanya frasa bagussholah sendiri maksudnya adalah kebanyakan *tweet* yang memberitakan berita keunggulan khofifah emil atas gus ipul putih ditulis oleh barisan gus sholah yang merupakan pendukung dari pasangan khofifah dan emil.

Pada *cluster* kedua yang berwarna hijau, terdapat empat frasa yang berhubungan satu sama lain yaitu khofifahip, emildardak, wiswayahe, dukung khofifah emil, dan baguss. Hubungan yang ada dari frasa-frasa tersebut adalah baguss atau singkatan dari barisan gus sholah mendukung pasangan khofifah dan emil dengan semboyan khasnya yaitu wiswayahe. Adapun khofifahip dan emildardak merupakan akun twitter resmi dari Khofifah Indar Parawansa dan Emil Elestianto Dardak.

Pada *cluster* ketiga yang berwarna ungu, terdapat dua frasa yaitu bu khofifah dan barisangussholah. Hubungan yang ada dari frasa-frasa tersebut yaitu barisan gus sholah merupakan pendukung dari bu khofifah.

Pada *cluster* keempat yang berwarna biru, terdapat empat frasa yaitu khofifah, emil, blusukan, dan aliansi santri pemuda ekonom kiai. Hubungan yang ada dari frasa-frasa tersebut adalah aliansi santri pemuda ekonom kiai merupakan salah satu pendukung Khofifah Indar Parawansa untuk mencalonkan diri sebagai gubernur

detikcom banyak memberitakan tentang gus ipul dengan program andalan yaitu seribu dewi atau seribu desa wisata.

Pada *cluster* kedua yang berwarna merah, terdapat enam frasa yang saling berhubungan yaitu gusipul mbak puti, bp pemilu pdip, puti guntur fighter, kabeh sedulur kabeh makmur, gus ipul mbak puti menang, dan pdip surabaya. Hal ini mengindikasikan bahwa pdip melalui bp pemilu mendukung penuh pasangan gus ipul puti. Selain itu, menurut pendukung pasangan gus ipul mbak puti yang mempunyai semboyan kabeh sedulur kabeh makmur, puti merupakan seorang fighter atau pejuang yang membuat para pendukung optimis akan kemenangan pada pemilihan gubernur.

Pada *cluster* ketiga yang berwarna ungu, terdapat tiga frasa yang saling berhubungan yaitu gus ipul puti, tidak punya rasa ego, dan sedulur ngaji ponorogo. Hal ini mengindikasikan bahwa menurut pendukung gus ipul dan mbak puti yang salah satunya adalah sedulur ngaji ponorogo, pasangan ini tidak memiliki ego sehingga merupakan calon pemimpin yang baik.

Pada *cluster* keempat yang berwarna hijau, terdapat enam frasa yang saling berhubungan yaitu menggalang semangat untuk memenangkan gus ipul, final sedulur volly hari minggu, gus ipul dan mbak puti, sedulur volly jokowi, dukung gus ipul, dan cakiminow. Hal ini menunjukkan bahwa final sedulur volly diadakan untuk menggalang semangat untuk memenangkan gus ipul. Selain itu, para pendukung setia gus ipul melakukan klaim bahwa jokowi merupakan salah satu tokoh nasional yang mendukung gus ipul untuk menjadi gubernur jawa timur. Di

- Maulana, I. (2012). *Istilah-Istilah Dasar Dalam Twitter (Pemula)*. Diakses pada 10 April 2018, dari Irfan Maulana's Journal: <https://mazipanneh.wordpress.com>
- Perdana, R. S. (2017). *Pengukuran Akurasi Menggunakan Precision dan Recall*. Diakses pada 10 April 2018, dari Rizal Setya Perdana: <http://www.rizalespe.com>
- Prasetyo, E. (2014). *DATA MINING Mengolah Data Menjadi Informasi Menggunakan Matlab*. Yogyakarta: ANDI.
- Putri, R. K., & Mukhlash, I. (2013). Penerapan Algoritma Klasifikasi Berbasis Aturan Asosiasi untuk Data Meteorologi. *Jurnal Sains dan Seni POMITS*, 1(1), 1-6.
- Rini, D. C., Farida, Y., & Puspitasari, D. (2016). Klasifikasi Menggunakan Metode Hybrid Bayesian-Neural Network (Studi Kasus: Identifikasi Virus Komputer). *JURNAL MATEMATIKA "MANTIK"*, 01(02), 38-43.
- Rodiyansyah, S. F., & Winarko, E. (2012). Klasifikasi Posting Twitter Kemacetan Kota Bandung Menggunakan Naive Bayes Classification. *IJCCS*, 6(1), 91-100.
- Saputro, B. (2016). *Klasifikasi dan Pemetaan Posdata Tematik berbasis Masjid menggunakan Naive Bayes Classifier*. Malang: UIN Maulana Malik Ibrahim.
- Witanto. (2018). *Posko Pemenangan Gus Ipul Tidak Pernah Tolak Tumpeng, Meski Tidak Tahu Maksudnya*. Diakses pada 20 Juni 2018, dari merdeka.com: <https://www.merdeka.com>
- Wong, A. H., & Abednego, L. (2015). *Pengelompokan Dokumen Otomatis dengan menggunakan TFIDf Classifier, Naive Bayes Classifier, dan KNN*. Bandung: Program Studi Teknik Informatika - UNPAR.
- Wongso, R., Luwinda, F. A., Trisnajaya, B. A., Rusli, O., & Rudi. (2017). News Article Text Classification in Indonesian Language. *2nd International Conference on Computer Science and Computational Intelligence* (hal. 137-143). Bali: Elsevier.
- Yulian, E. (2018). Text Mining dengan K-Means Clustering pada Tema LGBT dalam Arsip Tweet Masyarakat Kota Bandung. *JURNAL MATEMATIKA "MANTIK"*, 04(01), 53-58.
- Zhang, W., Yoshida, T., & Tang, X. (2011). A comparative study of TF-IDF, LSI, and Multi-words for text classification. *Expert System with Application*, 38(3), 2758-2765.