

Identifikasi *Tweet Hoax* yang Berhubungan dengan Pemilihan Presiden 2019 Menggunakan *Naïve Bayes Classifier*

Tugas Akhir
diajukan untuk memenuhi salah satu syarat
memperoleh gelar sarjana
dari Program Studi S1 Informatika
Fakultas Informatika
Universitas Telkom

1301178533

Azizah Zain



Program Studi Sarjana Informatika
Fakultas Informatika
Universitas Telkom
Bandung
2019

LEMBAR PENGESAHAN

**Identifikasi *Tweet Hoax* yang Berhubungan dengan Pemilihan Presiden 2019
Menggunakan *Naïve Bayes Classifier***

**Identification of *Tweet Hoax* Associated with Presidential Election 2019 Using *Naïve*
*Bayes Classifier***

NIM: 1301178533

Azizah Zain

Tugas akhir ini telah diterima dan disahkan untuk memenuhi sebagian syarat
memperoleh gelar pada Program Studi Sarjana Informatika
Fakultas Informatika Universitas Telkom

Bandung, 14 Januari 2020 Menyetujui

Pembimbing I,



Anisa Herdiani, S.T., M.T.
15850002-1

Pembimbing II,



Indra Lukmana Sardi, S.T.,
M.T.
18890120

Ketua Program Studi
Sarjana Informatika



Niken Dwi Wahyu Cahyani,
S.T., M.Kom., Ph.D.

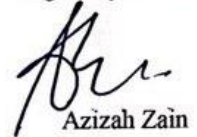
00750052

LEMBAR PERNYATAAN

Dengan ini saya, Azizah Zain, menyatakan sesungguhnya bahwa Tugas Akhir saya dengan judul Identifikasi *Tweet Hoax* yang Berhubungan dengan Pemilihan Presiden 2019 Menggunakan *Naïve Bayes Classifier* beserta dengan seluruh isinya adalah merupakan hasil karya sendiri, dan saya tidak melakukan penjiplakan yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan. Saya siap menanggung resiko/sanksi yang diberikan jika di kemudian hari ditemukan pelanggaran terhadap etika keilmuan dalam Laporan TA atau jika ada klaim dari pihak lain terhadap keaslian karya,

Bandung, 14 Januari 2020

Yang Menyatakan



Azizah Zain

Identifikasi *Tweet Hoax* yang Berhubungan dengan Pemilihan Presiden 2019

Menggunakan *Naïve Bayes Classifier*

Azizah Zain¹, Anisa Herdiani², Indra Lukmana Sardi³

^{1,2,3}Fakultas Informatika, Universitas Telkom, Bandung

¹azizaahzain@students.telkomuniversity.ac.id, ²anisaherdiani@telkomuniversity.ac.id,

³indraluk@telkomuniversity.ac.id

Abstrak

Hoaks merupakan berita yang tersebar tanpa tahu kebenarannya atau faktanya. Hoaks biasanya tersebar melalui media yang mudah diakses seperti media sosial facebook, twitter dan lain-lain. Terutama pada masa pemilihan presiden (pilpres), para pasangan calon presiden dan wakil presiden tersebut memanfaatkan media sosial untuk melakukan kampanye. Hal ini tentu dimanfaatkan oleh orang-orang yang tidak bertanggung jawab untuk menyebarkan berita hoaks. Oleh karena itu perlu adanya pengidentifikasian *tweet-tweet* hoaks yang berhubungan dengan pemilihan presiden 2019 agar masyarakat dapat terhindar dari berita hoaks. Metode yang digunakan untuk identifikasi hoaks adalah *Naïve Bayes Classifier* (NBC). Metode NBC dipilih karena telah terbukti efektif untuk kategorisasi teks, prosesnya sederhana, cepat dan akurasi klasifikasi yang tinggi. Data tweet yang digunakan sebanyak 500 tweet dengan 143 tweet hoaks dan 357 tweet bukan hoaks. Hasil menunjukkan bahwa NBC dengan tambahan seleksi fitur *Mutual Information* (MI) dan *Information Gain* (IG) dapat menghasilkan nilai *precision* sebesar 0.9146, *recall* sebesar 0.9146, dan *F1 score* sebesar 0.9433 sedangkan NBC tanpa seleksi fitur menghasilkan nilai *precision* sebesar 0.6932, *recall* sebesar 0.9531, dan *F1 score* sebesar 0.8026, dengan kata lain seleksi fitur dapat menambah nilai *F1-score* sebesar 0.1408.

Kata kunci: hoaks, pemilihan presiden, *tweet*, *Naïve Bayes Classifier*.

Abstract

Hoax are news that is spread without knowing the truth or the facts. Hoax are usually spread through easily accessible media such as Facebook, Twitter and others. Especially during the presidential election, the pair of candidates for president and vice president used social media to carry out campaigns. This is certainly used by people who are not responsible for spreading hoax news. Therefore, it is necessary to identify hoax tweets related to the 2019 presidential election so that people can avoid hoax. The method used for hoax identification is *Naïve Bayes Classifier* (NBC). The NBC method was chosen because it has been proven effective for text categorization, the process is simple, fast and has high classification accuracy. The tweet data used was 500 tweets with 143 hoax tweets and 357 non-hoax tweets. The results show that NBC with additional *Mutual Information* (MI) and *Information Gain* (IG) feature selection can produce a precision value of 0.974, a recall of 0.9146, and an *F1 score* of 0.9433 while an NBC without a feature selection produces a precision value of 0.6932, a recall of 0.9531, and *F1 score* of 0.8026, in other words feature selection can increase the value of *F1-score* of 0.1408.

Keywords: hoaks, presidential elections, tweets, *Naïve Bayes Classifier*.

1. Pendahuluan

1.1. Latar Belakang

Hoaks adalah berita yang tersebar tanpa tahu kebenaran atau faktanya [1]. Tak sedikit berita-berita bohong (hoaks) digunakan untuk membentuk opini publik yang mengarah pada terjadinya kehebohan, ketidakpastian informasi, dan ketakutan [2]. Bahkan, beberapa penyedia berita hoaks berusaha untuk mendukung ideologi yang diusungnya dengan menyerang kelompok oposisi yang menjadi rivalnya [3]. Menurut Hunt Allcott dan Matthew Gentzkow, 62% orang dewasa di Amerika Serikat mendapatkan berita dari media sosial [4] dimana terdapat banyak berita bohong yang tersebar. Terutama pada masa pemilihan presiden (pilpres), para pasangan calon presiden dan wakil presiden tersebut memanfaatkan media sosial untuk melakukan kampanye. Namun kampanye pilpres yang disampaikan melalui media sosial tidak selalu benar. Khususnya melalui Twitter, media sosial yang berisi banyak pendapat dari penggunanya, sehingga tidak dapat dipastikan kebenarannya [5]. Twitter juga memiliki struktur komunitas yang unik [6] cara desas-desus atau gosip menyebar melalui Twitter juga memiliki beberapa perbedaan dengan berita yang tersebar melalui mulut ke mulut dengan jejaring sosial konvensional [7]. Hal ini tentu dimanfaatkan oleh orang-orang yang tidak bertanggung jawab untuk menyebarkan berita hoaks. Oleh karena itu perlu adanya pengidentifikasian *tweet-tweet* hoaks yang berhubungan dengan pemilihan presiden 2019 agar masyarakat dapat terhindar dari berita hoaks.

Metode yang digunakan untuk mengidentifikasi *tweet* hoaks yang berhubungan dengan pemilihan presiden 2019 adalah *Naïve Bayes Classifier* (NBC). Metode NBC dipilih karena dapat mengidentifikasi berita hoaks dengan akurasi sebesar 91.36% dengan tambahan metode *mutual information* dan *information gain* sebagai seleksi fitur dan unigram sebagai ekstraksi fiturnya [8]. Selain itu, NBC telah terbukti efektif untuk kategorisasi teks,

prosesnya sederhana, cepat dan akurasi klasifikasi yang tinggi. Data yang digunakan dalam penelitian ini akan diambil dari Twitter dengan *crawling* data menggunakan Bahasa Pemrograman Python.

1.2. Topik dan Batasan Masalah

Berdasarkan latar belakang yang sudah dijabarkan sebelumnya, maka didapatkan rumusan masalah sebagai berikut:

1. Bagaimana mengidentifikasi *tweet* yang berkaitan dengan pemilihan presiden 2019 termasuk hoaks atau berita bohong menggunakan metode *Naïve Bayes Classifier*?
2. Bagaimana mengukur tingkat performansi sistem dalam mengidentifikasi *tweet* hoaks dengan menggunakan *Naïve Bayes Classifier*?

Dalam penelitian dibutuhkan batasan masalah agar penelitian yang dilakukan dapat berjalan sesuai dengan hasil yang diharapkan. Adapun batasan masalah tersebut yaitu:

1. Data diambil dari Twitter selama masa kampanye (23 September 2018 sampai 13 April 2019)
2. *Tweet* yang diidentifikasi adalah *tweet* berbahasa Indonesia.
3. Hoaks yang teridentifikasi adalah hoaks yang datanya terdapat di dalam data latih.

1.3. Tujuan

Pada penelitian ini diimplementasikan pengidentifikasian *tweet* hoaks yang berhubungan dengan pemilihan presiden 2019 menggunakan *Naïve Bayes Classifier* dan pengukuran performansi sistem menggunakan *FI-score*. Data diambil dari *tweet* yang berhubungan dengan pemilihan presiden 2019. Dari data *tweet* tersebut diklasifikasikan ke dalam 2 kelas yaitu hoaks dan bukan hoaks.

1.4. Organisasi Tulisan

Dalam melakukan penelitian ini, langkah pertama yang dilakukan adalah melakukan studi literatur mengenai teori-teori yang akan digunakan. Kemudian, dilakukan *crawling* data *tweet* yang berkaitan dengan pemilihan presiden 2019. Langkah selanjutnya adalah melakukan pengklasifikasian data hoaks dan bukan hoaks dengan metode *Naïve Bayes Classifier*. Setelah didapatkan hasil prediksi klasifikasi dengan menggunakan *Naïve Bayes Classifier*, maka bisa dilakukan analisis perhitungan evaluasi performansi sistem dengan menggunakan *FI Score*.

2. Studi Terkait

2.1 Hoaks

Hoaks adalah berita yang beredar tanpa tahu kebenaran atau faktanya [1]. Hoaks seringkali muncul dalam pemberitaan dalam suatu media, media cetak maupun media sosial. Hoaks sendiri berasal dari bahasa Inggris yaitu *hoax* yang artinya tipuan, menipu, berita bohong, berita palsu dan kabar burung. Menurut Robert Nares, kata *hoax* muncul sejak abad 18 yang berarti "permainan sulap" [9]. Pada umumnya hoaks tersebar akibat ketidaklengkapan seseorang dalam membaca atau menangkap sebuah berita yang ada [10]. Atau bisa jadi hoaks sengaja disebarkan untuk menguntungkan beberapa pihak.

2.2 Mutual Information

Mutual information (MI) menunjukkan seberapa banyak informasi mengenai ada atau tidaknya sebuah kata memberikan kontribusi dalam membuat keputusan klasifikasi secara benar atau salah.

$$I(U; C) = \sum_{ec \in \{1,0\}} P(U = et, C = ec) \log_2 \frac{P(U = et, C = ec)}{P(U = et) P(C = ec)} \quad (1)$$

Dengan keterangan, U merupakan sebuah *term* atau fitur dan C adalah kelas yang ada. Nilai et atau ec bernilai 0 atau 1 dimana, jika nilai et = 1, maka dokumen berisi *term* atau fitur t dan jika nilai et = 0, maka dokumen yang tidak mengandung *term* atau fitur t. Kemudian untuk nilai ec = 1, maka dokumen ada dikelas c dan jika nilai ec = 0, maka dokumen tidak ada dikelas c.

2.3 Information Gain

Information gain (IG) digunakan untuk menghitung pengaruh suatu fitur terhadap keseragaman kelas pada data yang dipecah menjadi sub data dengan nilai fitur tertentu. Untuk mendapatkan nilai IG dibutuhkan perhitungan entropy sebelum data dipisah dan sesudah data dipisah.

$$Entropy(S) = \sum_{i=1}^k (P_i) \log_2 (P_i) \quad (2)$$

Dengan nilai Pi adalah probabilitas data S didalam kelas i. K adalah jumlah kelas untuk Variable S.

$$Entropy(S, A) = \sum_{i=1}^v \left(\frac{Sv}{S} * Entropy(Sv) \right) \quad (3)$$

Dengan nilai v adalah semua nilai yang mungkin dari atribut A , S_v adalah subset dari S dimana atribut A bernilai v . Nilai information gain dihitung dari persamaan berikut:

$$Gain(S, A) = Entropy(S) - Entropy(S, A) \quad (4)$$

Dengan nilai Gain (S, A) adalah nilai information gain. Entropy (S) adalah nilai entropy sebelum pemisahan. Entropy (S, A) adalah nilai entropy setelah pemisahan. Besarnya nilai information gain menunjukkan seberapa besar pengaruh suatu atribut terhadap pengklasifikasian data.

2.4 Naïve Bayes Classifier

Naïve Bayes Classifier (NBC) merupakan metode *supervised learning* yang digunakan untuk pengklasifikasian teks berdasarkan teorema Bayes dengan asumsi “naif” tentang independensi bersyarat antara setiap pasangan fitur yang diberi nilai variabel kelas [11]. Pada langkah pertama dalam penggunaan metode ini, akan digunakan metode TF-IDF untuk merepresentasikan kata-kata tersebut ke bentuk angka. Guna dilakukannya perhitungan ini adalah untuk mengidentifikasi seberapa penting sebuah kata dalam teks berdasarkan frekuensinya. Selanjutnya akan dilakukan tahap pengklasifikasian data menggunakan *Naïve Bayes Classifier* dengan rumus sebagai berikut [11]:

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y) \quad (5)$$

Dimana $P(x_i|y)$ adalah probabilitas kondisional dari t_k yang berada pada dokumen yang dimiliki kelas y . Dari persamaan diatas dapat diketahui bahwasannya $P(x_i|y)$ merupakan *likelihood probability* dari x_i yang terdapat pada kelas y , di sisi lain $P(y)$ merupakan *prior probability* dokumen yang berada pada kelas y . Hasil dari *posterior probability* akan dibandingkan untuk menentukan kelas, kelas yang memiliki nilai *posterior probability* terbesar merupakan kelas yang dipilih sebagai hasil prediksi [12].

Rumus *Prior probability* :

$$P(y) = \frac{N_y}{N} \quad (6)$$

N_y merupakan jumlah dari kategori dari y . N jumlah kategori keseluruhan. Untuk *likelihood probability* akan dihitung kata atau fitur x_i pada seluruh dokumen latih pada y , dengan menggunakan *LaPlace Smoothing* memiliki Rumus [13]:

$$P(x_i|y) = \frac{N_y + 1}{|V| + N'} \quad (7)$$

2.5 Twitter

Twitter merupakan merupakan situs microblogging yang dibatasi 140 karakter yang dapat menyertakan link ke website atau media (foto, video, suara) yang di broadcast kepada pengikutnya (followers) [14]. *Tweet* adalah penulisan teks yang yang berada di halaman publik twitter, pengguna bisa membatasi pengiriman tweet atau pesan yang hanya dapat dilihat oleh teman-teman atau pengikut (Followers). Pengguna Twitter dapat mengelompokkan *tweet* yang dibuat menjadi satu tema yang disebut dengan *thread*. Jadi jika pengguna ingin membuat *tweet* yang panjang, Twitter menyediakan fasilitas *thread* agar setiap *tweet* yang berkaitan langsung terhubung. Atau bisa juga pengguna menggunakan tagar “#” (*hashtag*) untuk mempermudah pengguna dalam pencarian. Selain itu pengguna juga bisa menyebut pengguna Twitter lain dalam *tweet* yang akan dipublikasi menggunakan tanda “@” yang dilanjutkan dengan *username* pengguna yang ingin ditambahkan. Setelah menambahkan pengguna lain dalam *tweet* yang dipublikasikan, pengguna Twitter yang lain pun dapat membalas *tweet* yang sudah dipublikasikan yang disebut dengan *reply*. Twitter juga mempunyai mekanisme yang mendukung untuk menyebarkan informasi yang didapat dari hasil tulisan orang lain yaitu *retweet*.

2.6 Evaluasi

Tahap evaluasi bertujuan untuk mengukur performansi dari sistem yang telah dibangun. Evaluasi akan dilakukan menggunakan tabel yang berisi prediksi yang disebut *confusion matrix*. Berikut adalah tabel dari *confusion matrix* [15]:

Tabel 1 *Confusion Matrix* [15]

	Actual Positive	Actual Negative
Predicted Positive	True Positive (TP)	False Positive (FP)
Predicted Negative	False Negative (FN)	True Negative (TN)

Dari tabel tersebut terdapat 4 kategori, yaitu [15]:

- *True Positive* (TP) yaitu hasil yang positif dan dilabeli positif.
- *False Negative* (FN) yaitu hasil yang negatif dan dilabeli negatif.
- *False Positive* (FP) yaitu hasil yang negatif tetapi dilabeli positif.
- *True Negative* (TN) yaitu hasil yang positif tetapi dilabeli negatif.

Dalam *Precision and Recall* terdapat 2 definisi *confusion matrix* yaitu *Precision* pada sumbu y dan *Recall* pada sumbu x, berikut adalah rinciannya [15]:

- *Precision* digunakan untuk mengukur bagian-bagian kecil dari contoh-contoh yang diklasifikasikan sebagai positif dan mendapatkan label yang positif. *Precision* didapatkan dengan rumus sebagai berikut [15]:

$$Precision = \frac{TP}{TP + FP} \quad (8)$$

- *Recall* digunakan untuk mengukur bagian-bagian kecil dari contoh-contoh yang diklasifikasikan sebagai positif dan mendapatkan label yang benar. *Recall* didapatkan dengan rumus sebagai berikut [15]:

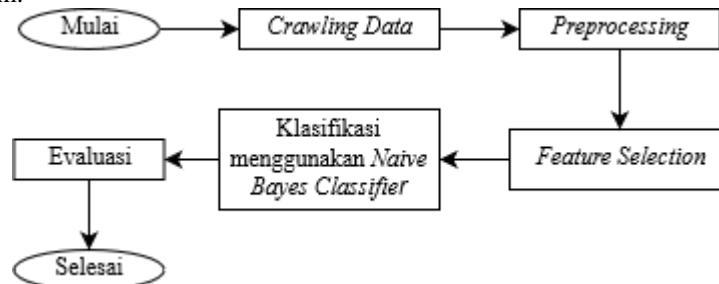
$$Recall = \frac{TP}{TP + FN} \quad (9)$$

Setelah didapatkan nilai *precision* dan *recall*, maka dapat dihitung nilai kombinasi dari *precision* dan *recall* menggunakan rumus *F1 Score* sebagai berikut [15]:

$$F1\ Score = \frac{2(Precision \times Recall)}{Precision + Recall} \quad (10)$$

3. Metodologi

Sistem yang akan dibangun adalah sebuah sistem yang mampu mengidentifikasi *tweet-tweet* hoaks yang berhubungan dengan Pemilihan Presiden 2019. Dalam sistem ini terdapat 5 bagian besar yaitu pengumpulan data, *preprocessing*, perhitungan MI dan IG, klasifikasi menggunakan NBC dan evaluasi. Berikut adalah gambaran umum alur kerja sistem:



Gambar 1 Gambaran Umum Alur Kerja Sistem

Berikut adalah penjelasan dari Gambar 1:

3.1. *Crawling Data*

Pada tahap pengumpulan pada Twitter atau *crawling* data ini memiliki beberapa tahap, yaitu pembuatan akun Dev Twitter dan Apps Twitter untuk memperoleh akses token dan *consumer key*, serta pembuatan kode program untuk *crawling* data menggunakan Bahasa Pemrograman Python. Data yang diambil hanya yang menggunakan Bahasa Indonesia dengan kata kunci yang berasal dari akun Twitter dari masing-masing pasangan calon presiden dan wakil presiden yaitu, @jokowi, @prabowo, @sandiuno, dan @Kiyai_MarufAmin.

Total data yang dikumpulkan adalah 500 *tweets* yang divalidasi oleh jurnalis. Data dengan label tidak benar ada 143 *tweets* dan data dengan label benar ada 357 *tweets*.

Tabel 2 Contoh Data Set

<i>Tweet</i>
@jokowi Dukungan dari kepala daerah di Sumatera Barat terhadap capres Joko Widodo di Pilpres 2019 bertambah. Dukungan terbaru disampaikan Wali Kota Padang Panjang Fadhly Amran. #2019JokowiKyaiMaruf #2019TetapJokowi #BerSATUikutKyai #KyaiPemimpin #KyaiTeladan

3.2. Preprocessing

Preprocessing merupakan tahapan pengolahan data untuk menghilangkan *data noise* sehingga dapat menghasilkan *clean word*. Ada beberapa tahapan yang dilakukan pada *preprocessing* yaitu *data cleaning*, *case folding*, *tokenization*, dan *stopword removal*.

a. Data cleaning

Proses untuk penghapusan *website link*, simbol, *username*, dan angka. Hal ini bertujuan untuk mengurangi data kotor yang akan masuk ke dalam perhitungan klasifikasi.

Tabel 3 Contoh proses *data cleaning*

Output
Dukungan dari kepala daerah di Sumatera Barat terhadap capres Joko Widodo di Pilpres 2019 bertambah Dukungan terbaru disampaikan Wali Kota Padang Panjang Fadhly Amran

b. Case folding

Tahap konversi teks menjadi suatu bentuk standar, biasanya huruf kecil atau disebut juga *lowercase*. *Case folding* dilakukan untuk meminimalisir perbedaan nilai pada pembobotan kata. Karena jika suatu data atau fitur memiliki huruf kapital, meskipun tulisannya sama dengan data atau fitur yang sudah menjadi *lowercase* akan memiliki bobot nilai yang berbeda.

Tabel 4 Contoh proses *case folding*

Output
dukungan dari kepala daerah di sumatera barat terhadap capres joko widodo di pilpres 2019 bertambah dukungan terbaru disampaikan wali kota padang panjang fadhly amran

c. Penanganan singkatan dan *Tokenization*

Kata-kata yang disingkat mempengaruhi tahap pengklasifikasian, sehingga harus dilakukan konversi ke bentuk aslinya. Setelah itu dilakukan proses mengubah suatu kalimat, dalam hal ini *tweet* ke dalam bentuk *token* per kata yaitu *tokenize*. Hasil dari proses ini yaitu berupa *unigram*.

Tabel 5 Contoh proses penanganan singkatan dan *tokenization*

Output
['dukungan'], ['dari'], ['kepala'], ['daerah'], ['di'], ['sumatera'], ['barat'], ['terhadap'], ['capres'], ['joko'], ['widodo'], ['di'], ['pilpres'], ['2019'], ['bertambah'], ['dukungan'], ['terbaru'], ['disampaikan'], ['wali'], ['kota'], ['padang'], ['panjang'], ['fadhly']

d. Stopword removal

Penghapusan kata-kata yang bersifat umum namun tidak memiliki makna atau informasi yang dibutuhkan. Hal ini dilakukan agar pada saat identifikasi *tweet*, dapat mengurangi fitur yang kurang relevan.

Tabel 6 Contoh proses *stopword removal*

Output
['dukungan'], [di], ['kepala'], ['daerah'], [di], ['sumatera'], ['barat'], ['terhadap'], ['capres'], ['joko'], ['widodo'], [di], ['pilpres'], ['2019'], ['bertambah'], ['dukungan'], ['terbaru'], ['disampaikan'], ['wali'], ['kota'], ['padang'], ['panjang'], ['fadhly']

3.3. Perhitungan *Feature Selection*

Seleksi fitur atau *feature selection* digunakan untuk memilih fitur yang spesifik terkait dengan pengidentifikasian *tweet* hoaks.

3.3.1. Perhitungan *Mutual Information*

Hasil *preprocessing* yang berupa *unigram*, akan dilakukan perhitungan *term* berdasarkan kemunculan di setiap dokumen. Berikut adalah contoh data yang akan digunakan dalam perhitungan MI:

Tabel 7 Contoh data setelah *preprocessing*

Dokumen	<i>Tweet</i>	Kelas	<i>Unigram</i>
1	prabowo pecat dari tni	Tidak Benar	['prabowo'], ['pecat'], ['dari'], ['tni']
2	padahal prabowo pindah jabatan tni	Benar	['padahal'], ['prabowo'], ['pindah'], ['jabatan'], ['tni']
3	seluruh alumni guru tinggi dukung jokowi	Tidak Benar	['seluruh'], ['alumni'], ['guru'], ['tinggi'], ['dukung'], ['jokowi']
4	tni pecat prabowo	Tidak Benar	['tni'], ['pecat'], ['prabowo']
5	jabatan jokowi adalah presiden	Benar	['jabatan'], ['jokowi'], ['adalah'], ['presiden']

Kemudian setelah data melalui proses tokenisasi, data yang tadinya berupa kalimat, sudah dipecah menjadi kata per kata atau fitur yang kemudian akan dihitung frekuensi kemunculannya agar dapat dihitung bobot MI nya.

Tabel 8 Frekuensi kemunculan kata pada kelas “Tidak Benar”

No	<i>Unigram</i>	Frekuensi kemunculan
1	prabowo	2
2	pecat	2
3	dari	1
4	tni	2
5	seluruh	1
6	alumni	1
7	guru	1
8	tinggi	1
9	dukung	1
10	jokowi	1

Tabel 9 Frekuensi kemunculan kata pada kelas “Benar”

No	<i>Unigram</i>	Frekuensi kemunculan
1	padahal	1
2	prabowo	1
3	pindah	1
4	jabatan	2
5	tni	1
6	jokowi	1
7	adalah	1
8	presiden	1

Setelah frekuensi kemunculan setiap *term* diketahui, maka akan dilakukan penyeleksian *term*. Caranya dengan menghitung berapa banyak jumlah *term* x di kelas y dan kelas non y. Lalu dimasukkan ke dalam rumus *mutual information*. Contoh untuk “pecat” jumlah *term* “pecat” pada kelas Tidak Benar ada 2, sedangkan pada kelas Benar ada 0. Sehingga nilai N11 dan N10 pada rumus MI secara berurutan adalah 2 dan 0. Untuk mencari nilai N01 dengan cara menghitung jumlah *term* selain *term* “pecat” pada kelas Tidak Benar, sedangkan N00 menghitung jumlah *term* selain *term* “pecat” pada kelas Benar. Nilai N01 dan N00 secara berurutan yaitu 43 dan 30. Lalu setiap variabel tersebut akan dimasukkan ke dalam rumus MI. Tujuan dari MI yaitu memilih *term* yang mempunyai peranan penting

dalam setiap kelas.

3.3.2. Perhitungan *Information Gain*

Untuk melakukan perhitungan IG dilakukan perhitungan entropy sebelum data dipisah dan sesudah data dipisah. Berikut adalah contoh perhitungannya:

Tabel 10 Contoh data *tweet*

Dokumen	<i>Tweet</i>	Kelas
1	prabowo pecat dari tni	Tidak Benar
2	padahal prabowo pindah jabatan tni	Benar
3	seluruh alumni guru tinggi dukung Jokowi	Tidak Benar
4	tni pecat prabowo	Tidak Benar
5	jabatan jokowi adalah presiden	Benar

Langkah-langkah untuk melakukan perhitungan IG:

- Melakukan perhitungan entropy sebelum data dipisah. Misal perhitungan entropy fitur Prabowo.

$$Entropy (prabowo) = \frac{2}{16} \log_2 \frac{2}{16} + \frac{1}{16} \log_2 \frac{1}{16} = 0.625$$

- Melakukan perhitungan setelah data dipisah, misal ambil *threshold* = 0.4, jadi terbagi dua yaitu ada 2 dan 3 data.

$$Entropy (prabowo) = \frac{1}{7} \log_2 \frac{1}{7} + \frac{1}{7} \log_2 \frac{1}{7} = 0.8$$

$$Entropy (prabowo) = \frac{1}{12} \log_2 \frac{1}{12} + \frac{0}{12} \log_2 \frac{0}{12} = 0.298$$

- Kemudian hitung entropy setelah dipisahny data untuk melihat bobot informasi data.

$$Entropy split = 0.2 * 0.8 + 0.3 * 0.298 = 0.249$$

- Selanjutnya dilakukan perhitungan untuk mendapatkan nilai IG.

$$Gain = 0.625 - 0.249 = 0.376$$

Semakin besar nilai dari IG, maka semakin penting informasi yang dimiliki oleh fitur tersebut. IG juga dipengaruhi oleh *threshold* yang digunakan.

3.4. Klasifikasi

Pada tahap klasifikasi, akan dilakukan klasifikasi tweet menggunakan *Naïve Bayes Classifier* (NBC). Tahapan awal dari klasifikasi ini adalah pembangunan model yang dilakukan oleh NBC menggunakan data latih yang ada untuk kemudian akan digunakan mengklasifikasi data pada data tes. Pembangunan model akan diawali dengan melakukan perhitungan *prior probability*. Berikut contoh data yang akan digunakan.

Tabel 11 Contoh data latih dan data uji

	Dokumen	<i>Tweet</i>	Kelas
Data Latih	1	prabowo adalah presiden indonesia	Tidak Benar
	2	jokowi jadi presiden	Benar
	3	prabowo angkat jokowi jadi menteri	Tidak Benar
Data Uji	4	jokowi calon presiden	?

Sebelum melakukan perhitungan NBC, data latih dan data uji terlebih dahulu dipilih fiturnya berdasarkan hasil seleksi fitur MI dan IG.

Tabel 12 Proses perhitungan kemunculan fitur pada masing-masing data

Dokumen	Prabowo	Pecat	Tni	Jabatan	Jokowi	Presiden
1	1	0	0	0	0	1
2	0	0	0	0	1	1
3	1	0	0	0	1	0
4	0	0	0	0	1	1

Setelah dipilih fitur yang ada pada data latih dan data uji, hanya fitur itu saja yang akan dihitung ke dalam perhitungan NBC.

Langkah-langkah untuk membangun model menggunakan NBC sebagai berikut.

- Menghitung *prior probability* dari kedua kelas yang ada.

$$P(\text{tidak benar}) = \frac{2}{3} = 0.67 \qquad P(\text{benar}) = \frac{1}{3} = 0.33$$

- Setelah mendapatkan nilai *prior probability* untuk setiap kelas, selanjutnya menghitung *conditional probability*.

$$P(\text{jokowi}|\text{tidak benar}) = \frac{1+1}{8+9} = 0.117 \qquad P(\text{jokowi}|\text{benar}) = \frac{1+1}{3+9} = 0.167$$

$$P(\text{presiden}|\text{tidak benar}) = \frac{1+1}{8+9} = 0.117 \qquad P(\text{presiden}|\text{benar}) = \frac{1+1}{3+9} = 0.167$$

- Setelah didapatkan nilai *conditional probability* data pada masing-masing kelas, akan dilakukan perkalian antara nilai *prior probability* dan *conditional probability* untuk menentukan data tersebut termasuk ke dalam kelas Tidak Benar atau Benar.

$$P(\text{tidak benar}|\text{Dokumen4}) = 0.67 \times 0.117 \times 0.117 = 9.17 \times 10^{-3}$$

$$P(\text{benar}|\text{Dokumen4}) = 0.33 \times 0.167 \times 0.167 = 9.20 \times 10^{-3}$$

Hasil perhitungan akhir dari NBC sudah didapatkan yang kemudian akan dibandingkan nilai probabilitas mana yang paling besar diantara kedua kelas tersebut. Untuk hasil perhitungan dokumen 4 terlihat bahwa nilai probabilitasnya lebih besar di kelas Benar daripada kelas Tidak Benar. Maka dapat disimpulkan dokumen 4 termasuk kelas Benar.

3.5. Evaluasi

Tahap terakhir dari perancangan sistem ini adalah tahap evaluasi. Evaluasi akan menggunakan *F1 Score* untuk menentukan performansi sistem. Sebelum mendapatkan nilai *F1 Score*, dilakukan dulu perhitungan untuk mendapatkan nilai *precision* dan *recall* menggunakan tabel *confusion matrix*. *Precision* digunakan untuk mengukur berapa banyak contoh yang positif dan berlabel positif, sedangkan *recall* digunakan untuk mengukur berapa banyak contoh yang positif dan mendapatkan label yang benar.

4. Evaluasi

Pengujian sistem pada penelitian ini dilakukan untuk melihat performansi *Naïve Bayes Classifier* dalam mengidentifikasi *tweet* hoaks yang berhubungan dengan pemilihan presiden 2019. Pengukuran performansi sistem diukur menggunakan *precision*, *recall*, dan *F1-score*.

4.1 Hasil Pengujian

Hasil pengujian dalam penelitian ini didapat dari mengukur tingkat performansi sistem identifikasi *tweet hoax* menggunakan metode NBC dengan tambahan seleksi fitur MI dan IG. Dalam pengujian ini digunakan sebanyak 400 data latih dan 100 data uji dari total 500 data yang dimiliki. *Threshold* yang digunakan untuk masing-masing seleksi fitur MI dan IG yaitu 0.3. Berikut adalah hasil dari perhitungan *precision*, *recall*, dan *F1-score*.

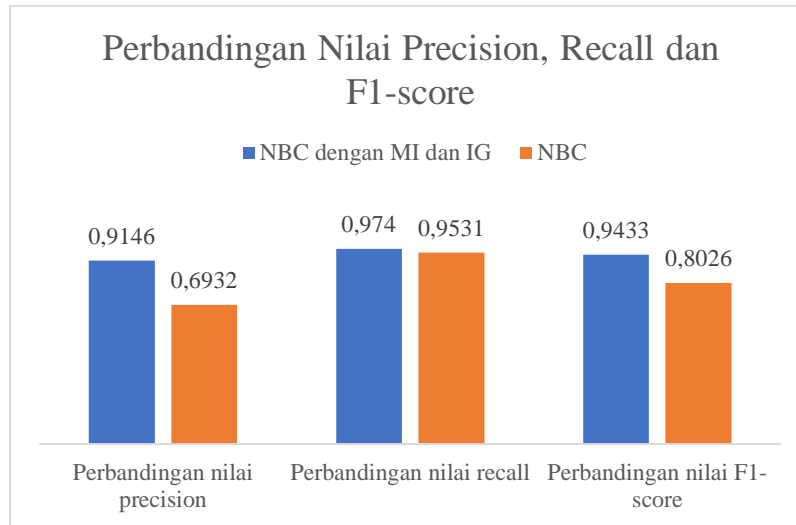
Tabel 13 Hasil Perhitungan Precision, Recall, dan F1 score

<i>Precision</i>	<i>Recall</i>	<i>F1-score</i>
0.9146	0.974	0.9433

Berdasarkan Tabel 6, dapat dilihat bahwa nilai *precision*, *recall*, dan *F1-score* cukup besar dengan rata-rata 90%, hal ini disebabkan oleh penggunaan metode NBC dengan tambahan MI dan IG.

4.2 Analisis Hasil Pengujian

Perbandingan nilai *precision*, *recall*, dan *F1-score* yang dihasilkan oleh sistem yang menggunakan metode NBC dengan tambahan seleksi fitur MI dan IG dan metode NBC yang tidak menggunakan seleksi fitur dapat dilihat pada Gambar 2.



Gambar 2 Perbandingan nilai *precision*, *recall*, dan *F1-score*

Berdasarkan Gambar 2, secara umum nilai yang dihasilkan oleh metode NBC dengan tambahan seleksi fitur MI dan IG lebih besar dibandingkan dengan metode NBC yang tidak menggunakan seleksi fitur dalam mengidentifikasi *tweet* hoaks. Hal ini dikarenakan seleksi fitur dapat memisahkan fitur apa saja yang memiliki informasi lebih penting untuk pengidentifikasian *tweet* hoaks. Jadi, fitur yang digunakan oleh metode NBC tanpa tambahan seleksi fitur memiliki fitur yang tidak berkaitan dengan pilpres lebih banyak dibandingkan dengan metode NBC dengan tambahan seleksi fitur sehingga terjadi kesalahan identifikasi *tweet*. Namun metode NBC dengan tambahan seleksi fitur MI dan IG juga masih memiliki kesalahan identifikasi *tweet*. Hal ini terjadi ketika suatu fitur yang dihasilkan MI dan IG memiliki frekuensi yang besar di kelas Tidak Benar. sehingga ketika perhitungan klasifikasi dengan metode NBC, data yang seharusnya teridentifikasi di kelas Benar secara makna, berakhir teridentifikasi di kelas Tidak Benar dikarenakan frekuensi kemunculan fitur yang ada pada data tersebut lebih besar di kelas Tidak Benar.

5. Kesimpulan dan Saran

Berdasarkan hasil pengujian dan pembahasan analisis yang telah dipaparkan, terdapat beberapa hal yang dapat disimpulkan, yaitu:

1. Metode NBC dapat digunakan untuk mengidentifikasi *tweet* hoaks yang berhubungan dengan pemilihan presiden 2019 dengan tambahan seleksi fitur MI dan IG.
2. Nilai *precision*, *recall*, dan *F1-score* yang dihasilkan dari identifikasi *tweet* hoaks menggunakan metode NBC dengan tambahan seleksi fitur MI dan IG yaitu nilai *precision* 0.9146, nilai *recall* 0.974, dan nilai *F1-score* 0.9433.

Sedangkan untuk pengembangan penelitian selanjutnya, ada beberapa hal yang direkomendasikan, yaitu:

1. Menggunakan tambahan metode lexical chain disertai metode word sense disambiguation yang berguna untuk mengidentifikasi makna yang ada di dalam sebuah data, sehingga pengidentifikasian hoaks dapat dilakukan lebih baik.

Daftar Pustaka

- [1] P. G. Lind, L. R. da Silva, J. S. Andrade Jr. and H. J. Herrmann, "Spreading gossip in social networks," 2008.
- [2] A. Budiman, "Berita Bohong (Hoax) DI Media Sosial Dan Pembentukan Opini Publik," in *Majalah Info Singkat Pemerintahan Dalam Negeri Isu Aktual* 9, 2017, p. 17.
- [3] C. Dewey, "Facebook Fake-News Writer: "I Think Donald Trump is in the White House because of Me"," in *Washington Post*, 2016.
- [4] H. Allcott and M. Gentzkow, "Social Media and Fake News in the 2016 Election," *Journal of Economic Perspectives*, vol. 31, pp. 211-236, 2017.
- [5] L. Zhang, R. Ghosh, M. Dekhil, M. Hsu and B. Liu, "Combining Lexicon-based and Learning-based Methods for Twitter Sentiment Analysis," 2011.
- [6] A. Maulana and H. Situngkir, "Some Inquiries to Spontaneous Opinions: A Case with Twitter in Indonesia," *SSRN Electronic Journal*, 2010.
- [7] H. Situngkir, "Spread of hoax in Social Media ," pp. 1-7, 2010.
- [8] E. Rasywir and A. Purwarianti, "Eksperimen pada Sistem Klasifikasi Berita Hoax Berbahasa Indonesia Berbasis Pembelajaran Mesin," *Cybermatika*, vol. 3, pp. 1-8, 2015.
- [9] K. M. Lhaksana, F. Nhita and A. Budhiarto, "Klasifikasi Pengguna Media Sosial Twitter Dalam Persebaran Hoax Menggunakan Metode Backpropagation," pp. 1-9, 2017.
- [10] P. K. L. Utama, "Identifikasi Hoaks pada Media Sosial dengan Pendekatan Machine Learning," *Jurnal Ilmiah Ilmu Agama dan Ilmu Sosial Budaya*, vol. 13, pp. 69-76, 2018.
- [11] H. Zhang, "The Optimality of Naive Bayes," 2004.
- [12] C. D. Manning, P. Raghavan and H. Schütze, *An Introduction to Information Retrieval*, Cambridge: Cambridge University Press, 2009.
- [13] M. W. Berry and J. Kogan, *Text Mining: Application and Theory*, USA: John Wiley & Sons, 2010.
- [14] P. A. Gayatri, "Citizen Journalism di Twitter," pp. 1-134, 2012.
- [15] J. Davis and M. Goadrich, "The Relationship Between Precision-Recall and ROC Curves," *Appearing in Proceedings of the 23 rd International Conference on Machine Learning*, pp. 233-240, 2006.

