

**YAYASAN LEMBAGA PENDIDIKAN ISLAM DAERAH RIAU
UNIVERSITAS ISLAM RIAU
FAKULTAS TEKNIK**

**ANALISIS SENTIMEN PADA TWEET DENGAN TAGAR
#KPUJANGANCURANG MENGGUNAKAN METODE NAÏVE
BAYES
SKRIPSI**

Diajukan Untuk Memenuhi Salah Satu Syarat
Penyusunan Skripsi Pada Fakultas Teknik
Universitas Islam Riau Pekanbaru

DIAN INDRIANI
153510523

PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNIK
UNIVERSITAS ISLAM RIAU
PEKANBARU
2019

KATA PENGANTAR

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

Assalamu'alaikum Warahmatullahi Wabarakatuh.

Dengan mengucapkan Alhamdulillah Robbil 'Alamin, berkat rahmat dan hidayah Allah SWT serta nikmat yang tak terhingga, penulis dapat menyelesaikan laporan skripsi ini dengan judul "ANALISIS SENTIMEN PADA TWEET DENGAN TAGAR #KPUJANGANCURANG MENGGUNAKAN METODE NAÏVE BAYES" sebagai salah satu syarat untuk memperoleh gelar sarjana teknik pada Fakultas Teknik Program Studi Teknik Informatika Universitas Islam Riau.

Dalam penyusunan laporan skripsi ini, penulis menyadari banyak mendapat hambatan dan tantangan. Namun, dalam penyelesaian penulisan ini penulis tidak terlepas dan terwujud tanpa bimbingan, pengarahan, saran dan bantuan moril maupun non-moril dari berbagai pihak. Untuk itu, dalam kesempatan ini penulis ingin menyampaikan ucapan terima kasih yang sebesar-besarnya kepada yang terhormat :

1. Ibu Ause Labellapansa, ST., M.Cs., M.Kom selaku ketua Program Studi Teknik Informatika atas bimbingan dan motivasi yang telah diberikan yang telah membantu dalam kelancaran penyelesaian laporan skripsi ini.
2. Bapak Dr.Arbi Haza Nasution, M.IT selaku dosen pembimbing skripsi yang telah ikhlas dan sabar memberikan bimbingan, saran dan motivasi yang bermanfaat dalam penyusunan laporan skripsi.
3. Seluruh Dosen Teknik Informatika yang telah banyak memberikan ilmu

akademik maupun non-akademik selama berada di bangku pendidikan.

4. Teristimewa untuk kedua orang tua tercinta, kakak, abang dan adik yang senantiasa selalu mendoakan, serta memberikan dukungan baik moril maupun materil dalam menyelesaikan laporan skripsi ini.
5. Rekan-rekan mahasiswa/i khususnya kelas C angkatan 2015 serta semua pihak yang telah memberikan bantuan, semangat dan motivasi selama penyusunan laporan skripsi ini.

Penulis menyadari bahwa dalam penyusunan laporan skripsi ini masih banyak kekurangan, oleh karena itu penulis mengharapkan adanya kritik dan saran yang bersifat membangun untuk memperbaiki laporan skripsi ini.

Akhir kata penulis berharap penyusunan laporan skripsi ini dapat bermanfaat bagi semua pihak yang membacanya dan dapat dikembangkan lebih lanjut, Amin.

Wassalammu'alaikum. Warahmatullahi Wabarakatuh

Pekanbaru, September 2019

Dian Indriani

ABSTRAK

Twitter adalah salah satu media sosial paling berpengaruh yang memiliki jutaan pengguna aktif. Ini biasanya digunakan untuk microblogging yang memungkinkan pengguna untuk berbagi pendapat, ide pemikiran dan banyak lagi. Dengan demikian, jutaan interaksi seperti pesan pendek atau tweet mengalir di antara para pengguna twitter yang membahas berbagai topik yang telah terjadi di seluruh dunia. Penelitian ini bertujuan untuk menganalisis sentimen pengguna terhadap simula topik tren tertentu yang telah dibahas secara aktif dan masif pada waktu itu. Peneliti memilih investasi hashtag #kpujangancurang yang menjadi trending topic selama pemilihan presiden Indonesia pada tahun 2019. Peneliti menggunakan tagar tersebut untuk memperoleh serangkaian data dari Twitter untuk menganalisis dan menyelidiki lebih jauh sentimen positif atau negatif dari pengguna dari tweet mereka. Penelitian ini menggunakan rapid miner untuk menghasilkan data twitter dan metode klasifikasi Nave Bayes untuk mengklasifikasikan sentimen data twitter. Keseluruhan ada 200 pelatihan data dan 100 pengujian data dalam percobaan ini. Keakuratan terakhir dari penganalisa sentimen adalah 0,89.

Kata Kunci : Twitter, Analisis Sentimen, Naive Bayes.

ABSTRAK

Twitter is one of the top influenced social media which has million number of active users. It is commonly used for microblogging that allows users to share messages, ideas thoughts and many more. Thus, millions interaction such as short message or tweet are flowing around among the twitter users discussing various topics that has been happening world-wide. This research aims to analyse a sentiment of the users towards simila a particular trending topic that has been actively and massively discussed at that time. Researcher chose a hashtag #kpujangancurang invest that was the trending topic during the Indonesia presidential election in 2019. Researcher use that hashtag to obtain a set of data from Twitter to analyse and investigate further the positive or the negative sentiment of the users from their tweets. This research utilizes rapid miner tool to generate the twitter data and Nave Bayes classification method to classify the sentiment of the twitter data. There are overall 200 data training and 100 data testing in this experiment. The final accuracy of the sentiment analyzer is 0.89.

Keyword: Twitter, Sentiment Analysis, Naive Bayes.

DAFTAR ISI

	Hal
LEMBAR PERNYATAAN BEBAS PLAGIARISME	i
LEMBAR IDENTITAS PENULIS	ii
HALAMAN PERSEMBAHAN	iii
KATA PENGANTAR	v
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI	ix
DAFTAR TABEL	xii
DAFTAR GAMBAR	xiii
DAFTAR DIAGRAM	xv
BAB I PENDAHULUAN	1
1.1 Latar Belakang Masalah.....	1
1.2 Identifikasi Masalah.....	3
1.3 Rumusan Masalah	3
1.4 Batasan Masalah.....	3
1.5 Tujuan	3
1.6 Manfaat	4
BAB II LANDASAN TEORI	5
2.1 Studi Kepustakaan.....	5
2.2 Dasar Teori.....	7
2.2.1 Twitter	7

2.2.2 Klasifikasi	9
2.2.3 <i>Term Frequency-Inverse Document Frequency (TF-IDF)</i>	9
2.2.4 <i>Evaluation Measure</i>	11
2.2.5 <i>Text Mining</i>	12
2.2.6 <i>Pre-processing</i> Teks	13
2.2.7 Analisis Sentimen	14
2.2.8 <i>Supervised Learning</i>	15
2.2.9 Naïve Bayes	15
2.2.10 <i>Binary Classification</i>	18
2.2.11 Rapid Miner	19
2.2.12 Python	26
2.2.13 Jupyter Notebook	28
2.2.12 Flowchart	28
BAB III METODOLOGI PENELITIAN	30
3.1 Metodologi Penelitian.....	30
3.1.1 Metode Penelitian.....	30
3.1.2 Spesifikasi Perangkat Keras (<i>Hardware</i>).....	31
3.1.3 Spesifikasi Perangkat Lunak (<i>Software</i>).....	32
3.2 Perancangan Sistem	32
3.2.1 Gambaran Umum.....	32
3.2.2 Proses <i>Dataset</i>	33
3.2.3 <i>Preprocessing</i>	35
3.2.4 <i>Term Frequency-Inverse Document Frequency (TF-IDF)</i>	36

3.2.5 Multinomial Naïve Bayes	39
BAB IV HASIL DAN PEMBAHASAN	44
4.1 Model Analisis Sentimen	44
4.1.1 Pengujian Pada Data Latih	45
4.1.2 Pengujian Pada Data Uji	47
4.2 <i>Evaluation Measure</i>	49
4.2.1 <i>Accuracy</i>	52
4.2.2 <i>Precision</i>	54
4.2.3 <i>Recall</i>	56
4.2.4 <i>F-Measure</i>	58
4.3 Antarmuka Pada Analisis Sentimen.....	62
4.4 Pengujian Kepada <i>User</i>	64
4.4.1 Implementasi <i>User</i>	64
4.4.2 Kesimpulan Implementasi Sistem.....	66
BAB V KESIMPULAN DAN SARAN	68
DAFTAR PUSTAKA	69
LAMPIRAN	

DAFTAR TABEL

	Hal
Tabel 2.1 Tabel <i>Confusion Matrix</i>	11
Tabel 2.2 Simbol <i>Flowchart</i>	29
Tabel 3.1 Proses <i>Preprocessing</i>	35
Tabel 3.2 Contoh Dokumen	36
Tabel 3.3 TF-IDF Doc A	37
Tabel 3.4 TF-IDF Doc B	38
Tabel 3.5 TF-IDF Doc C	39
Tabel 3.6 Tabel Bobot	40
Tabel 3.7 Likelihood	42
Tabel 4.1 Pengujian Data Latih	45
Tabel 4.2 Hasil Pengujian Data Latih	46
Tabel 4.3 Tabel <i>Confusion Matrix</i>	50
Tabel 4.4 <i>Confusion Matrix</i> Data A	51
Tabel 4.5 <i>Confusion Matrix</i> Data B	51
Tabel 4.6 <i>Confusion Matrix</i> Data C	51
Tabel 4.7 Hasil <i>Accuracy</i>	52
Tabel 4.8 Hasil <i>Precision</i>	54
Tabel 4.9 Hasil <i>Recall</i>	56
Tabel 4.10 Hasil <i>F-Measure</i>	58
Tabel 4.11 Hasil Persentase Kuisisioner	66

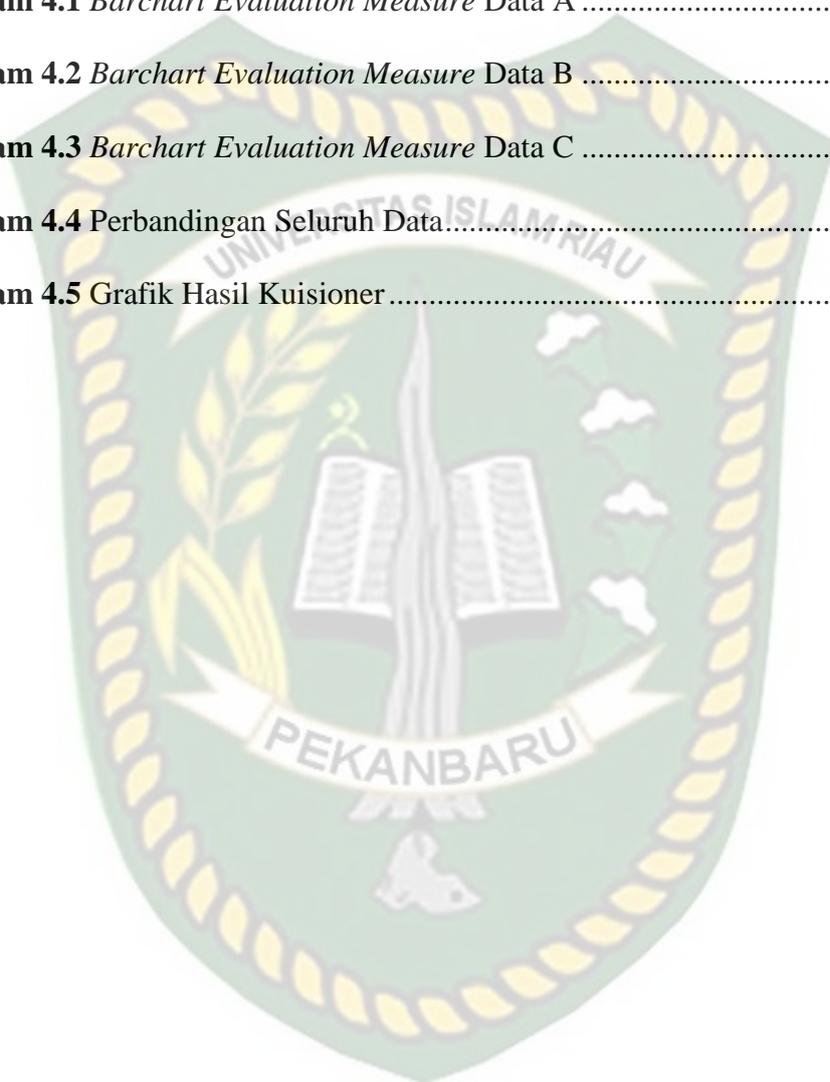
DAFTAR GAMBAR

	Hal
Gambar 2.1 <i>Search</i> Twitter	20
Gambar 2.2 <i>Search Code</i> Twitter	21
Gambar 2.3 <i>Open</i> Url	21
Gambar 2.4 Izin Akses	22
Gambar 2.5 <i>Code</i> API Twitter.....	22
Gambar 2.6 <i>Copy Code</i> API.....	23
Gambar 2.7 <i>Test Code</i> API.....	23
Gambar 2.8 <i>Write</i> Excel	24
Gambar 2.9 Simpan Folder Excel.....	24
Gambar 2.10 <i>Search</i> Twitter dan <i>Write</i> Excel.....	25
Gambar 2.11 Koneksi Twitter	25
Gambar 2.12 <i>Result</i> Data.....	25
Gambar 2.11 Koneksi Twitter	25
Gambar 3.1 Gambaran Umum Analisis Sentimen	32
Gambar 3.2 Diagram Alir <i>Dataset</i>	34
Gambar 4.1 Tampilan Data Latih	46
Gambar 4.2 Tampilan <i>Form Input</i> Data Uji	47
Gambar 4.3 Tampilan Ketepatan Prediksi.....	49
Gambar 4.4 Perhitungan <i>Accuracy</i> Data A	53
Gambar 4.5 Perhitungan <i>Accuracy</i> Data B.....	53
Gambar 4.6 Perhitungan <i>Accuracy</i> Data C.....	53

Gambar 4.7 Perhitungan <i>Precision</i> Data A	55
Gambar 4.8 Perhitungan <i>Precision</i> Data B	55
Gambar 4.9 Perhitungan <i>Precision</i> Data C	55
Gambar 4.10 Perhitungan <i>Recall</i> Data A	57
Gambar 4.11 Perhitungan <i>Recall</i> Data B.....	57
Gambar 4.12 Perhitungan <i>Recall</i> Data C.....	57
Gambar 4.13 Perhitungan <i>F-Measure</i> Data A.....	59
Gambar 4.14 Perhitungan <i>F-Measure</i> Data B	59
Gambar 4.15 Perhitungan <i>F-Measure</i> Data C	59
Gambar 4.16 Form <i>Input User</i>	63
Gambar 4.17 Form <i>Output User</i>	63

DAFTAR DIAGRAM

	Hal
Diagram 4.1 <i>Barchart Evaluation Measure Data A</i>	60
Diagram 4.2 <i>Barchart Evaluation Measure Data B</i>	60
Diagram 4.3 <i>Barchart Evaluation Measure Data C</i>	61
Diagram 4.4 Perbandingan Seluruh Data	62
Diagram 4.5 Grafik Hasil Kuisisioner	65



BAB I

PENDAHULUAN

1.1 Latar Belakang Masalah

Twitter adalah layanan jejaring social dan mikroblog daring yang memungkinkan penggunaanya untuk mengirim dan membaca pesan berbasis teks hingga 140 karakter akan tetapi pada tanggal 07 November 2017 bertambah hingga 280 karakter yang dikenal dengan sebutan kicauan. Chief Executive Officer Twitter, Dick Costolo, mengatakan bahwa Indonesia menjadi salah satu negara dengan pengguna Twitter terbanyak. Karena itulah Twitter pun akhirnya mendirikan kantor di Jakarta. Berdasarkan laporan Twitter di kuartal IV 2014, total pengguna aktifnya mencapai 288 juta per bulan. Menurut Costolo pengguna media sosial di Indonesia memiliki pengetahuan yang baik dengan dunia digital. Disamping itu pengguna Twitter di Indonesia dinilai sangat atraktif dan bersemangat dan dianggap sangat aktif menuliskan cuitan. Tidak jarang, hasil obrolan di lini masa menjadi trending topic atau topik yang paling banyak dibicarakan di seluruh dunia (Movementi, 2015).

Sesuai dengan Undang-Undang Republik Indonesia Nomor 15 Tahun 2011 tentang Penyelenggara Pemilihan Umum, yang dimaksud dengan Pemilu adalah sarana pelaksanaan kedaulatan rakyat yang diselenggarakan secara langsung, umum, bebas, rahasia, jujur, dan adil dalam Negara Kesatuan Republik Indonesia berdasarkan Pancasila dan Undang-Undang Dasar Negara Republik Indonesia

Tahun 1945. Komisi Pemilihan Umum (KPU) adalah lembaga Penyelenggara Pemilu yang bersifat nasional, tetap, dan mandiri yang bertugas melaksanakan Pemilu. KPU Provinsi dan KPU Kabupaten/Kota adalah Penyelenggara Pemilu di Provinsi dan Kabupaten/Kota. Wilayah kerja KPU meliputi seluruh wilayah Negara Kesatuan Republik Indonesia. KPU menjalankan tugasnya secara berkesinambungan dan dalam menyelenggarakan Pemilu, KPU bebas dari pengaruh pihak manapun berkaitan dengan pelaksanaan tugas dan wewenangnya. KPU berkedudukan di Ibu Kota Negara Republik Indonesia, KPU Provinsi berkedudukan di ibu kota provinsi, dan KPU Kabupaten/Kota berkedudukan di ibu kota kabupaten/kota.

Saat pemilu yang diadakan pada tahun 2019, warga twitter membuat tagar #kpujangancurang di twitter karna menganggap telah terjadi kesalahan pada saat perhitungan suara.

Berdasarkan latar belakang tersebut, penulis mengangkat judul skripsi “Analisis Sentiment pada Tweet dengan tagar #kpujangancurang menggunakan metode Naïve-Bayes”. Dengan input berupa data *tweet* dalam Bahasa Indonesia, akan dilakukan klasifikasi dengan algoritma Naïve-Bayes untuk menentukan apakah *tweet* tersebut bersentimen positif atau negatif.

1.2 Identifikasi Masalah

Adapun identifikasi masalah yang dapat diambil dari latar belakang tersebut adalah sebagai berikut “Dibutuhkan analisis sentimen untuk menganalisa opini masyarakat mengenai topik yang sedang populer dibicarakan di sosial media twitter”.

1.3 Rumusan Masalah

Berdasarkan dari latar belakang tersebut, maka dapat dirumuskan permasalahan yang ada di atas yaitu “Bagaimana mengamati dan menganalisa opini pengguna Twitter mengenai kasus kpu yang dianggap curang dan kemudian mengklasifikasikan sentimen data *tweets* tersebut”.

1.4 Batasan Masalah

Adapun batasan masalah dalam penelitian ini adalah sebagai berikut :

1. Data yang digunakan untuk penelitian ini adalah data *tweets* dalam Bahasa Indonesia.
2. Metode yang digunakan untuk mengklasifikasikan dalam penelitian ini adalah metode Naïve-Bayes.
3. Data yang digunakan sesuai topik yang diangkat yaitu kpu jangan curang.
4. Hanya menampilkan sentimen positif dan negatif dalam bentuk tabel.

1.5 Tujuan

Adapun tujuan dari penelitian yaitu mengimplementasikan algoritma Naïve-Bayes dalam menganalisa sentimen pengguna Twitter dengan topik kpu yang dianggap curang dalam penghitungan suara.

1.6 Manfaat

Adapun manfaat dari penelitian ini adalah sebagai berikut “Penelitian ini diharapkan dapat memberikan manfaat baik kepada penulis maupun pembaca tentang gambaran sentimen para pengguna Twitter terhadap kasus kpu yang dianggap curang oleh masyarakat”.



Dokumen ini adalah Arsip Miik :

Perpustakaan Universitas Islam Riau

BAB II

LANDASAN TEORI

2.1 Studi Kepustakaan

Sejumlah Penelitian telah dilakukan sebelumnya yang menjadi rujukan yaitu penelitian, oleh (Akhmad Debiyanto, 2018) dengan penelitian yang berjudul “*Penerapan Sentiment Analysis pada Pengguna Twitter Menggunakan Metode K-Nearest Neighbor*”. Penelitian ini bertujuan untuk melakukan analisis sentimen para pengguna Twitter terhadap ketiga pasangan kandidat pada Pilkada DKI 2017. Dengan *input* berupa data *tweet* dalam Bahasa Indonesia, akan dilakukan klasifikasi dengan algoritma KNN (K-Nearest Neighbor) untuk menentukan apakah *tweet* tersebut bersentimen positif atau negatif.

Penelitian kedua dilakukan oleh (Noviah Dwi Putranti & Edi Winarko, 2014), dengan penelitian berjudul “*Analisis Sentiment Twitter untuk Teks Berbahasa Indonesia dengan Maximum Entropy dan Support Vector Machine*”. Penelitian ini bertujuan untuk mengetahui sentiment publik mengenai sesuatu dengan menggunakan pendekatan dalam *machine learning* yang dikenal dengan nama *Support Vector Machine* dan *Maximum Entropy Part of Speech Tagging* yang di khususkan pada dokumen teks berbahasa Indonesia dengan fitur unigram. Pendekatan model *Maximum Entropy* (ME) dipilih dalam *Part of Speech* karena terbukti memiliki cara yang sangat efisien untuk mengintegrasikan satu set fitur yang sangat besar dalam model dengan mudah dan telah berhasil digunakan

dalam tugas-tugas seperti *Natural Language Processing* (NLP) sebagai bagian dari penandaan pidato atau informasi ekstraksi.

Penelitian selanjutnya dilakukan oleh (Arsaningtyas, Bijaksana, & Faraby, 2018) dengan penelitian berjudul “*Analisis Sentiment Tingkat Kepuasan Pengguna Penyedia Layanan Telekomunikasi Seluler Indonesia Pada Twitter Dengan Metode Support Vector Machine dan Lexicon Based Features*”, Penelitian ini mengelompokkan polaritas dari teks dalam kalimat atau dokumen untuk mengetahui apakah opini pada kalimat atau dokumen tersebut apakah termasuk apakah termasuk positif atau negatif. Penelitian ini dilakukan agar para perusahaan penyedia layanan telekomunikasi seluler dapat melihat bagaimana tingkat kepuasan masyarakat terhadap produk layanan mereka sehingga mereka akan terpacu untuk meningkatkan layanan dan mutu pada produk layanan mereka.

Penelitian berikutnya dilakukan oleh (Robet Habibi, Djoko Budiyanto Setyohadi, & Ernawati, 2016) dengan penelitian berjudul “*Analisis Sentimen Pada Twitter Mahasiswa Menggunakan Metode Backpropagation*”. Penelitian ini bertujuan untuk mengetahui kecenderungan emosi mahasiswa dengan analisis sentimen pada twitter menggunakan metode backpropagation. Hasil analisis sentimen merupakan kecenderungan emosi mahasiswa. Kecenderungan emosi mahasiswa dapat digunakan sebagai acuan untuk menentukan perlakuan yang sesuai terhadap mahasiswa pada saat proses belajar.

Pada penelitian lain yang dilakukan oleh (M.Ali Fauzi, 2018) dengan judul penelitian “*Pendekatan Metode Random Forest Pada Analisis Sentimen Dalam Bahasa Indonesia*”. Dalam penelitian ini, peneliti mengeksplorasi penggunaan

Random Forest untuk klasifikasi sentimen dalam bahasa Indonesia. *Random Forest* adalah teknik pembelajaran ensemble berdasarkan algoritma pohon keputusan. Dalam penelitian ini, peneliti juga akan mengeksplorasi penggunaan fitur kata-kata (BOW) dengan beberapa istilah variasi metode pembobotan seperti Binary TF, Raw TF, Logarithmic TF dan TF.IDF.

2.2 Dasar Teori

2.2.1 *Twitter*

Lesmana (2012: 61) twitter pertama kali resmi digunakan pada tanggal 13 juli 2006. Twitter adalah sebuah situs 4 web yang menawarkan jaringan sosial berupa mikroblog sehingga memungkinkan penggunanya untuk mengirimkan dan membaca pesan yang disebut kicaua (tweets), yang bebas mengekspresikan sesuatu seperti curhat/kritik terhadap kebijakan pemerintah. Kicauan adalah berupa teks tulisan hingga 140 karakter yang ditampilkan pada halaman profil penggunanya. Kelebihan twitter dibanding dengan media sosial lainnya menurut Putra (2014: 33) diantaranya adalah jangkauannya luas, tidak hanya teman, tetapi juga mampu menjangkau publik figur, potensi periklanan di masa mendatang lebih besar, komunikasi terjadi sangat cepat (update), multilink (terhubung dengan banyak jaringan) dan lebih terukur dari facebook. Twitter membantu penyebaran informasi secara lebih cepat yang kemudian akan menjadi sebuah topik yang dibahas oleh para penggunanya. Media massa seperti televisi, koran, majalah, tabloid pun menggunakan twitter sebagai penyampai berita-beritanya. Hal ini mempermudah masyarakat memperoleh informasi secara cepat dan update karena berita dapat di update setiap saat oleh media massa melalui twitter.

Menurut Sedarmayanti (2009: 313) perolehan dan penyebarluasan informasi dapat difasilitasi melalui media internet, penggunaan internet dalam rangka meningkatkan kinerja governance sudah menjadi perhatian banyak pihak, termasuk pejabat publik. Twitter juga digunakan pejabat publik untuk berinteraksi dengan masyarakat. Dalam konteks ini masyarakat dapat secara langsung memberikan pendapatnya ataupun komentarnya terhadap sebuah berita yang dikeluarkan oleh pejabat publik tertentu melalui twitter. Begitu juga sebaliknya, pejabat publik dapat mengetahui secara langsung dan cepat tanggapan dari para pembacanya. Karakteristik yang paling populer dari berita online adalah sifatnya yang real time, mendapatkan pendalaman dan titik pandang yang lebih luas bahkan berbeda. Interaktivitas juga dapat dilihat dari adanya pemberian feed back atau umpan balik dari pembaca yang membaca sebuah berita melalui kolom komentar. Berita, kisah-kisah, maupun peristiwa, bisa langsung dipublikasikan pada saat kejadian sedang berlangsung.

Penelitian-penelitian yang menggunakan twitter sebagai data sangatlah banyak. Karena memudahkan dalam mengambil data. Salah satu penelitian yang menggunakan twitter sebagai data adalah analisis sentimen.

2.2.2 Klafisikasi

Klasifikasi merupakan suatu pekerjaan menilai objek data untuk memasukkannya ke dalam kelas tertentu dari sejumlah kelas yang tersedia. Dalam klasifikasi ada dua pekerjaan utama yang dilakukan, yaitu : pertama, Pembangunan model sebagai *prototype* untuk disimpan sebagai memori dan kedua, Penggunaan model tersebut untuk melakukan pengenalan/ klasifikasi/

prediksi pada suatu objek data lain agar diketahui di kelas mana objek data tersebut dalam model yang mudah disimpan (Prasetyo, Eko. 2012).

Contoh aplikasi yang sering ditemui adalah pengklasifikasian jenis hewan, yang mempunyai sejumlah atribut. Dengan atribut tersebut, jika ada hewan baru, kelas hewannya bisa langsung diketahui. Contoh lain adalah bagaimana melakukan diagnosis penyakit kulit kanker melanoma (Prasetyo, Eko. 2012), yaitu dengan melakukan pembangunan model berdasarkan data latih yang ada, kemudian menggunakan model tersebut untuk mengidentifikasi penyakit pasien baru sehingga diketahui apakah pasien tersebut menderita kanker atau tidak.

2.2.3 *Term Frequency – Inverse Document Frequency (TF-IDF)*

Metode *Term Frequency-Inverse Document Frequency* (TF-IDF) adalah cara pemberian bobot hubungan suatu kata (*term*) terhadap dokumen. Untuk dokumen tunggal tiap kalimat dianggap sebagai dokumen. Metode ini menggabungkan dua konsep untuk perhitungan bobot, yaitu *Term frequency* (TF) merupakan frekuensi kemunculan kata (*t*) pada kalimat (*d*). *Document frequency* (DF) adalah banyaknya kalimat dimana suatu kata (*t*) muncul.

Frekuensi kemunculan kata di dalam dokumen yang diberikan menunjukkan seberapa penting kata itu di dalam dokumen tersebut. Frekuensi dokumen yang mengandung kata tersebut menunjukkan seberapa umum kata tersebut. Bobot kata semakin besar jika sering muncul dalam suatu dokumen dan semakin kecil jika muncul dalam banyak dokumen. Pada algoritma TF-IDF digunakan rumus untuk menghitung bobot (W) masing masing dokumen terhadap kata kunci dengan persamaan (Winata, Rainarli, Informatika, & Indonesia, 2016) :

$$tf = 0.5 + 0.5 * \frac{ft,d}{\max(ft,d)} \dots\dots\dots(2.1)$$

$$idf = \log \left(\frac{N}{dft} \right) \dots\dots\dots(2.2)$$

$$W = tf * idf \dots\dots\dots(2.3)$$

Keterangan:

d : dokumen

t : kata pada dokumen

ft,d : frekuensi kata pada d

tf : banyaknya kata i pada sebuah dokumen

N : total jumlah dokumen

dft : banyak dokumen yang mengandung kata i

idf : *Inversed Document Frequency*

W : bobot dokumen ke-d terhadap katake-t

2.2.4 *Evaluation Measure*

Evaluasi bertujuan untuk menilai performansi yang dapat dicapai oleh sistem. Evaluasi dalam tugas akhir ini digunakan untuk mengetahui apakah suatu sistem telah optimal dalam mendeteksi halaman yang terindikasi memiliki kemiripan semantik terhadap halaman yang lain. Evaluasi yang digunakan adalah *precision*, *recall*, akurasi dan *F-Measure(F1-Score)*.

Precision mengidentifikasi kualitas dari klasifikasi sistem, sedangkan *recall* mengidentifikasi kuantitas dari sistem, dan akurasi merupakan rasio prediksi benar (positif dan negatif) dengan keseluruhan data. *F-Measure*

merupakan salah satu perhitungan evaluasi dalam informasi temu kembali yang mengkombinasikan recall dan precission.

Tabel 2. 1 Tabel *Confusion Matrix*

Kelas	Terklasifikasi Positif	Terklasifikasi Negatif
Positif	TP (True Positive)	FP (False Positive)
Negatif	FN (False Negative)	TN (True Neagtive)

Rumus perhitungan evaluasinya adalah sebagai berikut(Jadhira, Bijaksana, & Wahyudi, 2018b):

$$Precision = \frac{TP}{TP+FP} \dots\dots\dots(2.4)$$

$$Recall = \frac{TP}{TP+FN} \dots\dots\dots(2.5)$$

$$Accuracy = \frac{TP+TN}{(TP+FP+FN+TN)} \dots\dots\dots(2.6)$$

Keterangan:

True Positive (TP) = suatu kondisi dimana sistem mendekteksi kelas positif dan faktanya pun positif

True Negative (TN) = suatu kondisi dimana sistem mendekteksi kelas negatif dan faktanya pun negatif

False Positive (FP) = suatu kondisi dimana sistem mendekteksi kelas positif namun faktanya negatif

False Negative (FN) = suatu kondisi dimana sistem mendekteksi kelas negatif namun faktanya positif.

F-Measure (*F1-Score*) adalah harmonic mean dari precision dan recall dengan rumus sebagaiberikut:

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \dots\dots\dots(2.7)$$

2.2.5 *Text Mining*

Text mining yang juga dikenal dengan *text data mining* atau pencarian pengetahuan di basis data textual yaitu sebuah proses yang semi otomatis melakukan ekstraksi dari pola yang ada di *database*. Dari hasil ekstraksi tersebut munculah pengetahuan baru yang bisa dimanfaatkan untuk kepentingan pengambilan keputusan. *Text mining* mempunyai kesamaan dengan *data mining*. Keduanya memiliki tujuan yang sama yaitu untuk memperoleh informasi dan pengetahuan dari sekumpulan data yang sangat besar. Data tersebut bisa berbentuk sebuah *database*. Namun keduanya memiliki perbedaan jenis data. *Data mining* memiliki input data dari data yang sudah terstruktur sedangkan *text mining* dimulai dengan data yang tidak terstruktur.

Pemanfaatan dari *text mining* secara nyata sangatlah luas. Areanya seluas data tekstual yang terbentuk seperti di area hukum dengan data putusan pengadilan, penelitian dengan data artikel penelitian, keuangan dengan data laporan triwulan, teknologi dengan data arsip paten, pemasaran dengan data komentar konsumen, dan di area lainnya. Sebagai contoh sebuah perusahaan membuat formulir yang biasa diisi apabila konsumennya ingin memuji, komplain,

ataupun klaim garansi. Dari kartu formulir tersebut terbentuklah data yang sangat besar dan bisa digunakan untuk mengidentifikasi secara objektif produk dan layanan dari suatu perusahaan menggunakan text mining. Selain itu proses *text mining* yang dilakukan secara otomatis adalah dibidang komunikasi elektronik dan email. *Text mining* tidak hanya mengklasifikasikan dan menyaring email sampah, tetapi bisa juga memprioritaskan *email* secara otomatis berdasarkan tingkat kepentingannya.

2.2.6 *Pre-processing Teks*

Data yang didapat dari hasil *crawling* belum bisa langsung diklasifikasikan karena data tersebut masih terdapat banyak simbol dan kata-kata yang tidak diperlukan, karena itu kita memerlukan *pre-processing* data agar data lebih terstruktur dan bersih sehingga bisa diklasifikasikan. *Text pre-processing* adalah bagian dimana data yang sudah didapat selanjutnya diolah dengan tahapan-tahapan *case folding*, *tokenizing*, *filtering*, *stopword* dan *stemming*.

1. *Case folding*, pada tahap ini dilakukan proses perubahan dan huruf besar menjadi huruf kecil dan menghilangkan seluruh tanda baca pada kalimat.
2. *Tokenizing*, pada setiap kata akan dipisahkan berdasarkan spasi yang ditentukan.
3. *Filtering*, yaitu pembuangan kata-kata tidak penting dan daftar *stopwords* bahasa indonesia (Tala, F.Z, 2003).
4. *Stemming* mengubah kata menjadi kata dasar, misalnya “menolong” menjadi “tolong”.
5. *Stopword* adalah kata-kata yang tidak deskriptif (tidak penting) yang

dapat dibuang dengan pendekatan bag-of-words.

2.2.7 Analisis Sentimen

Sejarah analisis sentimen pertama kali muncul pada sebuah jurnal karya Das, Chen, dan Tong pada tahun 2001, bahasan yang mereka angkat sesuai dengan minat mereka yaitu menganalisis sentiment pasar. Analisis sentimen adalah mengekstraksi pendapat, sentimen, evaluasi, dan emosi orang tentang suatu topik tertentu yang tertulis menggunakan teknik pemrosesan bahasa alami. Sejumlah karya-karya besar lainnya menyebutkan analisis sentimen fokus pada aplikasi spesifik yang mengklasifikasikan mengenai sifat yang berlawanan (antara positif dan negatif). Dari pengertian tersebut menjadi sebuah fakta yang menyebabkan beberapa penulis bahwa istilah analisis sentimen mengacu pada tugas yang sempit atau spesifik. Namun saat ini banyak yang menafsirkan istilah analisis sentiment lebih luas lagi yang berarti cara pengkomputasian pendapat, sentimen, dan subjektivitas pada teks.

2.2.8 *Supervised Learning*

Supervised learning adalah tugas pembelajaran mesin untuk mempelajari fungsi yang memetakan input ke output berdasarkan contoh pasangan input-output. Ini menyimpulkan fungsi dari data pelatihan berlabel yang terdiri dari serangkaian contoh pelatihan. Algoritma pembelajaran yang diawasi menganalisis data pelatihan dan menghasilkan fungsi yang disimpulkan, yang dapat digunakan untuk memetakan contoh-contoh baru (data uji). Skenario optimal akan memungkinkan algoritma menentukan label kelas dengan benar. Ini membutuhkan algoritma pembelajaran untuk menggeneralisasi dari data pelatihan

dengan cara yang masuk akal. *Supervised learning* memiliki berbagai macam metode yang dapat digunakan. Metode yang digunakan dalam penelitian ini adalah naïve bayes.

2.2.9 Naïve Bayes

Pengklasifikasi Bayesian (Tan, Steinbach, & Kumar, 2006), adalah pengklasifikasi yang didasarkan pada teorema Bayes. Sebuah studi yang membandingkan algoritma klasifikasi telah menemukan sebuah pengklasifikasi Bayesian sederhana yang dikenal dengan pengklasifikasi naïve Bayesian (*naïve Bayesian classifier*) yang secara performa sebanding dengan pohon keputusan dan pengklasifikasi jaringan saraf tertentu. Pengklasifikasi Bayesian juga menunjukkan akurasi dan kecepatan yang tinggi saat diterapkan pada basis data yang besar.

Pengklasifikasi naïve bayes mengansumsikan bahwa pengaruh dari nilai atribut pada kelas tertentu tidak bergantung pada nilai atribut lainnya. Asumsi ini disebut *class conditional independence*. Asumsi ini dibuat untuk menyederhanakan perhitungan yang rumit dan dalam arti ini dianggap “naïve”.

1. Multinomial Naïve Bayes atau Multinomial NB

Multinomial NB merupakan model pengembangan dari algoritma bayes yang cocok dalam pengklasifikasian teks atau dokumen. Pada formula Multinomial Naïve Bayes, kelas tidak hanya ditentukan dengan kata yang muncul tetapi juga jumlah kemunculannya. Model multinomial mengambil jumlah kata yang muncul pada sebuah dokumen, dalam model multinomial sebuah dokumen, dalam model multinomial sebuah dokumen terdiri dari

beberapa kejadian kata dan diasumsikan panjang dokumen tidak bergantung pada kelasnya. Dengan menggunakan asumsi Bayes yang sama bahwa kemungkinan tiap kejadian kata dalam sebuah dokumen adalah bebas tidak terpengaruh dengan konteks kata dan posisi kata dalam dokumen.

Pada saat klasifikasi, algoritma akan mencari probabilitas tertinggi dari semua kategori dokumen yang diujikan (vmap). Hal yang dilakukan pertama adalah menghitung probabilitas dengan rumus sebagai berikut (Manning, Raghavan, & Schutze, 2008):

$$P(c_i) = \frac{N_{ci}}{N} \dots\dots\dots(2.8)$$

Keterangan :

- $P(C_i)$: probabilitas kelas
 N_{ci} : jumlah kelas c pada seluruh dokumen
 N : jumlah seluruh dokumen

Untuk memperkirakan *conditional probability* persamaan yang digunakan yaitu:

$$P(w|c) = \frac{W_{ct}+1}{(\sum_{w' \in V} W_{ct'})+B'} \dots\dots\dots(2.9)$$

Keterangan:

- $P(w|c)$: *conditional probability*
 W_{ct} : W dari *term* t di kategori c
 $\sum_{w' \in V} W_{ct'}$: jumlah total W dari keseluruhan term yang berada di kategori c
 B' : nilai idf pada seluruh dokumen

Untuk menghilangkan nilai nol, digunakan add-one atau Laplace smoothing. Proses ini menambahkan nilai satu (1) pada setiap nilai W_{ct} dari perhitungan *conditional probability*.

Lalu hitung V_{MAP} pada kalimat dengan rumus sebagai berikut :

$$V_{MAP} = \operatorname{argmax} P(x_1, x_2, x_3, \dots, x_n) P(C) \dots \dots \dots (2.10)$$

Keterangan:

$P(x_1)$: probabilitas pada kata

$P(C)$: probabilitas pada kelas

Langkah terakhir setelah dapatkan V_{MAP} adalah menentukan kelas pada kalimat. Jika nilai V_{MAP} positif lebih tinggi dari V_{MAP} negatif, maka nilai dari kalimat tersebut adalah positif, begitu juga sebaliknya.

2.2.10 *Binary Classification*

Klasifikasi biner atau binomial adalah tugas untuk mengklasifikasikan elemen-elemen dari himpunan yang diberikan ke dalam dua kelompok (memprediksi kelompok mana yang masing-masing dimiliki) berdasarkan aturan klasifikasi. Konteks yang membutuhkan keputusan apakah suatu item memiliki sifat kualitatif atau tidak, beberapa karakteristik tertentu atau beberapa klasifikasi biner khas. Klasifikasi biner adalah dikotomisasi yang diterapkan untuk tujuan praktis, dan dalam banyak masalah klasifikasi biner praktis, kedua kelompok tidak simetris dari akurasi keseluruhan, proporsi relatif dari berbagai jenis kesalahan yang menarik. Misalnya dalam pengujian medis, *false positive* (mendeteksi penyakit ketika tidak ada) dianggap berbeda dari *false negative* (tidak mendeteksi penyakit ketika hadir).

Klasifikasi statistik adalah masalah yang dipelajari dalam pembelajaran mesin. Ini adalah jenis pembelajaran terawasi, metode pembelajaran mesin dimana kategori sudah ditentukan sebelumnya, dan digunakan untuk mengkategorikan pengamatan probabilistik baru ke dalam kategori tersebut. Ketika hanya ada dua kategori masalahnya dikenal sebagai klasifikasi biner statistik. Beberapa metode yang biasa digunakan untuk klasifikasi biner adalah *decision trees*, *random forests*, *bayesian networks*, *support vector machines*, *neural networks*, *logistic regression*, *probit model*. Setiap classifier terbaik hanya dalam domain pilih berdasarkan jumlah pengamatan, dimensi vektor fitur, kebisingan dalam data dan banyak faktor lainnya. Misalnya *random forest* berkinerja lebih baik dari pada pengklasifikasi SVM untuk awan titik 3D. Ada banyak metrik yang dapat digunakan untuk mengukur kinerja pengklasifikasi atau prediksi. Bidang yang berbeda memiliki preferensi yang berbeda untuk metrik tertentu karena tujuan yang berbeda. Misalnya, dalam sensitivitas dan spesifisitas obat sering digunakan, sedangkan dalam pengambilan informasi ketepatan dan daya ingat lebih disukai. Perbedaan penting adalah antara metrik yang tidak tergantung pada prevalensi (seberapa sering setiap kategori terjadi dalam populasi), dan metrik yang bergantung pada prevalensi. Kedua jenis berguna, tetapi keduanya memiliki sifat yang sangat berbeda.

2.2.11 *Rapid Miner*

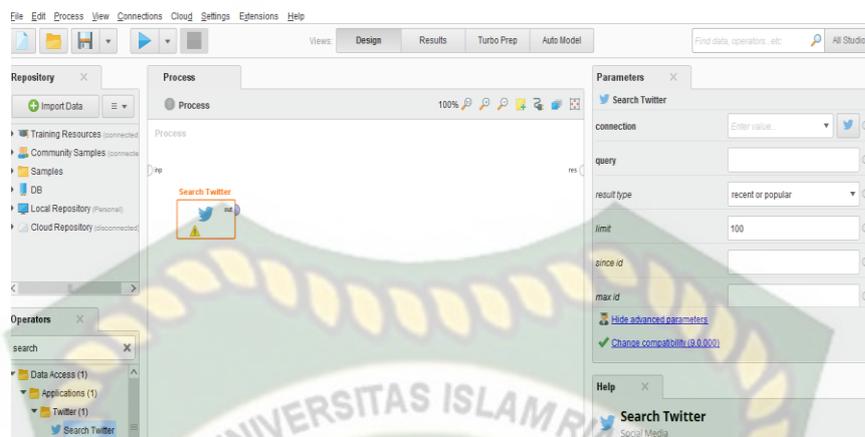
Rapid miner merupakan perangkat lunak yang bersifat terbuka. Rapid miner adalah sebuah solusi untuk melakukan analisis terhadap *data mining*, *text mining* dan analisis prediksi. Rapid miner menggunakan berbagai teknik deskriptif dan

prediksi dalam memberikan wawasan kepada pengguna sehingga dapat membuat keputusan yang paling baik. Rapid miner memiliki kurang lebih 500 operator *data mining*, termasuk operator untuk *input*, *output*, *data preprocessing* dan visualisasi. Rapidminer merupakan *software* yang berdiri sendiri untuk analisa data dan sebagai mesin *data mining* yang dapat diintegrasikna pada produknya sendiri.

Rapid miner menyediakan GUI (*Graphic User Interface*) untuk merancang sebuah *pipeline analitis*. GUI ini akan menghasilkan file XML(*Extensible Markup Language*) yang mendefenisikan proses analitis keinginan pengguna untuk diterapkan ke data. File ini kemudian dibaca oleh *rapidminer* untuk menjalankan analisa secara otomatis.

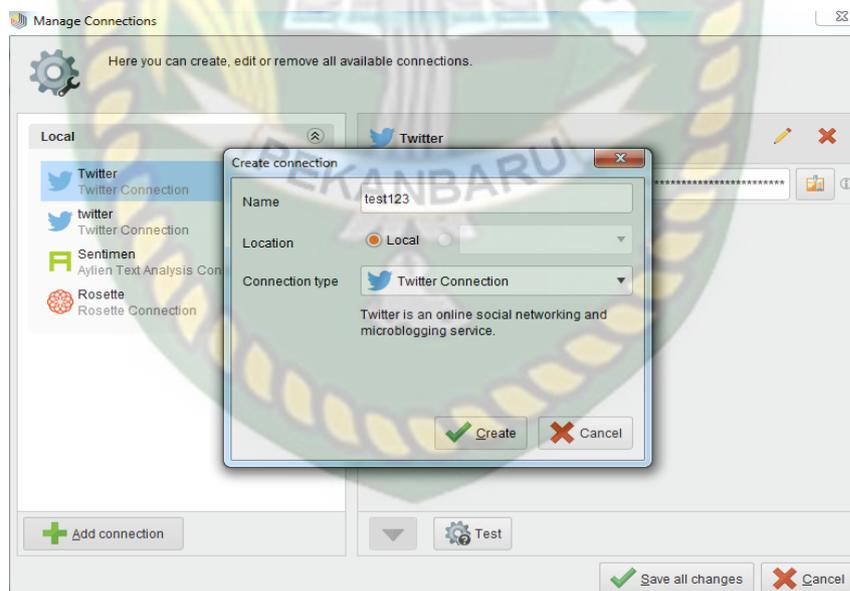
Peneliti menggunakan *rapid miner* sebagai aplikasi pengambilan data. Karna rapid miner sangat memudahkan dalam mengambil data, lebih efektif dan efisien. Data dapat diambil sebanyak yang diinginkan, dan dapat diatur batas waktu data yang dibutuhkan. *Rapid miner* harus terkoneksi dengan twitter terlebih dahulu. Data yang diambil akan tersimpan dalam bentuk excel. Berikut proses pengambilan data pada rapid miner.

1. Seret operator **search twitter** ke dalam **panel proses**. Pilih *connection* untuk mengkoneksi kan ke twitter.



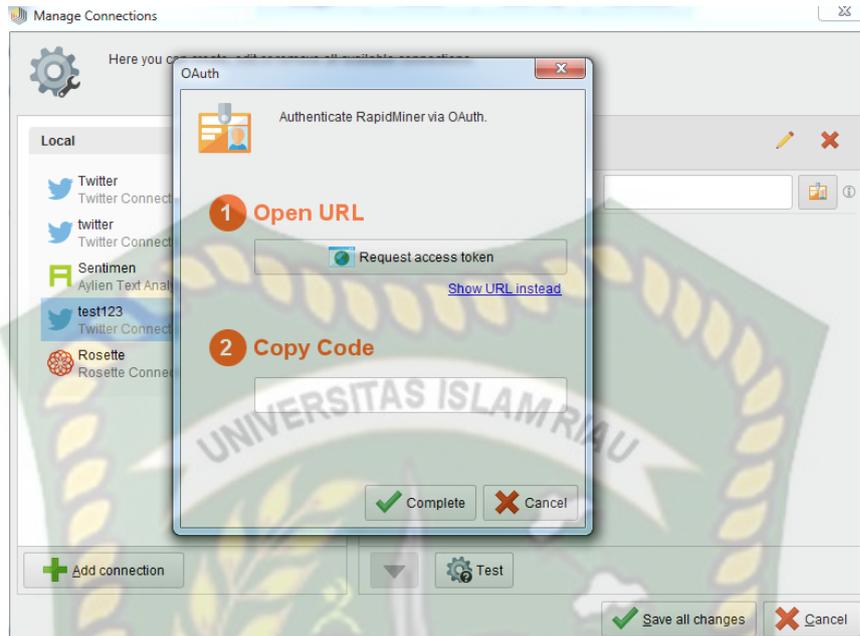
Gambar 2. 1 Search Twitter

2. Klik *add connection*, lalu ganti nama sesuai dengan yang diinginkan sebagai contoh test123. Setelah itu klik button create.



Gambar 2. 2 Search Code Twitter

3. Klik *button access token*, lalu klik *request access token* untuk login ke twitter. Setelah login, rapid miner akan meminta izin untuk mengakses akun yang sudah login seperti pada gambar.



Gambar 2. 3 Open URL

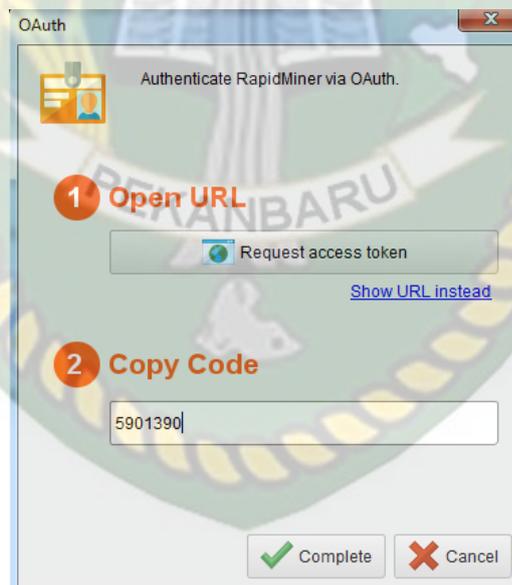


Gambar 2. 4 Izin Akses

4. Setelah akses diizinkan, maka akan keluar kode api twitter untuk bisa mengakses twitter dari rapid miner. Lalu copy kode tersebut ke dalam rapid miner seperti pada gambar 1.6

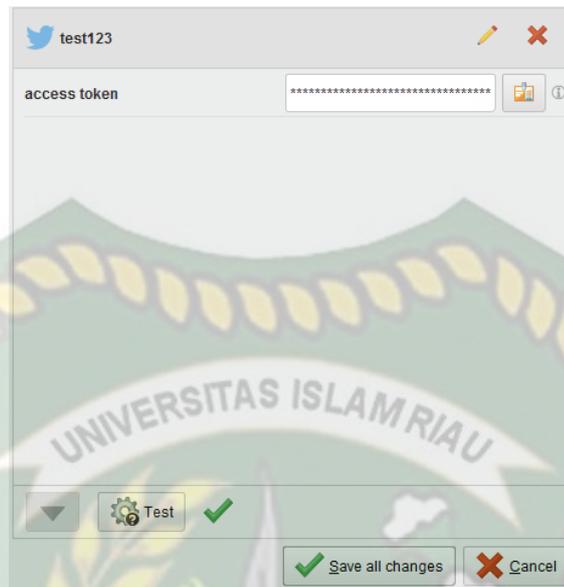


Gambar 2. 5 Code API Twitter



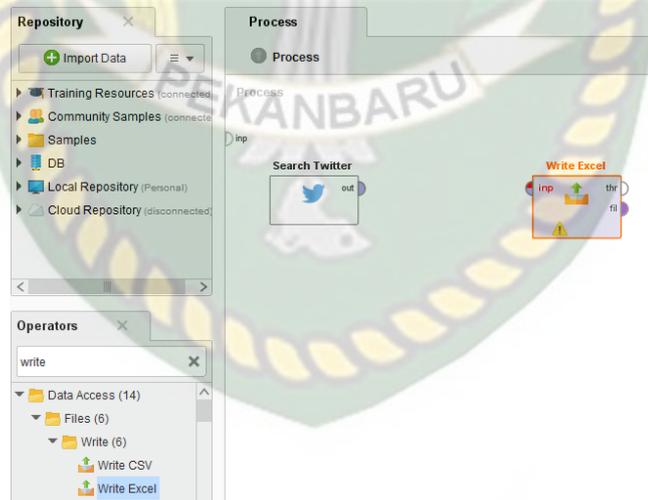
Gambar 2. 6 Copy Code API

5. Test code yang sudah diinputkan. Jika ada tanda centang hijau, berarti code terverifikasi. Selanjutnya klik button save all changes.



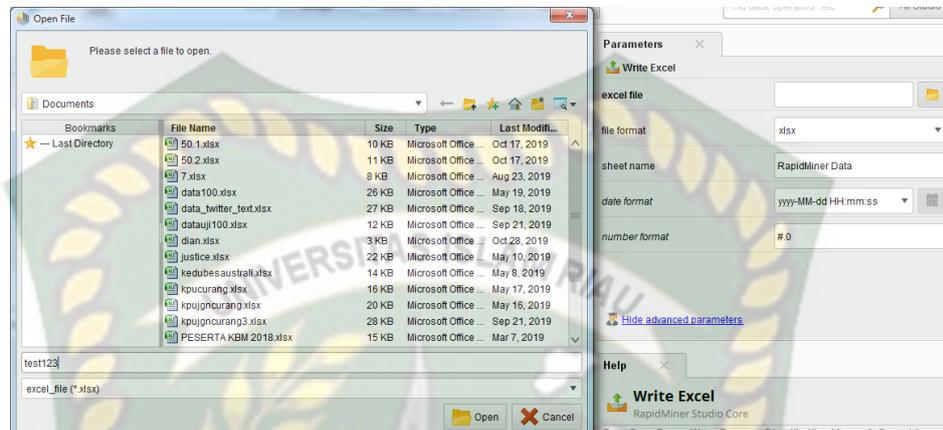
Gambar 2. 7 Test Code API

6. Seret write excel ke dalam panel proses.



Gambar 2. 8 Write Excel

- Klik excel file pada write excel yang sudah diletakkan pada panel proses, lalu simpan data yang akan diambil di dalam folder yang diinginkan.



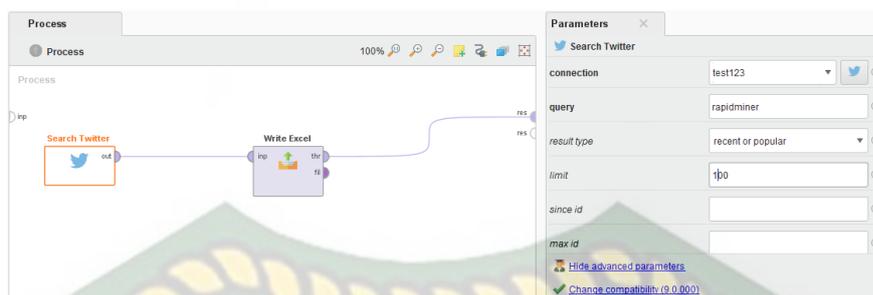
Gambar 2. 9 Simpan Folder Excel

- Hubungkan antara search twitter dengan write excel. Lalu write excel dengan result.



Gambar 2. 10 Search Twitter dan Write Excel

- Inputkan pada kolom query topik yang ingin diambil pada twitter, dengan range banyak data yang diambil, dan language in karna data yang diambil bahasa indonesia.



Gambar 2. 11 Koneksi Twitter

10. Klik button start, maka akan keluar hasil dari query yang sudah diinputkan pada *search* twitter.

The image shows the 'Result History' window in Orange software. It displays a table of search results for the query 'rapidminer'. The table has 12 rows and 11 columns: Row No., Id, Created-At, From-User, From-User-Id, To-User, To-User-Id, Language, Source, Text, and Geo-Lo. The results show various tweets mentioning 'rapidminer' or 'RapidMiner'.

Row No.	Id	Created-At	From-User	From-User-Id	To-User	To-User-Id	Language	Source	Text	Geo-Lo
1	1123919326...	May 2, 2019 1...	Scott Gencer	37734959	var143	151528259	en	<a href="http://...	@var143 @R...	?
2	1123836245...	May 2, 2019 8...	DataScience...	3106817804	?	-1	en	<a href="http://...	RT @RapidM...	?
3	1123714639...	May 2, 2019 1...	Mario Peshev	339991963	?	-1	en	<a href="http://...	RT @RalfKlin...	?
4	1123699403...	May 1, 2019 1...	Ralf Klinkenb...	123850117	?	-1	en	<a href="http://...	RT @RapidMin...	?
5	1123696062...	May 1, 2019 1...	RapidMiner	104211492	var143	151528259	en	<a href="http://...	@var143 @in...	?
6	1123694616...	May 1, 2019 1...	Renie Almira	1122608327	?	-1	en	<a href="http://...	RT @CRN: T...	?
7	1123693983...	May 1, 2019 1...	Varun	151528259	RapidMiner	104211492	en	<a href="http://...	@RapidMiner...	?
8	1123693459...	May 1, 2019 1...	RapidMiner	104211492	var143	151528259	en	<a href="http://...	@var143 @in...	?
9	1123688455...	May 1, 2019 1...	Ralf Klinkenb...	123850117	?	-1	en	<a href="http://...	RT @ai4u_jo...	?
10	1123684913...	May 1, 2019 1...	Machine Lear...	8697941962	?	-1	en	<a href="http://...	RT @CRN: T...	?
11	1123684576...	May 1, 2019 1...	Helios22	135356716	?	-1	en	<a href="http://...	RT @Progra...	?
12	1123682453...	May 1, 2019 1...	Carl M...C...	2277561378	?	-1	en	<a href="http://...	RT @Progra...	?

Gambar 2. 12 Result Data

2.2.12 Python

Python merupakan bahasa pemrograman dinamis yang mendukung pemrograman berbasis objek. *Python* didistribusikan dengan beberapa lisensi yang berbeda dari beberapa versi. Namun pada prinsipnya *python* dapat diperoleh dan dipergunakan secara bebas, bahkan untuk kepentingan komersial. Karena lisensi *python* tidak bertentangan baik menurut definisi *Open Source* maupun *General Public License* (GPK).

Python dikembangkan oleh Guido van Rossum pada tahun 1990-an di CWI, Amsterdam sebagai kelanjutan dari bahasa pemrograman ABC. Versi terakhir

yang dikeluarkan CWI adalah 1.2 tahun 1995, Guido pindah ke CNRI sambil terus melanjutkan pengembangan python. Versi terakhir yang dikeluarkan adalah 1.6 tahun 2000, Guido dan para pengembang inti python pindah ke BeOpen.com yang merupakan sebuah perusahaan komersial dan membentuk *BeOpen pythonlabs* pindah ke *Digital Creations*. Saat ini pengembangan python terus dilakukan oleh sekumpulan pemrograman yang dikoordinir Guido dan *Python Software Foundation*. *Python Software Foundation* adalah sebuah organisasi non-profit yang dibentuk sebagai pemegang hak cipta intelektual python sejak versi 2.1 dan dengan demikian mencegah python dimiliki oleh perusahaan komersial. Saat ini distribusi python sudah mencapai versi 2.7.13 dan versi 3.6.0. Nama python dipilih oleh Guido sebagai nama bahasa ciptaannya karena kecintaan Guido pada acara televisi *Monty Python's Flying Circus*. Oleh karena itu sering kali ungkapan-ungkapan khas dari acara tersebut muncul dalam korespondensi antar pengguna python.

Python merupakan salah satu bahasa pemrograman tingkat tinggi dengan tipe bahasa interpreted karena program-program *Python* langsung dieksekusi oleh interpreter tanpa harus melalui tahap komplikasi. *Python* dapat dijalankan dengan dua cara, yakni (Ema Utami, Suwanto Raharjo, 2004):

1. Mode command-line

Dengan mode ini, program dapat dilakukan dengan memanggil interpreter, kemudian memberi statement *Python* dan interpreter akan menampilkan hasil.

2. Mode script

Mode ini dilakukan dengan cara menuliskan keseluruhan program dalam file, kemudian interpreter akan mengeksekusi seluruh isi dari file. File seperti ini dinamakan script. Sebagai contoh, sebuah file yang ditulis dengan text editor dan diberi nama *script01.py*.

Bahasa ini muncul pertama kali pada tahun 1991, dirancang oleh seorang bernama Guido van Rossum. Sampai saat ini Python masih dikembangkan oleh Python Software Foundation. Bahasa Python mendukung hampir semua sistem operasi, bahkan untuk sistem operasi Linux, hampir semua distronya sudah menyertakan Python di dalamnya.

Hal yang membedakan python dengan bahasa lain adalah dalam hal aturan penulisan kode program. Bahasa yang digunakan juga mendukung hampir disemua sistem operasi, bahkan untuk sistem linux, hampir semua distronya sudah menyertakan python di dalamnya. Dengan kode yang simpel dan mudah diimplementasikan, seorang programmer dapat lebih mengutamakan pengembangan aplikasi yang dibuat. Selain itu python merupakan salah satu produk yang *opensource* juga multiplatform.

2.2.13 Jupyter Notebook

Jupyter notebook seperti namanya, digunakan untuk catatan untuk menunjukkan bagaimana alur programnya berjalan dan berhasil dijalankan. Karena kode program ada di sana dan bisa diutik-utik dan dijalankan ulang, mudah bagi sipembuat untuk memodifikasi jika ada kesalahan atau mau ditambahkan dan juga mudah bagi orang lain untuk memahami alur kerja program tersebut. Fungsi utamanya adalah menjelaskan program secara interaktif dan

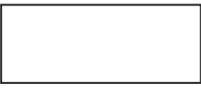
populer untuk ilmu data karena selain butuh kode program yang jalan, biasanya dilengkapi visualisasi yang interaktif juga.

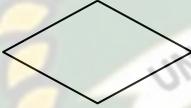
Biasanya di lingkungan penelitian atau di bagian perusahaan yang ada data scientistnya, *Jupyter notebook* cuma digunakan untuk *blueprint* pengolahan data dan sampel visualisasi yang diinginkan jika teknik atau ide tersebut mau diluncurkan kemesin produksi (ada engineer yang menggunakan kode yang ada di *Jupyter notebook* sebagai acuan bagaimana kode di mesin produksi bekerja), fungsi lainnya untuk presentasi menjelaskan ke orang lain bagaimana program bekerja.

2.2.14 Flowchart

Flowchart adalah penyajian yang sistematis tentang proses dan logika dari kegiatan penanganan informasi atau penggambaran secara grafik dari langkah-langkah dan urutan prosedur dari suatu program (Anharku, 2009). *Flowchart* membantu pekerjaan seorang analis dalam memecahkan masalah menjadi bagian-bagian yang lebih kecil dan menolong dalam menganalisis cara-cara lain dalam pengoperasian.

Tabel 2. 2 Simbol *Flowchart*

No.	Simbol	Nama	Fungsi
1		Terminator	Permulaan atau akhir program.
2		Proses	Proses pengolahan data.
3		Garis Alir	Arah aliran program.

4		Preparation	Proses inisialisasi.
5		Input/Output Data	Proses Input/Output data, parameter, informasi.
6		Predefined Process (Sub Program)	Permulaan sub program/ proses menjalankan sub program
7		Decision	Perbandingan pernyataan, penyeleksian data yang memberikan pilihan untuk langkah selanjutnya.
8		Off page Connector	Penghubung bagian-bagian flowchart yang berada pada halaman berbeda.

Dokumen ini adalah Arsip Miik :

Perpustakaan Universitas Islam Riau

BAB III

METODOLOGI PENELITIAN

3.1 Metodologi Penelitian

3.1.1 Metode Penelitian

Metode penelitian merupakan tahapan-tahapan yang dilalui oleh peneliti untuk memperoleh gambaran yang jelas mengenai penelitian. Tahapan yang dilalui dalam metode penelitian ini adalah sebagai berikut:

1. Pengumpulan Data

Data yang dikumpulkan yaitu data yang diambil dari twitter, tweet yang memiliki tagar sesuai dengan topik yang diangkat yaitu kpu jangan curang. Data tersebut diperoleh dengan cara melakukan pencarian data di twitter melalui rapid miner, aplikasi yang bisa terhubung ke twitter untuk mengambil data sesuai topik yang diinginkan.

2. Studi Literatur

Studi literatur dilakukan dengan cara mengumpulkan dan mempelajari segala macam informasi yang berhubungan dengan Analisis Sentimen Pada Twitter Menggunakan Metode Naïve Bayes. Pengambilan data menggunakan rapid miner yang bisa langsung terhubung dengan twitter. Pada penelitian ini, data yang diambil tidak berdasarkan waktu.

3. Perancangan

Pada tahap ini dilaksanakan perancangan terhadap perangkat lunak yang akan dibuat berdasarkan hasil studi literatur yang ada. Perancangan perangkat lunak ini meliputi desain struktur data, desain aliran informasi, desain algoritma dan pemrograman. Perancangan ini dilakukan dengan membuat alur program, menentukan algoritma yang sesuai agar dapat berjalan dengan baik dan sesuai dengan tujuan yang akan dicapai.

4. Implementasi

Tahap implementasi dilakukan secara bertahap dengan acuan studi literatur dan perancangan yang telah dibuat. Perancangan tersebut akan diimplementasikan pada bahasa pemrograman yang telah disepakati.

5. Pengujian dan Evaluasi

Pada tahap ini dilakukan uji coba untuk mencari permasalahan yang mungkin terjadi, mengevaluasi jalannya sistem dan melakukan perbaikan apabila dibutuhkan.

6. Penyusunan Laporan Penelitian

Penyusunan laporan dilakukan pada tahap akhir sebagai dokumentasi. Dokumentasi ini dibuat untuk mempermudah orang lain dalam mempelajari dan mengembangkan sistem lebih lanjut.

3.1.2 Spesifikasi Perangkat Keras (*Hardware*)

Personal komputer atau laptop digunakan untuk perancangan dengan spesifikasi *hardware* sebagai berikut :

1. *Processor* : Intel Pentium Inside™ i5-6200U 2.3GHz up to 2.8GHz

2. *RAM* : 8.00 GB

3. *System Type* : 32-bit *Operating Syatem*

3.1.3 Spesifikasi Perangkat Lunak (*Software*)

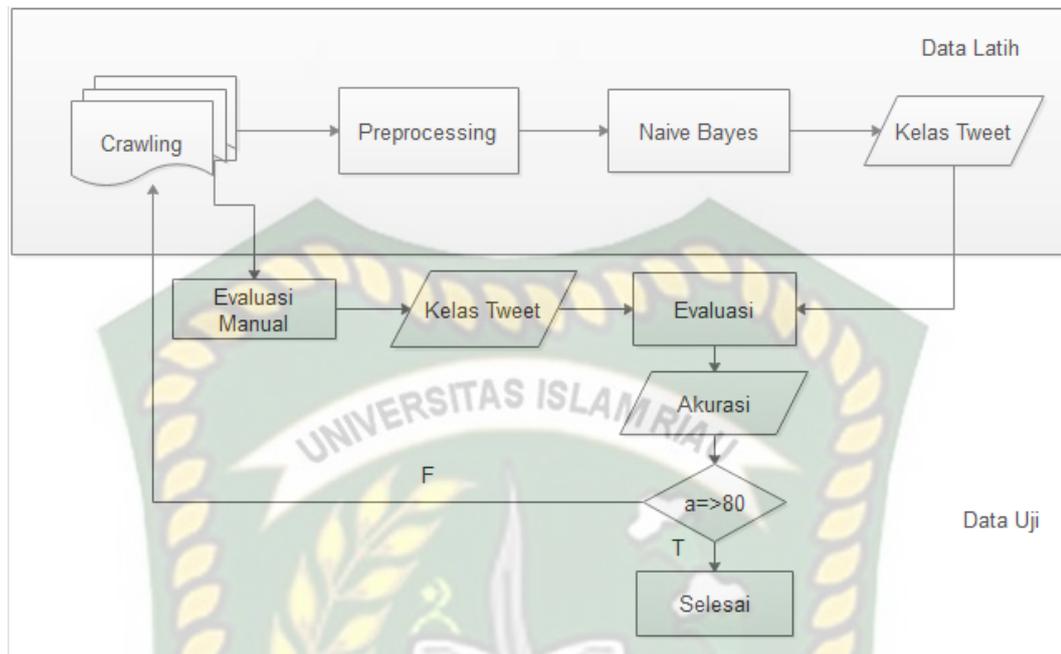
Perangkat lunak yang digunakan dalam pembuatan analisis sentimen pada twitter adalah sebagai berikut:

1. Sistem Operasi : Windows 7 Home
2. Bahasa Pemograman : Python
3. Design Logika Program : Edraw Max
4. Chrome dan Jupyter Notebook

3.2 Perancangan Sistem

3.2.1 *Gambaran Umum*

Analisis Sentimen yang akan dibangun dapat digambarkan secara detail melalui perancangan sistem yang bisa dilihat pada gambar 3.1.



Gambar 3.1 Gambaran Umum Analisis Sentimen

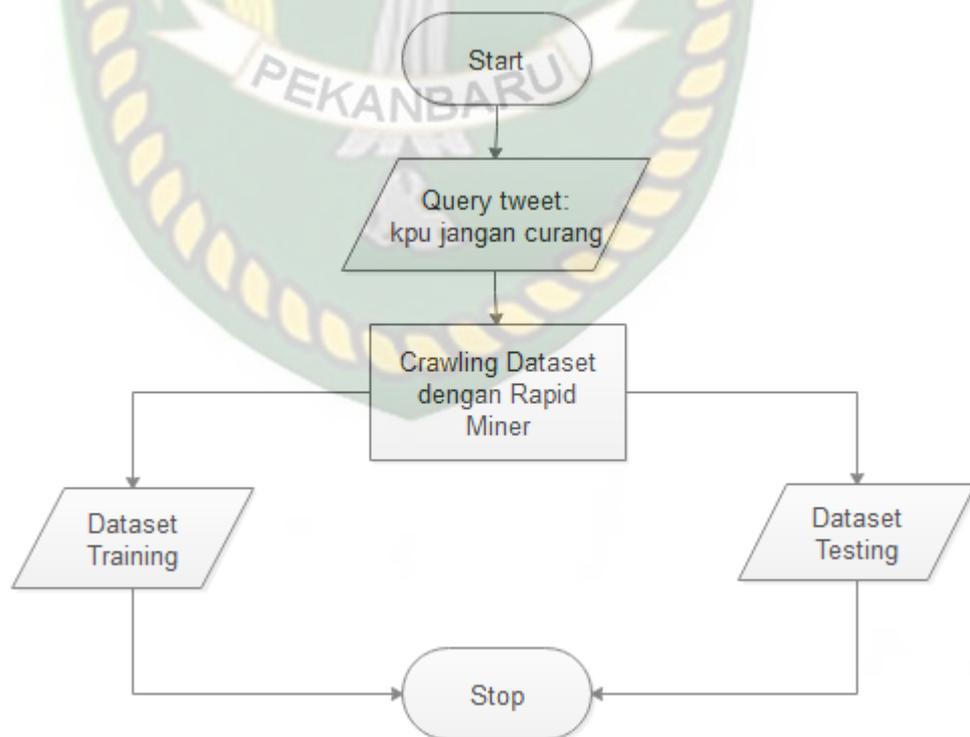
Berikut penjelasan dari gambar 3.1:

1. Crawling data: Mengambil data dari twitter menggunakan rapid miner dengan tagar #kpujangancurang. Data yang diambil berupa data latih dan data uji.
2. Kemudian akan melakukan *preprocessing* terhadap data latih, yaitu dengan melakukan *stopwords*.
3. Setelah *preprocessing* didapatkan data bersih yang kemudian akan diterapkan metode *naïve bayes*.
4. Setelah menerapkan metode *naïve bayes*, didalam metode tersebut akan melakukan perhitungan bobot pada sentimen yang nanti nya akan menghasilkan kelas *tweet*.

5. Tahap selanjutnya setelah kelas *tweet* ditentukan adalah evaluasi.
6. Kelas *tweet* yang dilakukan secara manual dari data uji di evaluasi kembali dengan data latih yang sudah terproses.
7. Baru setelah itu mendapatkan akurasi yang tepat pada data.
8. Jika tingkat akurasi mencapai 80% , maka proses selesai.
9. Jika tidak mencapai 80% , maka mengulang ke proses pengambilan data.

3.2.2 Proses Dataset

Dataset berupa teks berbahasa Indonesia yang diambil dari twitter. Data yang diambil untuk penelitian ini menggunakan query ‘kpu jangan curang’. *Query* tersebut merupakan tagar yang sedang *trending topic* di twitter. Diagram alur proses pengambilan *dataset* seperti pada gambar 3.2.



Gambar 3. 2 Diagram Alir *Dataset*

Dataset dari hasil *crawling* ini akan dibagi menjadi dua bagian yaitu data *training* (data latih) dan data *testing* (data uji) yang dipresentasikan pada gambar 3.2. Data latih diklasifikasikan menggunakan naïve bayes dengan label sentimen negatif dan positif. Sedangkan data uji diklasifikasikan secara manual dengan label sentimen negatif dan positif. Data uji nantinya akan digunakan pada saat evaluasi untuk menentukan keakuratan data pada sistem. Pengambilan data *training* dan data *testing* dilakukan dengan waktu pengambilan yang berbeda.

3.2.3 Preprocessing

Proses *preprocessing* merupakan hal yang penting untuk tahap selanjutnya, yaitu mengurangi atribut yang kurang berpengaruh terhadap proses klasifikasi. Data yang dimasukkan pada tahap ini masih berupa data mentah yang masih kotor, sehingga hasil dari proses ini adalah dokumen berkualitas yang diharapkan mempermudah dalam proses klasifikasi. Berikut proses *preprocessing* yang terjadi:

Tabel 3. 1 Proses Preprocessing

No	Sebelum	Sesudah
1.	Benar juga, kpu yang membuat rakyat resah. Aduh kejamnya kecurangan.	Benar juga kpu yang membuat rakyat resah Aduh kejamnya kecurangan

2.	Benar juga kpu yang membuat rakyat resah Aduh kejamnya kecurangan	benar juga kpu yang membuat rakyat resah aduh kejamnya kecurangan
3.	benar juga kpu yang membuat rakyat resah aduh kejamnya kecurangan	benar juga kpu yang membuat rakyat resah aduh kejamnya kecurangan
4.	benar juga kpu yang membuat rakyat resah aduh kejamnya kecurangan	benar kpu membuat rakyat resah kejamnya kecurangan

Berikut penjelasan dari tabel 3.1 pada proses preprocessing :

1. Penghilangan tanda baca, tahap ini sebagai proses awal supaya didapatkan teks yang murni berisi kata-kata agar pemrosesan selanjutnya menjadi lebih mudah.
2. *Case Folding*, pada tahap ini dilakukan proses perubahan dan huruf besar menjadi huruf kecil.
3. *Tokenizing*, pada setiap kata akan dipisahkan berdasarkan spasi yang ditentukan.
4. *Filtering*, yaitu pembuangan kata-kata tidak penting dan daftar *stopwords* bahasa indonesia.

3.2.4 Term Frequency – Inverse Document Frequency (TF-IDF)

Setelah melakukan *preprocessing*, langkah selanjutnya adalah pembobotan kata menggunakan perhitungan tf-idf. Tfidf adalah cara pemberian bobot hubungan suatu kata (*term*) terhadap kata. Untuk kata tunggal tiap kalimat dianggap sebagai dokumen. Berikut contoh perhitungan tf-idf.

Tabel 3. 2 Contoh Dokumen

Doc A	Jangan ancam rakyat, rakyat indonesia pintar
Doc B	Rakyat tidak pernah gagal bernegara, pemerintah yang gagal bernegara
Doc C	Suara rakyat dicuri, bagaimana uang rakyat

1. Dari contoh dokumen pada tabel 3.2, hitung frekuensi kemunculan kata pada tiap dokumen sebagai langkah awal.
2. Setelah didapatkan frekuensi kata pada tiap dokumen, langkah selanjutnya adalah menghitung tf dengan rumus $tf = 0.5 * 0.5 \frac{ft,d}{\max(ft,d)}$, misalkan pada doc A pada kata ancam maka $tf = 0.5 * 0.5 \left(\frac{1}{1}\right) = 1$
3. Untuk mencari nilai **IDF** adalah dengan rumus $\log \left(\frac{N}{dft}\right)$, maka $\log 3/1 = 0.477$
4. Terakhir dengan mencari nilai hasil bobot setiap dokumen berdasarkan hasil dari tf-idf yaitu dengan $tf \times idf$, maka $1 * 0.477 = 0.477$. Berikut hasil dari perhitungan tf idf .

Tabel 3. 3 TF-IDF Doc A

Kata	Frekuensi Kemunculan Kata			TF ($0.5*0.5$ ($ft,d/\max(ftd)$))	IDF ($\log(\frac{N}{dft})$)	W (TF*IDF)
	A	B	C	A	A	A
Ancam	1	0	0	1	0.477	0.477
Bernegara	0	2	0	0	0.176	0
Gagal	1	0	0	0	0.176	0
Jangan	0	1	0	1	0.477	0.477
Rakyat	0	1	0	0.4	-0.2218	-0.0887
Indonesia	0	2	0	1	0.477	0.477
Pintar	2	1	2	1	0.477	0.477
Tidak	1	0	0	0	0.477	0
Pernah	0	1	0	0	0.477	0
Pemerintah	1	0	0	0	0.477	0
Dicuri	0	0	1	0	0.477	0
Bagaimana	0	0	1	0	0.477	0
Uang	0	0	1	0	0.477	0

Tabel 3. 4TF-IDF Doc B

Kata	Frekuensi Kemunculan Kata			TF ($0.5*0.5$ ($ft,d/\max(ftd)$))	IDF ($\log(\frac{N}{dft})$)	W (TF*IDF)
	A	B	C	B	B	B
Ancam	1	0	0	0	0.477	0

Bernegara	0	2	0	1	0.176	0.176
Gagal	1	0	0	1	0.176	0.176
Jangan	0	1	0	0	0.477	0
Rakyat	0	1	0	0.2	-0.2218	-0.044
Indonesia	0	2	0	0	0.477	0
Pintar	2	1	2	0	0.477	0
Tidak	1	0	0	1	0.477	0.477
Pernah	0	1	0	1	0.477	0.477
Pemerintah	1	0	0	1	0.477	0.477
Dicuri	0	0	1	0	0.477	0
Bagaimana	0	0	1	0	0.477	0
Uang	0	0	1	0	0.477	0

Tabel 3. 5 TF-IDF Doc C

Kata	Frekuensi Kemunculan Kata			TF ($0.5 \cdot 0.5$ ($\frac{ft,d}{\max(ftd)}$))	IDF ($\log(\frac{N}{aft})$)	W (TF*IDF)
	A	B	C	C	C	C
Ancam	1	0	0	0	0.477	0
Bernegara	0	2	0	0	0.176	0
Gagal	1	0	0	0	0.176	0
Jangan	0	1	0	0	0.477	0
Rakyat	0	1	0	0.4	-0.2218	-0.0887
Indonesia	0	2	0	0	0.477	0
Pintar	2	1	2	0	0.477	0

Tidak	1	0	0	0	0.477	0
Pernah	0	1	0	0	0.477	0
Pemerintah	1	0	0	0	0.477	0
Dicuri	0	0	1	1	0.477	0.477
Bagaimana	0	0	1	1	0.477	0.477
Uang	0	0	1	1	0.477	0.477

3.2.5 Multinomial Naive Bayes

Model multinomial memperhitungkan frekuensi setiap kata yang muncul pada dokumen. Contoh pada perhitungan berikut yaitu data yang ada pada tabel 3.2. Doc A memiliki kelas positif, doc B dan C adalah kelas negatif. Dari data TF-IDF sebelumnya, didapatkan hasil jumlah keseluruhan W untuk kelas positif adalah $W(+)=1.819$, jumlah keseluruhan W pada kelas negatif adalah $W(-)=2.127$, dan jumlah keseluruhan idf pada semua kelas adalah $B'=14.7$. Berikut bobot tfidf yang sudah di dapatkan sebelumnya beserta kelas nya masing-masing.

Tabel 3. 6 Tabel Bobot

Kata	W	
	Positif (C1)	Negatif (C2)
Ancam	0.477	0
Bernegara	0	0.176
Gagal	0	0.176
Jangan	0.477	0

Rakyat	-0.0887	-0.132
Indonesia	0.477	0
Pintar	0.477	0
Tidak	0	0.477
Pernah	0	0.477
Pemerintah	0	0.477
Dicuri	0	0.477
Bagaimana	0	0.477
Uang	0	0.477

Dari tabel bobot, langkah selanjutnya adalah menghitung probabilitas prior dengan rumus $P(c_i) = \frac{N_{c_i}}{N}$ adalah sebagai berikut:

$$P(c_1) = \frac{1}{3} = 0.333$$

$$P(c_2) = \frac{2}{3} = 0.666$$

Ket:

$P(C_i)$: probabilitas kelas

N_{c_i} : jumlah kelas c pada seluruh dokumen

N : jumlah seluruh dokumen

Setelah didapatkan probabilitas, hitung likelihood pada setiap kata dengan

rumus $P(w|c) = \frac{W_{ct}+1}{(\sum_{w' \in V_{W_{ct'}}} w') + B'}$, berikut contoh perhitungan pada kata ancam :

$$\text{ancam (c1)} = \frac{W_{ct}+1}{(\sum_{w' \in V_{W_{ct'}}} w') + B'} = \frac{0.477+1}{1.819+14.7} = \frac{1.477}{16.519} = 0.089$$

$$\text{ancam (c2)} = \frac{W_{ct}+1}{(\sum_{w' \in V_{W_{ct'}}} w') + B'} = \frac{0+1}{2.127+14.7} = \frac{1}{16.827} = 0.059$$

Ket:

W_{ct} : W dari *term* t di kategori c

$\sum_{w' \in V_{W_{ct'}}$: jumlah total W dari keseluruhan term yang berada di kategori c

B' : nilai idf pada seluruh dokumen

Hasil dari perhitungan likelihood dapat dilihat pada tabel 3.7 di bawah ini.

Tabel 3. 7Likelihood

P(w c)	Positif	Negatif
Ancam	0.089	0.059
Bernegara	0.06	0.069
Gagal	0.06	0.069
Jangan	0.089	0.059
Rakyat	0.055	0.05
Indonesia	0.089	0.059

Pintar	0.089	0.059
Tidak	0.06	0.087
Pernah	0.06	0.087
Pemerintah	0.06	0.087
Dicuri	0.06	0.087
Bagaimana	0.06	0.087
Uang	0.06	0.087

Lalu hitung V_{MAP} positif dan negatif pada kalimat dari data testing untuk menentukan apakah kalimat tersebut bernilai positif atau negatif. Sebagai contoh pada kalimat “Rakyat Indonesia Dicuri”, berikut perhitungannya:

$$V_{MAP}(c1) = P(\text{rakyat}) * P(\text{indonesia}) * P(\text{dicuri}) * P(C)$$

$$= 0.055 * 0.089 * 0.06 * 0.333$$

$$= 0.000097$$

$$V_{MAP}(c2) = P(\text{rakyat}) * P(\text{indonesia}) * P(\text{dicuri}) * P(C)$$

$$= 0.055 * 0.059 * 0.087 * 0.666$$

$$= 0.003418$$

Selanjutnya menentukan determinan dari kalimat tersebut adalah dengan cara membandingkan $V_{MAP}(c1)$ yang merupakan kelas positif dan $V_{MAP}(c2)$ yang merupakan kelas negatif. Pada kalimat “ rakyat indonesia dicuri” nilai $V_{MAP}(c1)$ lebih tinggi dari pada nilai $V_{MAP}(c2)$, maka determinan dari kalimat tersebut adalah negatif.



BAB IV

HASIL DAN PEMBAHASAN

4.1 Model Analisi Sentimen

Pada penelitian ini, penulis menggunakan jupyter notebook sebagai wadah dalam pembuatan analisis sentimen karena jupyter memiliki interface yang sederhana dan mudah untuk dipahami. Jupyter memiliki kelebihan dapat memperlihatkan alur program yang berhasil dijalankan atau error pada alur program tertentu. Berikut tahapan-tahapan analisis sentimen.

4.1.1 Pelatihan

Proses pelatihan bertujuan untuk membentuk model prediksi dari data yang telah diketahui kelasnya. Pada kasus ini penulis menggunakan data latih sebanyak 200. Awalnya yang digunakan hanya 100 data, tetapi akurasi yang dihasilkan oleh sistem untuk 100 data latih hanya 0.6. Lalu penulis menambahkan 100 data lagi menjadi 200 data untuk digunakan sebagai data latih dan akurasi yang dihasilkan adalah 0.807. Angka yang cukup bagus untuk digunakan sebagai data latih. Terdapat dua jenis kelas pada data latih dikasus ini yaitu “POSITIF” dan “NEGATIF”. Berikut contoh tampilan data latih pada sistem pada gambar 4.1.

Out [14] :

	Sentimen	Tweet
0	0	Berita mengenai dagelan pilpres 2019 di bawah ...
1	0	Ayolh pak rkyat indonesia jangan kira masih go...
2	0	Jangan ancam rakyat !!! TNI/POLRI dari rakyat ...
3	0	Polisi Membubarkan Massa Malam ini Praya Lombo...
4	0	Wakil Ketua PKS Sumsei Askweni menuntut Bawasl...
5	0	#KpuJanganCurang #JihadLawanKecurangan HASIL ...
6	0	@VIVAcoid Tolong diusut Bapak2, @KPU_RI keracu...
7	0	@rmlco Semoga laknat Allah bersama org2 yg cu...
8	0	Panik itulah yg tepat utk kubu 01,mereka sudah...
9	0	Pemilu 2019 merupakan pesta demokrasi paling k...

Gambar 4. 1 Tampilan Data Latih

Sentimen pada tabel yang terdapat pada gambar 4.1 adalah kelas dari kalimat dimana 0 adalah negatif dan 1 adalah positif. Sedangkan *tweet* adalah kalimat atau *tweet* yang telah diambil dari twitter. Data yang diinputkan berupa data dari excel dengan format csv. Gambar 4.1 adalah 10 contoh data yang ditampilkan dari data yang ada.

Pada tahapan terakhir setelah melakukan pengujian dengan data ujian melakukan perhitungan *evaluation measure*, data latih ditambahkan dengan data uji untuk melihat perbedaan akurasi data latih yang dihasilkan oleh mesin dari akurasi sebelumnya. Maka data latih terakhir berjumlah 300 data. Dari 300 data ini dihasilkan akurasi 0.882 oleh sistem, sedangkan akurasi yang dihasilkan

dengan 200 data adalah 0.807. Perbedaan yang dihasilkan tidak begitu jauh karena data yang digunakan memiliki kata yang tidak berbeda jauh dengan data sebelumnya. Sehingga mesin tidak mempelajari kata terlalu banyak karena sudah dipelajari pada data sebelumnya.

4.1.2 Pengujian Dengan Data Uji

Pengujian dengan data uji adalah proses puncak pada analisis. Pada proses ini *user* akan melakukan penginputan data dan memperoleh hasil berupa kelas dalam bentuk binner yaitu “0” sebagai “negatif” dan “1” sebagai “positif”. Hasil ini didapatkan oleh sistem melalui pembelajaran dari data latih yang sudah diinputkan. Untuk menginputkan data uji, dilakukan secara manual ditempat yang telah disediakan.

```
In [16]: data_uji=np.array([ "Mantap...berjuanglah selama masih hidup...klo udh mati ga bisa ngapa ngapain",
"@do_ra_dong Klo ga mau dimaki mundur!",
"Klo gk tahan sama tekanan mending keluar aja mas byarkan orang2 yg militan berjuangvuntuk prabowo sandi",
"Aceh Luar Biasa,Kata Jendral .. Selamat datang di Tanah Aceh kami Jendral, anda adalah presiden kami",
"apapun bacodnya ttp I LOVE YOU HABIB RIZIQ SIHAB #DoaKamiUntukPrabowoSandi",
"RT @dongdong4: Insya Allah siap !!!",
"Insya Allah siap !!!",
"RT @BARET_Riau: Aksi Damai ""Memperjuangkan Kedaulatalan Rakyat"" Di Prov. Riau, 17 Mei 2019",
"Beda tujuan. Yg ikhlas sama yg pamrih",
"Aksi Damai ""Memperjuangkan Kedaulatalan Rakyat"" Di Prov. Riau, 17 Mei 2019",
])
data_uji_vector = vectorizer.transform(data_uji)
print (clf.predict(data_uji_vector))
[1 1 0 1 1 1 1 1 0 1]
```

Gambar 4. 2 Tampilan *Form* Input Data Uji

Hasil dari sistem ini di uji dengan hasil yang didapatkan oleh user dengan cara yang manual untuk mendapatkan ketepatan prediksi. Hasil dari uji manual yang sudah dilakukan oleh *user* diinputkan pada sistem secara manual. Disini penulis menggunakan 50 data baru yang diambil secara *random*. Untuk menghasilkan ketepatan prediksi, penulis menggunakan pengujian *whitebox* untuk menghasilkan tabel prediksi antara uji pada sistem dan uji manual. Pengujian *whitebox* adalah salah satu metode pengujian perangkat lunak yang berfokus pada sisi *source code* yang digunakan, yaitu apakah algoritma program sudah berjalan dengan baik dan menghasilkan *output* yang dikehendaki.

Input : uji manual

Output: ketepatan prediksi

1. Input uji manual
2. Prediksi = uji system
3. For uji (ujimanual)
4. if ujimanual [uji] == prediksi [uji]
5. print "Tepat"
6. else
7. print "Tidak Tepat"
8. end if
9. end for
10. end

Hasil dari algoritma yang ada pada *source code* akan menampilkan tabel ketepatan prediksi seperti pada gambar 4.3 di bawah ini.

Out [116]:

	Tweet	Sentimen	Uji Manual	Ketepatan Prediksi
0	Indonesia sedang sakit!	0	0	Tepat
1	banyak banget yang post kecurangan	0	0	Tepat
2	Isu kematian massal petugas KPPS 2019 semakin ...	0	0	Tepat
3	Jokowi's government opposition was arrested by...	1	0	Tidak Tepat
4	Indonesia sedang sakit!	0	0	Tepat
5	Guru Besar Ekonomi UI: Indonesia Krisis Konsti...	0	0	Tepat
6	Suara rakyat aja kau curi, bgmn uang rakyat	0	0	Tepat
7	Polisi Membubarkan Massa Malam ini Praya Lombo...	0	0	Tepat
8	Ini orang harus yg paling duluan ditangkap	0	0	Tepat
9	Suara rakyat aja kau curi, bgmn uang rakyat br...	0	0	Tepat

Gambar 4.3 Tampilan Ketepatan Prediksi

Sentimen merupakan hasil uji dari sistem sedangkan uji manual merupakan hasil uji dari user. Jika hasil uji dari sistem sama dengan hasil uji dari *user* maka prediksi nya adalah tepat, dan jika hasil uji tidak sama maka prediksi nya tidak tepat. Dari pengujian yang penulis lakukan, dapat ditarik kesimpulan bahwa hasil pengujian *whitebox* ini menunjukkan hasil yang sesuai dengan yang diharapkan.

4.2 Evaluation Measure

Evaluation measure atau evaluasi performansi dilakukan untuk menguji hasil klasifikasi dengan mengukur nilai kebenaran dari sistem. Ketika dataset memiliki hanya dua kelas, salah satu sering dianggap sebagai positif dan yang lain sebagai negatif. Dalam kasus ini entri dalam dua baris dan kolom confusion

matrix dirujuk sebagai *true and false positives* dan *true and false negatives*, seperti pada tabel 4.1 di bawah ini:

Tabel 4. 1 Tabel *Confusion Matrix*

		Predicted	
		Positive	Negative
Actual	Positive	TP	FN
	Negative	FP	TN

True Positive (TP) adalah jumlah record positif yang diklasifikasikan sebagai positif, *false positive* (FP) adalah jumlah record negatif yang diklasifikasikan sebagai positif, *false negative* (FN) adalah jumlah record positif yang diklasifikasikan sebagai negatif, *true negative* (TN) adalah jumlah record negatif yang diklasifikasikan sebagai negatif. Untuk klasifikasi text, biasanya pengukuran akurasi menggunakan beberapa pengukuran, yaitu *accuracy*, *precision*, *recall*, dan *F-Measure*.

Dalam penelitian ini, penulis menggunakan 100 data uji, dimana percobaan pertama dilakukan dengan 50 data uji yang diberi nama data A, yang kedua dengan 50 data uji lainnya dengan nama data B, dan yang ketiga menggabungkan data uji menjadi 100 data yang disebut sebagai data C.

Tabel 4. 2 Tabel *Confusion Matrix* Data A

	Predicted	
	Positive	Negative
Positive	5	1
Negative	22	22

Tabel 4. 3 Tabel *Confusion Matrix* Data B

	Predicted	
	Positive	Negative
Positive	11	19
Negative	16	4

Tabel 4. 4 Tabel *Confusion Matrix* Data C

	Predicted	
	Positive	Negative
Positive	11	19
Negative	16	4

4.1.3 Accuracy

Accuracy merupakan persentase dokumen yang berhasil diklasifikasikan dengan tepat oleh sistem. *Accuracy* diperoleh dari hasil perhitungan yang ditunjukkan pada persamaan (2.6).

Dari data *testing* A yang sudah diproses ketepatan prediksinya, didapatkan hasil TP : 5, TN : 22, FP : 22, FN : 1. Perhitungan *accuracy* nya adalah sebagai berikut:

$$Accuracy = \frac{5+22}{(5+22+22+1)} = \frac{27}{50} = 0.54$$

Disini penulis menggunakan bilangan bulat karna pada perhitungan ini penulis menentukan *range* untuk hasilnya dengan nilai 1-100 maka hasilnya dikalikan dengan 100, sehingga didapatkan hasil dari *accuracy* antara uji manual dan uji pada sistem adalah 54. Berikut tabel nilai *accuracy* pada tiap data.

Tabel 4. 5 Hasil *Accuracy*

Tipe Data	Nilai <i>Accuracy</i>
Data A	54%
Data B	60%
Data C	57%

Pengujian yang dilakukan pada proses ini menggunakan pengujian *white box*, dengan algoritma sebagai berikut:

Output : Accuracy

1. read tp, tn, fp, fn
2. $accuracy = ((tp + tn) * 100) / \text{float}(tp + tn + fn + fp)$
3. print "accuracy = " , accuracy

Hasil yang diperoleh dari algoritma yang ada pada *source code* adalah sebagai berikut:

```
In [121]: print('Accuracy :', compute_accuracy(tp, tn, fn, fp))
Accuracy : 54.0
```

Gambar 4. 4 Perhitungan *Accuracy* Data A

```
In [134]: print('Accuracy :', compute_accuracy(tp, tn, fn, fp))
Accuracy : 60.0
```

Gambar 4. 5 Perhitungan *Accuracy* Data B

```
In [188]: print('Accuracy :', compute_accuracy(tp, tn, fn, fp))
Accuracy : 57.0
```

Gambar 4. 6 Perhitungan *Accuracy* Data C

Hasil dari pengujian yang dilakukan ini dapat disimpulkan bahwa perhitungan *accuracy* antara data *testing* A,B dan C memiliki hasil yang berbeda dan hasil *accuracy* tertinggi antara tiga pengujian ini adalah pengujian *accuracy* pada data B.

4.1.4 Precision

Precision berguna untuk mengukur tingkat ketepatan antara informasi yang diminta oleh pengguna dengan jawaban yang diberikan oleh sistem, seperti pada persamaan (2.4).

Maka, perhitungan untuk data *testing* A adalah sebagai berikut:

$$Precision = \frac{TP}{TP+FP} = \frac{5}{5+22} = \frac{5}{27} = 0.185$$

Sama seperti *accuracy*, untuk *precision* penulis menggunakan *range* 1-100 dalam menentukan hasil dari perhitungan, maka hasil dikalikan dengan 100 didapatkan *precision* dari uji manual dengan uji sistem adalah 18.5 pada data *testing* A, 40.7 pada data *testing* B, dan 29.6 pada data *testing* C.

Tabel 4. 6 Hasil *Precision*

Tipe Data	Nilai <i>Precision</i>
Data A	18.5%

Data B	40.7%
Data C	29.6%

Proses ini menggunakan pengujian *whitebox*. Pengujian yang dilakukan adalah pada perhitungan *precision*, dengan algoritma pada *source code* seperti ini di bawah ni:

Output : Precision

1. read tp, fp
2. $precision = (tp * 100) / float(tp + fp)$
3. print "precision = " , precision

Hasil dari algoritma tersebut adalah sebagai berikut:

```
In [123]: print('Precision :', compute_precision(tp, fp))
Precision : 18.51851851851852
```

Gambar 4. 7 Perhitungan *precision* Data A

```
In [136]: print('Precision :', compute_precision(tp, fp))
Precision : 40.74074074074074
```

Gambar 4. 8 Perhitungan *precision* Data B

```
In [190]: print('Precision :', compute_precision(tp, fp))
Precision : 29.62962962962963
```

Gambar 4. 9 Perhitungan *precision* Data C

Pada pengujian ini, dapat dilihat bahwa perhitungan *precision* pada data B memiliki nilai tertinggi diantara data lainnya, dan data A memiliki nilai terendah diantara kedua data lainnya.

4.1.5 Recall

Recall berguna untuk mengukur tingkat keberhasilan system dalam menemukan kembali sebuah informasi dan menampilkan informasi berupa analisis sentimen. Berikut persamaan *recall* dapat dilihat pada persamaan (2.5).

Dari persamaan tersebut, dapat dihitung dengan perhitungan di bawah ini:

$$Recall = \frac{TP}{TP+FN} = \frac{5}{5+1} = 0.83$$

Pada perhitungan *recall* penulis juga menggunakan *range* 1-100 dalam mengukur tingkat keberhasilan. Untuk mendapatkan nilai tersebut, hasil dari perhitungan dikali dengan 100 menjadi 83 pada data A, 73 pada data B, dan 76 pada data C.

Tabel 4. 7 Hasil *Recall*

Tipe Data	Nilai <i>Recall</i>
Data A	83%
Data B	73%
Data C	76%

Pada proses ini, pengujian yang dilakukan adalah pengujian *whitebox*. Yang perlu dilakukan *testing* adalah perhitungan untuk recall pada sistem, dengan algoritma pada *source code* sebagai berikut:

Output : Recall

1. read tp, fn
2. $\text{recall} = (\text{tp} * 100) / \text{float}(\text{tp} + \text{fn})$
3. print "recall = ", recall

Hasil dari algoritma tersebut adalah seperti berikut ini:

```
In [125]: print('Recall :', compute_recall(tp, fn))
Recall : 83.33333333333333
```

Gambar 4. 10 Perhitungan *Recall* Data A

```
In [138]: print('Recall :', compute_recall(tp, fn))
Recall : 73.33333333333333
```

Gambar 4. 11 Perhitungan *Recall* Data B

```
In [192]: print('Recall :', compute_recall(tp, fn))
Recall : 76.19047619047619
```

Gambar 4. 12 Perhitungan *Recall* Data C

Dari pengujian white box ini dapat disimpulkan bahwa hasil perhitungan *recall* tertinggi antara data A,B, dan C adalah pada data A.

4.1.6 F – Measure

F-measure akan memiliki nilai yang besar hanya ketika *precision* dan *recall* memiliki nilai yang besar, dan dapat dilihat sebagai cara untuk menemukan penyesuaian terbaik antara *precision* dan *recall*. Range untuk nilai *F-measure* ini adalah 0-1. Persamaan yang digunakan adalah persamaan (2.7).

Perhitungan yang dilakukan menggunakan hasil asli dari *precision* dan *recall*, tanpa dikalikan 100 karna pada perhitungan ini range yang digunakan adalah 0-1. Berikut perhitungan dari persamaan di atas pada data A:

$$F1 = 2 * \frac{0.185 * 0.833}{0.185 + 0.833} = 2 * \frac{0.154}{1.018} = 2 * 0.151 = 0.3$$

Hasil *F-measure* dari masing-masing data dapat dilihat pada tabel 4.8.

Tabel 4. 8 Hasil *F-measure*

Tipe Data	Nilai <i>F-measure</i>
Data A	0.30
Data B	0.52
Data C	0.426

Sama seperti pada perhitungan sebelum nya, pengujian yang digunakan pada perhitungan ini adalah *white box*, dengan algoritma sebagai berikut:

Output : `f1_score`

1. `read precision, recall`
2. `f1_score = (2 * precision * recall) / (precision + recall)`
3. `print "f1_score = ", f1_score`

Hasil dari algoritma di atas adalah sebagai berikut:

```
In [127]: print('F1 :', compute_f1_score(df.prediksi, df.ujimanual))
          F1 : 0.30303030303030304
```

Gambar 4. 13 Perhitungan F-Measure Data A

```
In [140]: print('F1 :', compute_f1_score(df.prediksi, df.ujimanual))
          F1 : 0.5238095238095238
```

Gambar 4. 14 Perhitungan F-Measure Data B

```
In [194]: print('F1 :', compute_f1_score(df.prediksi, df.ujimanual))
          F1 : 0.42666666666666667
```

Gambar 4. 15 Perhitungan F-Measure Data C

Hasil dari pengujian ini dapat dilihat bahwa perhitungan pada data B memiliki hasil lebih tinggi dibandingkan dengan hasil dari data lainnya.

Berikut *barchart* dari *evaluation measure* pada sistem sesuai dengan perhitungan yang sudah dilakukan, dapat dilihat pada diagram 4.1 di bawah ini.

Evaluation Measure

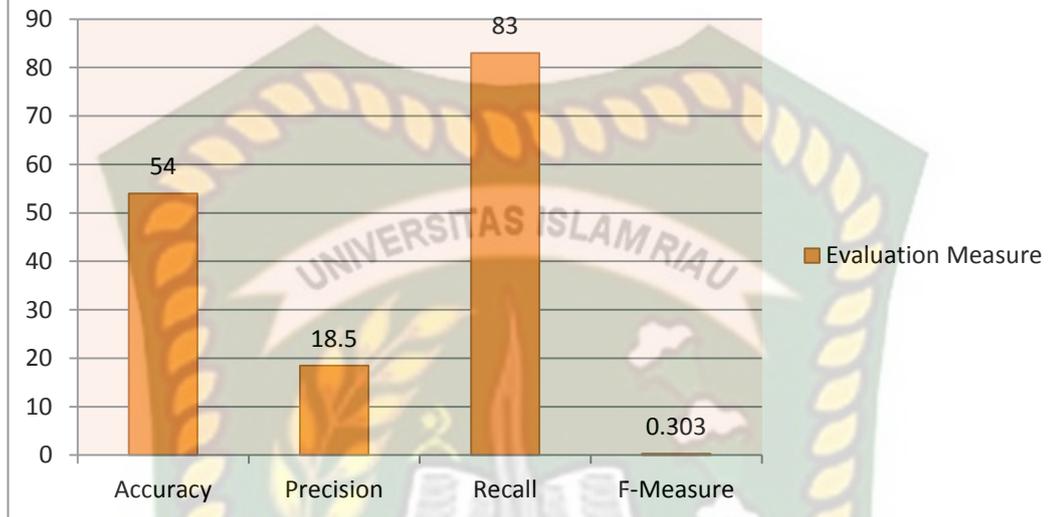


Diagram 4. 1 BarchartEvaluationMeasure Data A

Evaluation Measure

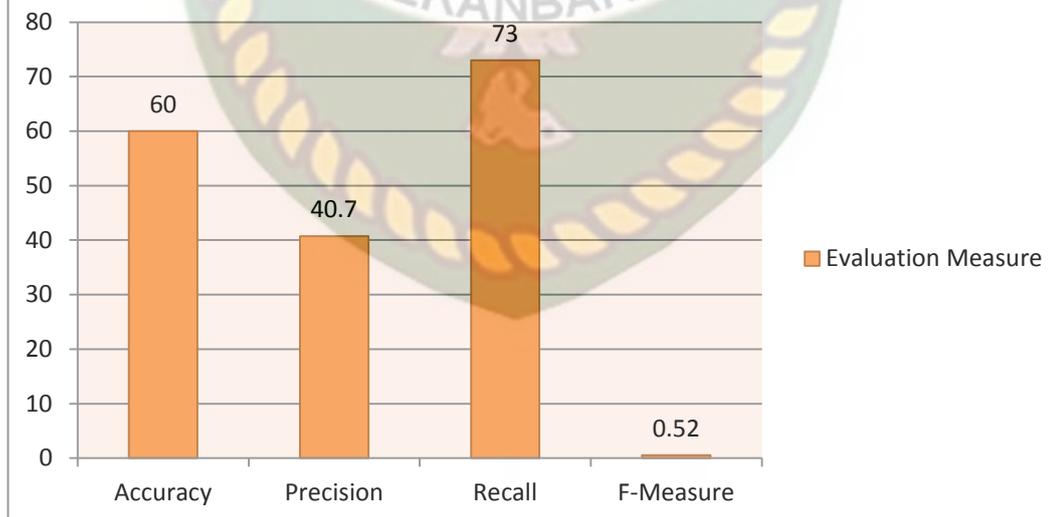


Diagram 4. 2 BarchartEvaluationMeasure Data B

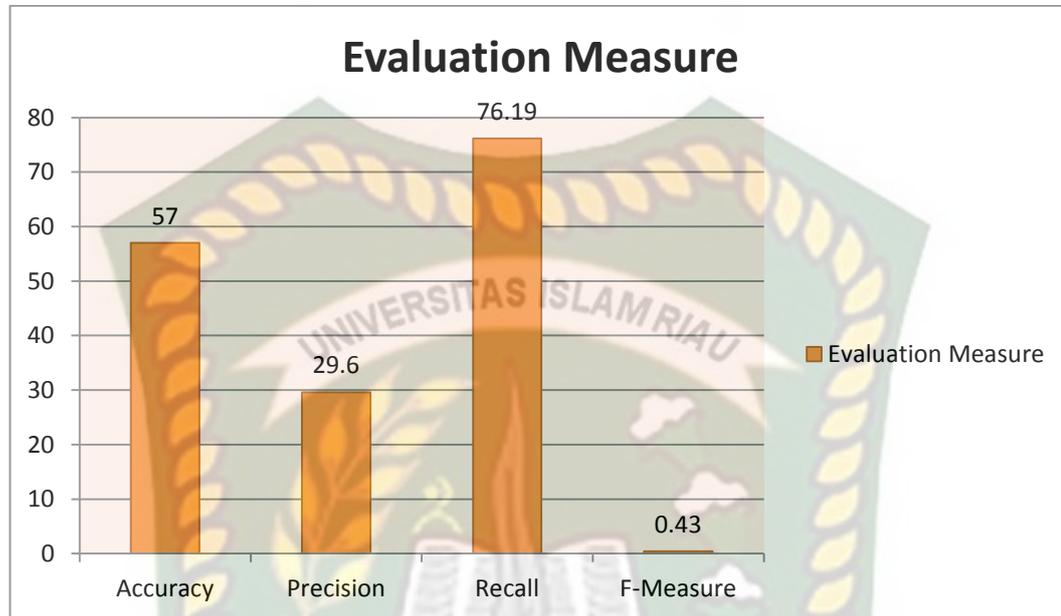


Diagram 4.3 Bar chart Evaluation Measure Data C

Evaluasi yang dilakukan pada data A,B dan C memiliki hasil yang berbeda-beda. Tetapi kesimpulan yang dapat ditarik adalah, perhitungan *recall* lebih menonjol dibandingkan dengan perhitungan lainnya. Ini berlaku pada data A, data B, dan data C.

4.3 Antarmuka Pada Analisis Sentimen

Pada penelitian ini, pengguna umum bisa menggunakan analisis sentimen untuk melihat apakah sebuah tweet memiliki kelas positif atau negatif. Antarmuka yang disediakan untuk pengguna umum adalah penginputan berupa text,yang nanti nya akan di proses dengan data training yang ada dan menghasilkan ouput berupa tabel yang berisi tweet yang diinputkan dan sentimen berupa positif atau negatif dari tweet tersebut. Berikut antarmuka untuk pengguna umum bisa dilihat

pada gambar 4.7 dan output dari antarmuka tersebut dapat dilihat pada gambar 4.8 di bawah ini.

Analisis Sentimen Pada Twitter
Dengan Tagar #KPUJanganCurang

Enter your text below

Biarin mereka nolak, yg penting #KPU adalah lembaga yang sah dlm penyelenggaraan pemilu

Result Clear

Tweet	Sentiment
Biarin mereka nolak, yg penting #KPU adalah lembaga yang sah dlm penyelenggaraan pemilu	Positif

Gambar 4. 16 Form *Input* User

Analisis Sentimen Pada Twitter
Dengan Tagar #KPUJanganCurang

Enter your text below

Result Clear

Tweet	Sentiment
Biarin mereka nolak, yg penting #KPU adalah lembaga yang sah dlm penyelenggaraan pemilu	Positif

Gambar 4. 17 Form *Output* User

4.4 Pengujian Kepada User

4.1.7 Implementasi User

Implementasi sistem yang dipakai adalah membuat kuisisioner dengan 5 pertanyaan dan 30 koresponden yang mana ditujukan kepada user yang ingin memakai sistem ini. kepada 30 koresponden diajukan pertanyaan yang terkait dengan kinerja atau *performance* dari sistem. Adapun kelima pertanyaan yang dimaksud adalah sebagai berikut :

1. Apakah informasi yang ditampilkan mudah dimengerti oleh user?
2. Apakah bahasa yang digunakan pada analisis sentimen mudah dimengerti?
3. Bagaimana pendapat anda tentang tampilan pada analisis sentimen ini?
4. Apakah analisis sentimen ini cukup mudah digunakan (dioperasikan)?
5. Menurut anda, apakah analisis sentimen ini sudah layak dipublikasikan?

Dari pertanyaan-pertanyaan diatas, maka hasil jawaban atau tanggapan dari koresponden terhadap kinerja atau *performance* dari sistem berdasarkan pertanyaan yang diajukan adalah sebagai berikut:

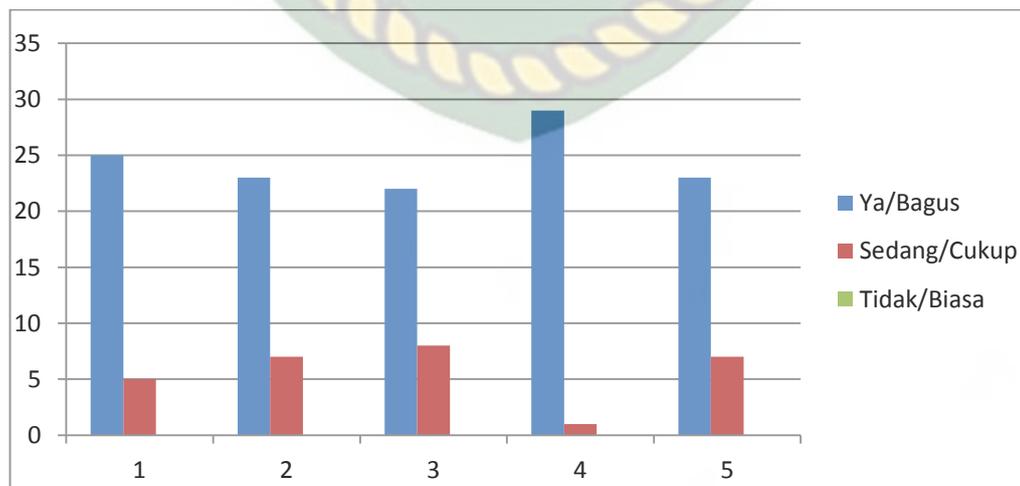


Diagram 4. 4 Grafik Hasil Kuisisioner

Keterangan:

1. Apakah informasi yang ditampilkan mudah dimengerti oleh user memiliki nilai YA : 25 koresponden, SEDANG : 5 koresponden, TIDAK : 0 koresponden.
2. Apakah bahasa yang digunakan pada analisis sentimen mudah dimengerti memiliki nilai YA : 23 koresponden, SEDANG : 7 koresponden, TIDAK : 0 koresponden.
3. Bagaimana pendapat anda mengenai tampilan pada analisis sentimen ini memiliki nilai BAGUS : 22 koresponden, CUKUP : 8 koresponden, BIASA : 0 koresponden.
4. Apakah analisis sentimen ini cukup mudah digunakan (dioperasikan) memiliki nilai YA : 29 koresponden, SEDANG : 1 koresponden, TIDAK : 0 koresponden.
5. Menurut anda, apakah analisis sentimen ini sudah layak dipublikasikan memiliki nilai YA : 23 koresponden, CUKUP : 7 koresponden, TIDAK : 0 koresponden.

4.1.8 Kesimpulan Implementasi Sistem

Berdasarkan hasil kuisioner tersebut maka dapat disimpulkan bahwa analisis sentimen pada tweet dengan tagar #KPUJanganCurang menggunakan metode naïve bayes ini memiliki persentase sebagai berikut:

Tabel 4. 9 Hasil Nilai Persentase Tiap Pertanyaan Kuisisioner

No	Pertanyaan	Jumlah Persentase Koresponden		
		Ya / Bagus	Sedang / Cukup	Tidak / Biasa
1	Apakah informasi yang ditampilkan mudah dimengerti oleh user?	83%	17%	0%
2	Apakah bahasa yang digunakan pada analisis sentimen mudah dimengerti?	77%	23%	0%
3	Bagaimana pendapat anda mengenai tampilan pada analisis sentimen ini?	73%	27%	0%
4	Apakah analisis sentimen ini cukup mudah digunakan (dioperasikan)	97%	3%	0%
5	Menurut anda, apakah analisis sentimen ini sudah layak dipublikasikan	77%	23%	0%

Dari hasil persentase tabel 4.2 nilai presentase tiap pertanyaan kuisisioner, analisis sentimen pada twitter dengan tagar #KPUJanganCurang memiliki *performance* baik dengan nilai persentase rata-rata sebesar 81.4%, sehingga aplikasi ini dapat diimplementasikan ke publik.

BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan hasil analisa yang dilakukan pada kasus di atas, maka dapat ditarik beberapa kesimpulan, yaitu :

1. Jumlah dan keanekaragaman data berpengaruh terhadap akurasi pada sistem dan perhitungan *evaluation measure*.
2. Saat dilakukan perhitungan *evaluation measure* pada data A,B dan C, hasil dari data B lebih tinggi dibandingkan dengan data A dan C untuk perhitungan accuracy, precision, dan F-1 sedangkan data A adalah nilai terendah. Untuk perhitungan recall, data A memiliki nilai lebih tinggi.
3. Metode naïve bayes cocok digunakan pada analisis sentimen ini karena menghasilkan akurasi yang tinggi dengan akurasi pada sistem 89% pada 300 data.

5.2 Saran

Untuk penelitian berikutnya, dapat menggunakan metode terbaru yang bisa digunakan pada analisis sentimen dan menambahkan data yang lebih banyak untuk akurasi yang lebih tinggi.

DAFTAR PUSTAKA

- Bahary, A. F., Sibaroni, Y., & Mubarok, M. S. 2019. *Sentiment Analysis Of Student Responses Related To Information System Service Using Multinomial Naïve Bayes*. Telkom University.
- Deviyanto, A., & Wahyudi, M. D. R. 2018. *Penerapan Analisis Sentimen Pada Pengguna Twitter Menggunakan Metode K-Nearest Neighbor*. JISKA (Jurnal Informatika Sunan Kalijaga).
- Falahah, & Dwiki Adriadi Nur, D. 2015. *Pengembangan Aplikasi Sentimen Analisis Menggunakan Metode Naïve Bayes (Studi Kasus Sentiment Analysis dari media Twitter)*. Seminar Nasional Sistem Informasi Indonesia.
- Farisi, A. A., Sibaroni, Y., & Faraby, S. A. 2019. *Analisis Sentimen Pada Ulasan Hotel Menggunakan Multinomial Naïve Bayes Classifier*. School of Computing, Telkom Universitas.
- Lestari, Agnes R.T., Perdana, Rizal S., & Fauzi, M.Ali. 2017. *Analisis Sentimen Tentang Opini Pilkada Dki 2017 Pada Dokumen Twitter Berbahasa Indonesia Menggunakan Naïve Bayes dan Pembobotan Emoji*. Fakultas Ilmu Komputer, Universitas Brawijaya. Vol.1. No.12.

- Mardalius. 2018. *Pemanfaatan Rapid Miner Studio 8.2 Untuk Pengelompokan Data Penjualan Aksesoris Menggunakan Algoritma K-Means*. Program Studi Sistem Informasi, STMIK Royal Kisaran. Vol.4. No. 2.
- Nugroho, D. G., Chrisnanto, Y. H., & Wahana, A. 2016. *Analisis Sentimen Pada Jasa Ojek Online*. Prosiding SNST Fakultas Teknik.
- Nurhuda, F., Widya Sihwi, S., & Doewes, A. 2016. *Analisis Sentimen Masyarakat terhadap Calon Presiden Indonesia 2014 berdasarkan Opini dari Twitter Menggunakan Metode Naive Bayes Classifier*. Jurnal Teknologi & Informasi ITSmart.
- Putranti, Noviah Dwi & Winarko, Edi. 2014. *Analisis Sentimen Twitter untuk Teks Berbahasa Indonesia dengan Maximum Entropy dan Support Vector Machine*. FMIPA UGM, Yogyakarta. Vol.8. No.1.
- Rahman, Amelia., Wiranto., & Doewes, Afrizal. 2017. *Online News Classification Using Multinomial Naïve Bayes*. ITSMART:Jurnal Ilmiah Teknologi dan Informasi. Vol.6. No.1
- Russell, Stuart & Norvig, Peter. 2010. *Artificial Intelligence: A modern Approach, 3rd ed.* America: Pearson Education.
- Santoso, Budi. 2007. *DATA MINING: Teknik Pemanfaatan Data untuk Keperluan Bisnis*. Yogyakarta: Graha Ilmu.
- Sunni, Ismail & Widyantoro, Dwi.H. 2012. *Analisis Sentimen dan Ekstraksi Topik Penentu Sentimen pada Opini Terhadap Tokoh Publik*. Fakultas Teknik Informatika, Institut Teknologi Bandung. Vol.1. No.2.
- Tama, V., Sibaroni, Y., & Adiwijaya. 2019. *Analisis Pelabelan Dalam Klasifikasi Sentimen*

Review Produk Dengan Menggunakan Multinomial Algoritma Naïve Bayes. Telkom University.

Taylor, John Shawe & Cristianini, Nello. 2004. *Kernel Methods for Pattern Analysis.* Cambridge University Press.

Y, Lu & C, Rasmussen. 2012. *Simplified markov random fields for efficient semantic labeling of 3D point clouds.* IROS.

Zulfa, Ira & Winarko, Edi. *Sentimen Analisis Tweet Berbahasa Indonesia dengan Deep Belief Network.* FMIPA UGM, Yogyakarta. Vol.11. No.2.

Tama, V., Sibaroni, Y., & Adiwijaya. 2019. *Analisis Pelabelan Dalam Klasifikasi Sentimen Review Produk Dengan Menggunakan Multinomial Algoritma Naïve Bayes.* Telkom University.

Habibi, Robet., Budiyanto Setyohadi, Djoko., & Ernawati. 2016. *Analisis Sentimen Pada Twitter Mahasiswa Menggunakan Metode Backpropagation.* Universitas Atma Jaya, Yogyakarta. Vol.12. No.1.

Fauzi,M.Ali. 2018. *Pendekatan Metode Random Forest Pada Analisis Sentimen Dalam Bahasa Indonesia.* Universitas Brawijaya, Malang. Vol.12. No.1.