

PREDIKSI TINGKAT KEPERCAYAAN MASYARAKAT TERHADAP PILPRES 2024 MENGGUNAKAN TF-IDF DAN BOW MENGGUNAKAN METODE SVM

Eka Rifut Nur Mustaqim¹, Usman Pagalay², Cahyo Crysdiyan³

^{1,2,3}Jurusan Magister Informatika, Fakultas Sains dan Teknologi, Universitas Islam Malik Maulana Ibrahim, Malang

email: ¹200605210020@student.uin-malang.ac.id,

²usman@mat.uin-malang.ac.id, ³cahyo@ti.uin-malang.ac.id

Kata kunci:

Big Data, Prediksi, TF-IDF, BOW, SVM

ABSTRAK

Dalam era modern ini, dunia maya telah menjadi salah satu aspek yang tak terpisahkan dari kehidupan sehari-hari kita. Dunia maya, atau internet, adalah hasil dari kemajuan teknologi informasi yang telah merevolusi dunia selama beberapa dekade terakhir. Namun, lebih dari sekadar teknologi, ini telah menjadi sebuah ekosistem yang hidup, dihuni oleh miliaran orang yang terhubung, menciptakan dan mengonsumsi informasi. Prediksi pada pemanfaatan big data ini dengan cara kerja mencari dan mengolah data dari segala bentuk ekspresi atau keadaan yang sedang atau telah dialami seseorang user yang diluapkan dalam bentuk teks kedalam media sosial, Prediksi tidak harus memberikan jawaban secara pasti kejadian yang akan terjadi, melainkan berusaha untuk mencari jawaban sedekat mungkin yang akan terjadi. Berdasarkan pada permasalahan yang telah dibahas beberapa teknik yang paling umum dan sering digunakan dalam feature extraction TF-IDF dan BOW, dikarenakan kedua teknik tersebut sangat bersaing serta berperan baik dan sama-sama digunakan untuk merepresentasikan numerik dari data teks serta memiliki kekurangan dan kelebihan masing. Pada penelitian kali ini akan membandingkan kedua metode tersebut yang dipadukan dengan menggunakan metode SVM, untuk model penelitian TF-IDF dengan menggunakan metode SVM mendapatkan hasil Accurasi sebesar 85%, hasil nilai precision sebesar 84%, hasil Recall sebesar 83% dan untuk hasil F1-Score sebesar 83%, sedangkan penelitian menggunakan teknik BOW dengan metode SVM mendapatkan hasil Accurasi sebesar 84%, hasil nilai precision sebesar 79%, hasil Recall sebesar 89% dan untuk hasil F1-Score sebesar 83%.

ABSTRACT

In this modern era, the cyber world has become an inseparable aspect of our everyday lives. The virtual world, or the Internet, is the result of the advances of information technology that have revolutionized the world over the last few decades. Predicting the use of this big data by the way it works to find and process data of all forms of expression or situation that is or has been experienced by a user that is carried out in the form of text into social media, Predictions do not have to give an exact answer to what is going to happen, but rather try to find answers as close as possible to what will happen. In this study we will compare the two methods combined using the SVM method, for the TF-IDF research model using the SVM method obtained the accuracy result of 85%, the recall result of 83% and for the F1-Score result of 83%, while the research using the BOW technique with SVM procedure obtain the accurate result of 84%, the precision result of 79%, the Recall result is 89% and for F1-score result is 83%.

Keywords:

Big Data, Prediction, TF-IDF, BOW, SVM

PENDAHULUAN

Dalam era modern ini, dunia maya telah menjadi salah satu aspek yang tak terpisahkan dari kehidupan sehari-hari kita. Dunia maya, atau internet, adalah hasil dari kemajuan teknologi informasi yang telah merevolusi dunia selama beberapa dekade terakhir. Namun, lebih dari sekadar teknologi, ini telah menjadi sebuah ekosistem yang hidup, dihuni oleh miliaran orang yang terhubung, menciptakan dan mengonsumsi informasi, dan menciptakan budaya digital yang unik. Pengaruh paling pesat dan kita biasa rasakan sekarang ini adalah dalam sektor media sosial, media sosial telah mengubah cara kita berinteraksi, berbagi informasi, dan berkomunikasi dengan dunia di sekitar kita, bertukar informasi bahkan sebagai tempat untuk bercerita dan meluapkan apa yang sedang dirasakannya. Dengan adanya sumber big data tersebut banyak penelitian-penelitian yang memanfaatkan data untuk diproses dan menghasilkan sesuatu yang lebih bermanfaat lagi, salah satunya dengan memanfaatkan machine learning untuk mengolahnya dengan hasil analisis sentimen.

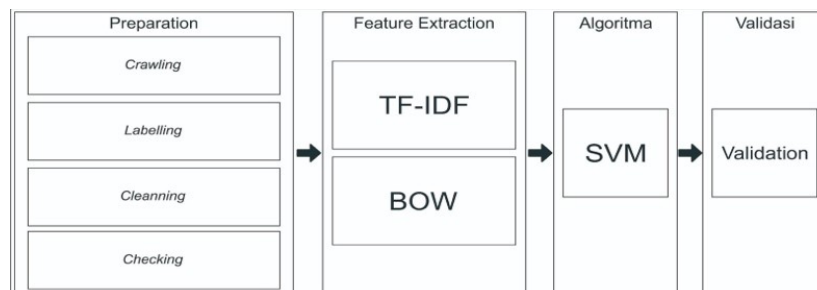
Prediksi pada pemanfaatan big data ini dengan cara kerja mencari dan mengolah data dari segala bentuk ekspresi atau keadaan yang sedang atau telah dialami seseorang user yang diluapkan dalam bentuk teks ke dalam media sosial. Kemudian penelitian menurut Mohammad. Menurut [1] salah satu media sosial yang paling banyak digunakan oleh masyarakat Indonesia saat ini adalah media sosial twitter dengan jumlah persentase 56% setelah youtube, whatsapp, facebook dan Instagram. Kemudian untuk penelitian analisis sentimen dengan data yang bersumber dari aplikasi media social twitter juga sedang ramai dilakukan dan sudah banyak diimplementasikan khususnya di Indonesia [2][3][4][5], prediksi adalah suatu proses memperkirakan secara sistematis tentang sesuatu yang paling mungkin terjadi di masa depan berdasarkan informasi masa lalu dan sekarang yang dimiliki, agar kesalahannya (selisih antara sesuatu yang terjadi dengan hasil perkiraan) dapat diperkecil, Prediksi tidak harus memberikan jawaban secara pasti kejadian yang akan terjadi, melainkan berusaha untuk mencari jawaban sedekat mungkin yang akan terjadi.

Berdasarkan pada permasalahan yang telah dibahas beberapa teknik yang paling umum dan sering digunakan dalam feature extraction yaitu TF-IDF dan BOW, dikarenakan kedua teknik tersebut sangat bersaing serta berperan baik dan sama-sama digunakan untuk merepresentasikan numerik dari data teks serta memiliki kekurangan dan kelebihan masing masing. Pada penelitian kali ini akan membandingkan kedua metode tersebut yang dipadukan dengan menggunakan metode SVM (Support Vector Machine) dengan tujuan untuk mendapatkan kinerja machine learning yang optimal dalam mengolah data teks untuk analisis sentiment dengan data bersumber dari media social twitter pada studi kasus pemilihan presiden 2024.

METODE

Pada penelitian analisis sentimen menggunakan machine learning ini mempunyai rancangan atau alur dari beberapa tahapan yaitu mulai dari tahapan pengumpulan data, pelabelan data, preprocessing data, feature extraction data, klasifikasi data, evaluasi dan pengujian data, perbandingan hasil pengujian dan kesimpulan.

Prediksi Tingkat Kepercayaan Masyarakat Terhadap Pilpres 2024 Menggunakan Tf-Idf dan Bow Menggunakan Metode SVM



Gambar 1. Alur penelitian

Dari diagram alir diatas menunjukkan beberapa langkah untuk melakukan proses analisis sentiment menggunakan machine learning dengan penjelasan berikut:

A. Preparation

Dalam fase ini adalah proses pengumpulan data mentah dari twitter yang diolah menjadi data siap di proses oleh mesin dengan beberapa tahap yang harus dilalui:

1. Crawling

Proses penambangan data atau pengambilan data yang bersumber dari aplikasi social media bernama twitter yang dilakukan dengan bantuan google colab dari penggunaan bahasa pemrograman python, data yang diambil berupa teks dari unggahan tweet para pengguna dengan mengelompokkan berdasarkan topik yang akan diangkat yaitu Pemilihan Presiden 2024 yang dikelompokkan berdasarkan Hastag (#pilpres2024) dengan data berjumlah 300 tweet.

2. Labelling

Proses labeling ini merupakan suatu pelabelan data yang sudah didapatkan dari aplikasi twitter dengan data berjumlah 300 tweet, pada fase ini dilakukan secara manual oleh orang yang berkompeten bidang tim survei pemilu. Pelabelan ini menggunakan dua kategori saja yaitu jika komentar tweet pada satu kalimat mengandung makna atau konotasi yang jelek maka akan diberi label negatif, begitupun sebaliknya jika makna yang baik akan diberi label positif.

3. Cleaning

Proses pembersihan data dari segala bentuk karakter pada kata yang tidak berguna dan tidak memiliki manfaat untuk menghasilkan sentimen maka akan dibuang dan dibersihkan Langkah pada fase ini yaitu case folding atau menyetarakan menjadi huruf kecil, remove punctuation atau penghilangan url, karakter, angka dll pada setiap kalimat, stopword removal atau penghilangan kata hubung seperti di, ke, yang dll dan terakhir adalah stemming yaitu proses menghilangkan kata imbuhan seperti “makan” jadi “makan”.

4. Data checking

Proses persiapan data terakhir adalah data checking yaitu pengecekan data setelah melalui tahap labeling dan cleaning dengan tujuan tidak terdapat data kosong atau blank data, tidak terdapat duplicate data atau data yang ganda serta imbalance data yaitu jumlah data yang seimbang antara kedua label positif maupun negatif.

B. Feature Extraction

Dalam fase ini adalah proses pengekstrakan kata menjadi sebuah angka berupa vektor, dikarenakan komputer hanya dapat mengenali dan memproses angka bukan sebuah kata dengan catatan masih mempunyai arti dari kata tersebut.

1. TF-IDF (*Term Frequency - Inverse Document Frequency*)

Merupakan proses perhitungan atau pengekstrakan kata menjadi sebuah angka berbentuk vektor yang digunakan untuk menentukan bobot dari sebuah kata dalam sebuah dokumen atau korpus. Bobot ini berguna untuk menentukan seberapa penting kata tersebut dalam sebuah dokumen [12]. Pada dasarnya untuk perhitungan atau rumus TF-IDF terbagi menjadi dua yaitu TF (*Term Frequency*) dan IDF (*Inverse Document Frequency*) dengan rumus dan cara kerja yang berbeda dan akan digabungkan di akhir perhitungan antara TF dan IDF sebagai berikut:

a. TF (*Term Frequency*)

Dengan cara kerja menghitung frekuensi jumlah kemunculan kata pada sebuah dokumen. Dalam setiap dokumen memiliki panjang yang berbeda-beda maka nilai TF akan dibagi dengan panjang dokumen:

$$tf_{t,d} = \frac{n_{t,d}}{\text{(Total number of term in document)}} \quad (1)$$

Keterangan

Tf = frekuensi kemunculan kata pada sebuah dokumen

b. IDF (*Inverse Document Frequency*)

Setelah berhasil menghitung nilai TF kita akan menghitung nilai IDF yang merupakan nilai untuk mengukur seberapa penting sebuah kata, dengan beracuan semakin kecil nilai IDF maka akan dianggap semakin tidak penting kata tersebut, begitupun sebaliknya:

$$idf_d = \log \frac{\text{Number of document}}{\text{(Total number of term in document)}} \quad (2)$$

Keterangan

Idf = mengukur penting/tidak sebuah kata dalam dokumen

c. TF-IDF (*Term Frequency - Inverse Document Frequency*)

Setelah mendapatkan nilai TF dan IDF selanjutnya akan menghitung nilai TF-IDF dengan mengalikannya:

$$tfidf_{t,d} = tf_{t,d} \times idf_d \quad (3)$$

Keterangan

TF-IDF = hasil penggabungan antara TF dan IDF

2. BOW (*Bag Of Word*)

Metode ini adalah salah satu metode yang cukup sederhana dalam memproses suatu data teks yang diubah menjadi angka berbentuk vektor agar dapat diproses dan diolah oleh komputer. Metode ini sejatinya hanya menghitung jumlah frekuensi kemunculan kata pada seluruh dokumen yang di proses [13].

- Langkah pertama membuat sekumpulan kata atau *vocabulary* dari seluruh dokumen yang akan di proses
- Dari setiap dokumen, akan dihitung frekuensi kemunculan setiap kata dari dokumen tersebut
- Masukan hasil perhitungan frekuensi kata untuk setiap dokumen dalam sebuah vektor.
- Vector tersebut merepresentasikan dokumen BOW

Dalam notasi matematis BOW dapat direpresentasikan jika d adalah sebuah dokumen dan v adalah sekumpulan kata atau *vocabulary* dari keseluruhan dokumen yang akan direpresentasikan sebagai berikut:

$$\text{BoW}(d) = [\text{count}(w_1,d), \text{count}(w_2,d) \dots \text{count}(w_n,d)] \quad (4)$$

Dimana $\text{count}(w_i,d)$ merupakan jumlah kemunculan kata w_i dalam dokumen d . sedang n adalah jumlah kata dalam V .

B. Algoritma SVM (*Support Vector Machine*)

Merupakan metode pembelajaran yang digunakan mesin untuk proses klasifikasi dan regresi, dengan cara kerja memisahkan dua kelas dengan memaksimalkan margin atau juga disebut jarak antar kelas dengan persamaan atau rumus dasar dari metode SVM adalah sebagai berikut [14].

$$y = f(x) = \text{sign}(w \cdot x + b) \quad (c=1 \quad g=1 \quad \text{karnel}=\text{linear}) \quad (5)$$

Dimana x adalah vektor atau fitur dari sebuah data, w adalah vektor bobot yang mengatur arah dan jarak garis atau *hyperplane*, b adalah bias, sedangkan y adalah label kelas prediksi (-1 atau +1), dan sign adalah fungsi dari signum yang menghasilkan -1 tau +1 tergantung kepada apakah $f(x) > 0$ atau $f(x) < 0$.

C. Validation

Langkah terakhir pada proses klasifikasi dalam memprediksi data menggunakan teks dari twitter ini adalah *validation*, untuk menguji dan mengukur keberhasilan dari teknik dan metode yang telah diterapkan. Pada pengukuran evaluasi model ini menggunakan teknik *confusion matrix* yang menggunakan acuan *accuracy*, *precision*, *recall* dan *f1-Score* dengan ketentuan table berikut.[15]

Tabel I. *confusion matrix*

TP (<i>True Positive</i>)	FP (<i>False Positive</i>)
FN (<i>False Negative</i>)	TN (<i>True Negative</i>)

Keterangan:

- a. TP yaitu jumlah data positif yang terklasifikasi benar
- b. TN yaitu jumlah data negatif yang terklasifikasi benar
- c. FN yaitu jumlah data negatif yang terklasifikasi salah
- d. FP yaitu jumlah data positif yang terklasifikasi salah

Accuracy (A) prediksi benar dari *true positif* dan *true negatif*.

$$A = \frac{(TP+TN)}{(TP+FP+FN+TN)} = 100\% \quad (6)$$

Precision (P) nilai *true positif* dari seluruh nilai *positif*.

$$P = \frac{(TP)}{(TP+FP)} = 100\% \quad (7)$$

Recall (R) persentase prediksi *positif* dengan *true positif*.

$$R = \frac{(TP)}{(TP+FN)} = 100\% \quad (8)$$

F1-Score (F) perbandingan rata-rata *precision* dan *recall*.

$$F = 2 \times (R \times P) : (R + P) \tag{9}$$

D. Komparasi Metode

Dari penerapan Teknik TF-IDF dan menggunakan algoritma SVM (*Support Vector Machine*) akan menghasilkan nilai-nilai yang ada pada *confusion matrix* dengan dibandingkan dengan teknik BOW dengan menggunakan algoritma masih sama dengan sebelumnya yaitu algoritma *Support Vector Machine* SVM, tidak hanya dengan membandingkan nilai ahir saja tetapi pada kesimpulannya juga akan mendapatkan kekurangan dan kelebihan pada masing-masing komparasi metode yang telah dilakukan dengan dua jumlah penelitian.

Tabel II. Komparasi hasil perhitungan *confusion matrix*

Metode	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-Score</i>
TF-IDF + SVM	-%	-%	-%	-%
BOW + SVM	-%	-%	-%	-%

HASIL DAN PEMBAHASAN

A. Persiapan Data

Proses mempersiapkan data mulai dari penambangan data twitter hingga mengolahnya, agar siap untuk dilakukan pemrosesan menggunakan *machine learning*.

1. Crawling twitter

Penambangan data twitter yang telah dilakukan didapat hasil 300 tweet dengan pengelompokan hastag #pilpres2024 pada priode rentan waktu mulai 01 maret 2022 sampai dengan 28 desember 2023, yang berisikan variabel nama akun *account name*, isi teks unggahan *tweet* dan waktu unggah *date* dengan hasil pada tabel berikut:

Tabel III. Hasil penambangan teks tweet

No	Account Name	Tweet	Date
1	@dedy_pram	Pastilah, Prabowo Anies nggak punya prestasi...wkwkwk	03/01/2023 04:04
2	@kumparan	@aniesbaswedan anies...BERBOHONG DAN NIPU TANPA RASA BERSALAH	03/01/2023 04:07
...
300	@yuliaaanidiya	Paslon 02, Prabowo-Gibran: Jangan galau,pilih yang pasti!"dekade08"	03/02/2023 04:30

2. Labelling

Dari hasil penambangan data twitter yang telah dilakukan dengan jumlah data 300 tweet akan dilakukan pelabelan secara manual oleh orang yang berkompeten bidang hukum dengan hasil sebagai berikut:

Tabel IV. Hasil *labelling* teks tweet

No	Tweet	L1	L2	L3	Keputusan
1	Pastilah, Prabowo Anies nggak punya prestasi...wkwkwk	+	-	-	Negatif
2	@aniesbaswedan anies...BERBOHONG DAN NIPU TANPA RASA BERSALAH	-	-	-	Negatif
...
3	Paslon 02, Prabowo-Gibran: Jangan galau,pilih yang pasti!"dekade08"	-	+	+	Positif

Keterangan :

Pelabelan dilakukan oleh 3 orang dengan kode L1,L2 dan L3 dan keputusan label didapatkan jika jumlah label idak kurang dari 2.

3. Cleanning

Pada tahap ini data akan dilakukan proses pembersihan atau penyederhanaan katomendaji baku, mulai dari *case folding* atau menyetarakan semua kata menjadi huruf kecil, *remove punctuation* atau penghilangan url, karakter, angka dll pada setiap kalimat, *stopword removal* atau penghilangan kata hubung dan *stemming* yaitu proses menghilangkan kataimbunan dengan hasil pada tabel berikut:

Tabel V. Hasil *cleannig* teks tweet

No	Tweet	Cleanning
1	Pastilah, Prabowo Anies nggak punya prestasi...wkwkwk	Pastilah Prabowo Anies nggak punya prestasi
2	@aniesbaswedan anies...BERBOHONG DAN NIPU TANPA RASA BERSALAH	Anies berbohong dan nipu tanpa rasa bersalah

...
300	Paslon 02, Prabowo-Gibran: Jangan galau,pilih yang pasti!"dekade08"	Prabowo Gibran Jangan galau pilih yang pasti

4. Data Checking

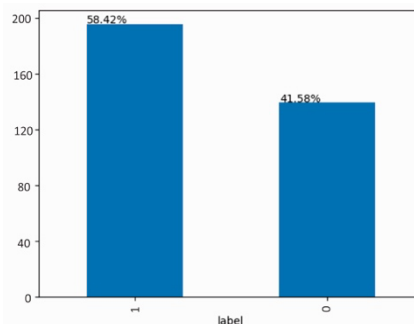
Data yang sudah siap untuk diproses akan dilakukan pengecekan lagi terlebih dahulu, agar data benar-benar valid. Pengecekan ini memiliki hasil yang baik dengan tidak ada data kosong dan duplikat serta jumlah antar kedua label positif dan negative yang seimbang, sebagai berikut:

```

df.dropna(inplace=True)
df['tweet'].isna()
0      False
1      False
2      False
3      False
4      False
...
295    False
296    False
297    False
298    False
299    False
Name: tweet, Length: 300, dtype: bool

df.drop_duplicates(inplace=True)
df.duplicated()
0      False
1      False
2      False
3      False
4      False
...
295    False
296    False
297    False
298    False
299    False
Length: 298, dtype: bool
    
```

Gambar 2. Hasil pengecekan *blank* data dan *duplicated* data



Gambar 3. Hasil pengecekan *imbalance* data

B. Feature Extraction

Proses perubahan dari data berbentuk teks menjadi sebuah vektor berupa angka untuk dapat diproses oleh computer dengan python dan dibuktikan perhitungan manual.

1. TF-IDF (*Term Frequency- Inverse Document Frequency*)

Perhitungan atau rumus TF-IDF terbagi menjadi dua yaitu TF (*Term Frequency*) dan IDF (*Inverse Document Frequency*) dengan rumus dan cara kerja yang berbeda dan akan digabungkan di akhir perhitungan antara TF dan IDF menjadi TF-IDF (*Term Frequency- Inverse Document Frequency*) sebagai berikut:

DATA UTAMA

Prediksi Tingkat Kepercayaan Masyarakat Terhadap Pilpres 2024 Menggunakan Tf-Idf dan Bow Menggunakan Metode SVM

Dokumen 1 (D1)	Pasti tidak mempunyai prestasi	Panjang Dokumen	4
Dokumen 2 (D2)	Punya prestasi	Panjang Dokumen	2
Dokumen 3 (D3)	Elaktibilitas selalu menurun	Panjang Dokumen	3
Total Dokumen			3

Gambar 4. Sempel data perhitungan manual

- TF (*Term Frequency*)

Hasil perhitungan TF dengan menetapkan rumus yang sudah ada pada sub bab sebelumnya:

Token	Frekuensi kemunculan kata			df	TF		
	D1	D2	D3		TF1	TF2	TF3
elaktibilitas	0	0	1	1	0	0	0,5
mempunyai	1	0	0	1	0,25	0	0
menurun	1	0	0	1	0,333333	0	0
pasti	1	0	0	1	0,333333	0	0
prestasi	1	1	0	2	0,333333	0,25	0
punya	0	1	0	1	0	0,25	0
selalu	0	0	1	1	0	0	0,25
tidak	1	0	0	1	0,5	0	0

Gambar 5. Hasil perhitungan *Term Frequency*

Keterangan:

D1,D2 dan D3 merupakan jumlah frekuensi kemunculan kata pada setiap dokumen dengan hasil df, sedangkan TF1,TF2 dan TF3 merupakan notasi untuk setiap riview TF

- IDF (*Inverse Document Frequency*)

Setelah berhasil menghitung nilai TF selanjutnya kita menghitung nilai IDF dari hasil nilai TF:

Token	Frekuensi kemunculan kata			df	TF			IDF
	D1	D2	D3		TF1	TF2	TF3	
elaktibilitas	0	0	1	1	0	0	0,5	0,477121
mempunyai	1	0	1	1	0,25	0	0	0,477121
menurun	1	0	0	1	0,333333	0	0	0,477121
pasti	1	0	0	1	0,333333	0	0	0,477121
prestasi	1	1	0	2	0,333333	0,25	0	0,176091
punya	0	1	0	1	0	0,25	0	0,477121
selalu	0	0	1	1	0	0	0,25	0,477121
tidak	1	0	0	1	0,5	0	0	0,477121

Gambar 6. Hasil perhitungan *Inverse Document Frequency*

Prediksi Tingkat Kepercayaan Masyarakat Terhadap Pilpres 2024 Menggunakan Tf-Idf dan Bow Menggunakan Metode SVM

Keterangan:

Nilai IDF adalah logaritma dari pembagian jumlah dokumen dengan nilai dari TF.

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

Setelah kedua nilai TF dan IDF didapatkan kemudian akan digabungkan menjadi nilai TF-IDF dengan menyusun angka -angka dengan susunan berupa vektor berikut:

Token	Frekuensi kemunculan kata			df	TF			IDF	TF-IDF		
	D1	D2	D3		TF1	TF2	TF3		TF-IDF1	TF-IDF2	TF-IDF3
elaktilabilitas	0	1	0	1	0	0	0,5	0,477121	0	0	0,238561
mempunyai	0	0	1	1	0,25	0	0	0,477121	0,11928	0	0
menurun	1	0	0	1	0,333333	0	0	0,477121	0,15904	0	0
pasti	1	0	0	1	0,333333	0	0	0,477121	0,15904	0	0
prestasi	1	0	1	2	0,333333	0,25	0	0,176091	0,058697	0,044023	0
punya	0	0	1	1	0	0,25	0	0,477121	0	0,11928	0
selalu	0	0	1	1	0	0	0,25	0,477121	0	0	0,11928
tidak	0	1	0	1	0	0,5	0	0,477121	0	0,238561	0

Gambar 7. Hasil perhitungan TF-IDF

Keterangan:

Hasil akhir untuk tf idf didapatkan dari penggabungan antara nilai TF dan nilai IDF.

HASIL VEKTOR TF-IDF	
Dokumen 1	[0, 0, 0.15904, 0.15904, 0.15904, 0, 0, 0]
Dokumen 2	[0.238561, 0, 0, 0, 0, 0, 0, 0.238561]
Dokumen 3	[0, 0.11928, 0, 0, 0.044023, 0.11928, 0.11928, 0]

Gambar 8. Hasil Vektorisasi TF-IDF

Untuk penerapan dengan menggunakan machine learning berbasis *python* adalah sebagai berikut:

```

[3] from sklearn.feature_extraction.text import TfidfVectorizer
2s

[4] tfidf = TfidfVectorizer()
    respons = tfidf.fit_transform(corpus)
    print (respons)

[5] tfidf.get_feature_names()

[6] respons.todense()
0s
matrix([[0.          , 0.          , 0.62276601, 0.62276601, 0.4736296 ,
         0.          , 0.          , 0.          ],
        [0.70710678, 0.          , 0.          , 0.          , 0.          ,
         0.          , 0.          , 0.70710678],
        [0.          , 0.52863461, 0.          , 0.          , 0.40204024,
         0.52863461, 0.52863461, 0.          ]])
    
```

Gambar 9. Hasil vektorisasi menggunakan python

2. BOW (*Bag of Word*)

Dalam implementasinya dari korpus yang ada hanya mengambil kata yang unik saja, setiap kata yang berulang akan ditulis sekali:

Prediksi Tingkat Kepercayaan Masyarakat Terhadap Pilpres 2024 Menggunakan Tf-Idf dan Bow Menggunakan Metode SVM

Review	elaktibilitas	mempunyai	menurun	pasti	prestasi	punya	selalu	tidak
Dokumen 1	0	1	0	1	1	0	0	1
Dokumen 2	0	0	0	0	1	1	0	0
Dokumen 3	1	0	1	0	0	0	1	0

Gambar 10. Hasil perhitungan *Bag of Word*

Keterangan:

Menghitung frekuensi setiap kemunculan kata pada korpus tersebut pada tiga riview sebelumnya. Jika kata tersebut muncul maka diberi nilai satu sebaliknya jika tidak muncul maka diberi nilai 0 dengan hasil vector berikut:

HASIL VEKTOR BAG OF WORD	
Dokumen 1	[0, 1, 0, 1, 1, 0, 0, 1]
Dokumen 2	[0, 0, 0, 0, 1, 1, 0, 0]
Dokumen 3	[1, 0, 1, 0, 0, 0, 1, 0]

Gambar 11. Hasil Vektorisasi BOW

Untuk hasil penerapan dengan menggunakan machine learning berbasis *python* adalah sebagai berikut:

```
[3] from sklearn.feature_extraction.text import CountVectorizer

[5] bow = CountVectorizer()
    responsbow = bow.fit_transform(corpus)
    print (responsbow)

[8] bow.get_feature_names()

[7] responsbow.todense()

matrix([[ 0,  1,  0,  1,  1,  0,  0,  1],
        [ 0,  0,  0,  0,  1,  1,  0,  0],
        [ 1,  0,  1,  0,  0,  0,  1,  0]])
```

Gambar 12. Hasil vektorisasi menggunakan python

C. Algorithm SVM (*Support Vector Machine*)

Setelah mendapatkan hasil vektor dari data yang sudah melewati beberapa tahapan, pada proses SVM dengan cara kerja memisahkan dua kelas dengan memaksimalkan margin. Pada proses metode SVM ini dilakukan untuk setiap percobaan feature extraction TF-IDF maupun BOW menggunakan konfigurasi standart sebagai berikut:

```
[ ] from sklearn.svm import SVC

[ ] model1 = SVC(C=1, gamma=1, kernel='linear')
    model1.fit(X_train, y_train)

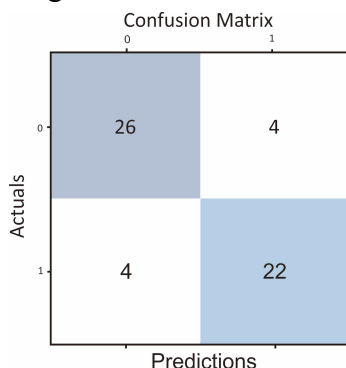
SVC(C=1, gamma=1, kernel='linear')
```

Gambar 13. Metode SVM di python

D. Validation

Proses menguji dan mengukur keberhasilan dari teknik dan metode yang telah diterapkan untuk masing-masing percobaan yaitu antara TF-IDF dengan metode SVM (*Support Vector Machine*) dan BoW juga dengan menggunakan metode SVM (*Support Vector Machine*):

- Berikut hasil validasi TF-IDF dengan metode SVM:



Gambar 14. *Confusion matrix* TF-IDF

Keterangan hasil:

- a. Berjumlah 26 data positif yang terklasifikasi benar
- b. Berjumlah 22 data negatif yang terklasifikasi benar
- c. Berjumlah 4 data negatif yang terklasifikasi salah
- d. Berjumlah 4 data positif yang terklasifikasi salah

Hasil dari *confusion matrix* akan digunakan sebagai acuan dalam mencari *Accuracy*, *Precision*, *Recall* dan *F1-score*:

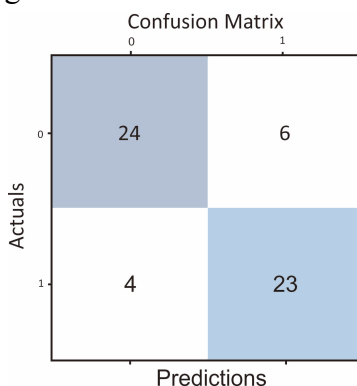
Accuracy (0.85) prediksi benar dari *TP* dan *TN*.

Precision (0.84) nilai *true positif* dari seluruh nilai *positif*.

Recall (0.83) persentase prediksi *positif* dengan *true positif*.

F1-Score (0.83) perbandingan rata-rata *precision* dan *recall*.

- Berikut hasil validasi BOW dengan metode SVM:



Gambar 15. *Confusion matrix* BOW

Keterangan hasil:

- a. Berjumlah 24 data positif yang terklasifikasi benar
- b. Berjumlah 23 data negatif yang terklasifikasi benar
- c. Berjumlah 4 data negatif yang terklasifikasi salah
- d. Berjumlah 6 data positif yang terklasifikasi salah

Hasil dari *confusion matrix* akan digunakan sebagai acuan dalam mencari *Accuacy*, *Precision*, *Recall* dan *F1-score*:

Accuracy (0.84) prediksi benar dari *TP* dan *TN*.

Precision (0.79) nilai *true positif* dari seluruh nilai *positif*.

Recall (0.89) persentase prediksi *positif* dengan *true positif*.

F1-Score (0.83) perbandingan rata-rata *precision* dan *recall*.

E. Komparasi Metode

Membandingkan hasil uji validasi guna untuk mendapatkan nilai hasil yang paling optimal antara kedua teknik yaitu TF-IDF dan BOW untuk klasifikasi analisis sentimen menggunakan metode SVM dengan hasil komparasi sebagai berikut.

Tabel VI. Komparasi hasil *confusion matrix*

Metode	<i>TP</i>	<i>TN</i>	<i>FN</i>	<i>FP</i>
TF-IDF + SVM	26	22	4	4
BOW + SVM	24	23	4	6

Didapatkan hasil komentar positif yang diprediksi benar lebih unggul pada TF-IDF + SVM daripada komentar negative yang di prediksi salah lebih besar pada BOW +SVM sedangkan untuk tipe eror FN dan FP pada kasus ini lebih berpihak kepada dengan komentar positif yang terdeteksi sebagai komentar negative.

Tabel VI. Komparasi hasil validasi

Metode	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-Score</i>
TF-IDF + SVM	85%	84%	83%	83%
BOW + SVM	84%	79%	89%	83%

Didapatkan hasil rasio prediksi benar baik komentar negative dan positif dengan selisih 1 angka dan prediksi benar positif di bandingkan dengan keseluruhan hasil yang diprediksi positif dengan selisih 6 angka dan perbandingan rata-rata dengan presisi yang dibobotkan selisih 1 angka yang semuanya unggul pada TF-IDF+BOW, Tetapi benar positif dibandingkan dengan keseluruhan data yang benar positif unggul pada BOW+SVM.

KESIMPULAN

Kesimpulan dari hasil perbandingan kedua teknik dan metode yang diterapkan yaitu teknik feature extraction TF-IDF dengan metode SVM dan teknik feature extraction BOW juga dengan metode SVM dalam memprediksi data terkait pemilihan presiden sebagai berikut:

- Untuk model penelitian TF-IDF dengan SVM mendapatkan hasil Accurasi sebesar 85%, hasil nilai precision sebesar 84%, hasil Recall sebesar 83% dan untuk hasil F1-Score sebesar 83%. Dengan nilai yang didapatkan pada komparasi ini sangat bagus yang menggunakan system mempertimbangkan frekuensi kemunculan kata dalam suatu dokumen.
- Untuk model penelitian BOW dengan SVM mendapatkan hasil Accurasi sebesar 84%, hasil nilai precision sebesar 79%, hasil Recall sebesar 89% dan untuk hasil F1-Score sebesar 83%. Dengan nilai yang didapatkan pada nilai Recall untuk komparasi teknik ini dapat mengungguli TF-IDF+BOW dikarenakan memiliki perhitungan dengan mempertahankan urutan kata pada setiap dokumen.

Dengan perolehan nilai-nilai tersebut untuk hasil yang maksimal dalam melakukan proses analisis sentimen dalam memprediksi data menggunakan teks twitter yaitu dengan menerapkan teknik feature extraction TF-IDF dan SVM yang memiliki keunggulan lebih baik pada setiap pengukurannya baik Accuracy, Precision maupun F1-Score sedangkan pada teknik BOW dan SVM hanya unggul untuk nilai Recall saja.

DAFTAR PUSTAKA

- [1] M. Fadilah Arfat et al., “Analisis Sentimen Masyarakat Indonesia Terkait Vaksin Covid-19 Pada Media Sosial Twitter Menggunakan Metode Support Vector Machine (SVM),” vol. 7, no. 2, 2022.
- [2] Y. Romadhoni, K. Fahmi, and H. Holle, “Analisis Sentimen Terhadap PERMENDIKBUD No.30 pada Media Sosial Twitter Menggunakan Metode Naive Bayes dan LSTM,” vol. 7, no. 2, 2022.
- [3] C. Ayunda et al., “Analisis Komparasi Algoritma Machine Learning untuk Sentiment Analysis (Studi Kasus: Komentar YouTube ‘Kekerasan Seksual’),” vol. 7, no. 2, 2022.
- [4] R. Aprillya, P. : Perbandingan, M. Klasifikasi, R. A. Putri, and N. S. Fatonah, “Perbandingan Metode Klasifikasi serta Analisis Faktor Akademis Pola Kelulusan Mahasiswa di Perguruan Tinggi,” vol. 7, no. 2, 2022.
- [5] A. Setyawinda, B. Setiyadi, and A. D. Hartanto, “Perbandingan Algoritma Word Matching dan Naive Bayes untuk Klasifikasi Sentimen Analisis Komentar Instagram,” vol. 5, no. 1, 2020.
- [6] O. I. Gifari, M. Adha, I. Rifky Hendrawan, F. Freddy, and S. Durrand, “Analisis Sentimen Review Film Menggunakan TF-IDF dan Support Vector Machine,” JIFOTECH (JOURNAL OF INFORMATION TECHNOLOGY), vol. 2, no. 1, 2022.
- [7] “JOURNAL OF INTELLIGENT SYSTEMS AND COMPUTATION 43.” [Online]. Available: <https://t.co/9Wl0aWpfd5>
- [8] C. H. Yutika, A. Adiwijaya, and S. al Faraby, “Analisis Sentimen Berbasis Aspek pada

- Review Female Daily Menggunakan TF-IDF dan Naïve Bayes,” JURNAL MEDIA INFORMATIKA BUDIDARMA, vol. 5, no. 2, p. 422, Apr. 2021, doi: 10.30865/mib.v5i2.2845.
- [9] M. M. Munir, M. A. Fauzi and R. S. Perdana, “Implementasi Metode Backpropagation Neural Network berbasis lexion Based Features dan Bag of Words Untuk Identifikasi Ujaran Kebencian Pada Twitter,” vol. 2, no. 10, 2018.
- [10] Hartanto, “TEXT MINING DAN SENTIMEN ANALISIS TWITTER PADA GERAKAN LGBT,” vol. 9, no. 1, 2017.
- [11] A. P. Wibawa, M. Guntur, A. Purnama, M. Fathony Akbar, and F. A. Dwiyanto, “Metode-metode Klasifikasi,” Prosiding Seminar Ilmu Komputer dan Teknologi Informasi, vol. 3, no. 1, 2018.
- [12] D. Farah Zhafira, B. Rahayudi, and P. Korespondensi, “ANALISIS SENTIMEN KEBIJAKAN KAMPUS MERDEKA MENGGUNAKAN NAIVE BAYES DAN PEMBOBOTAN TF-IDF BERDASARKAN KOMENTAR PADA YOUTUBE,” 2021.
- [13] J. Muara Sains, dan Ilmu Kesehatan, W. Trisari Harsanti Putri, and R. Hendrowati, “PENGALIAN TEKS DENGAN MODEL BAG OF WORDS TERHADAP DATA TWITTER,” vol. 2, no. 1, pp. 129–138, 2018.
- [14] N. Hendrastuty, A. Rahman Isnain, and A. Yanti Rahmadhani, “Analisis Sentimen Masyarakat Terhadap Program Kartu Prakerja Pada Twitter Dengan Metode Support Vector Machine,” vol. 6, no. 3, 2021, [Online]. Available: <http://situs.com>
- [15] D. Normawati and S. A. Prayogi, “Implementasi Naïve Bayes Classifier Dan Confusion Matrix Pada Analisis Sentimen Berbasis Teks Pada Twitter,” 2021.