

**TUGAS AKHIR**

**IMPLEMENTASI METODE ANALISIS SENTIMEN  
DAN DETEKSI SARKASME MENGGUNAKAN DATA  
MEDIA SOSIAL**



**SAFINA**

**D42116011**

**DEPARTEMEN TEKNIK INFORMATIKA**

**UNIVERSITAS HASANUDDIN**

**2020**

**TUGAS AKHIR**

**IMPLEMENTASI METODE ANALISIS SENTIMEN  
DAN DETEKSI SARKASME MENGGUNAKAN DATA  
MEDIA SOSIAL**



**SAFINA**

**D42116011**

**DEPARTEMEN TEKNIK INFORMATIKA  
UNIVERSITAS HASANUDDIN**

**2020**

LEMBAR PENGESAHAN SKRIPSI

"Implementasi Metode Analisis Sentimen dan Deteksi Sarkasme  
Menggunakan Data Media Sosial"

OLEH:

Safina

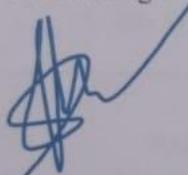
D42116011

Skripsi ini telah dipertahankan pada Ujian Akhir tanggal 30 November 2020.  
Diterima dan disahkan sebagai salah satu syarat memperoleh gelar Sarjana Teknik  
(S.T) pada Progran Studi Strata-1 Teknik Informatika Fakultas Teknik Universitas  
Hasanuddin

Makassar, 30 November 2020

Disetujui oleh:

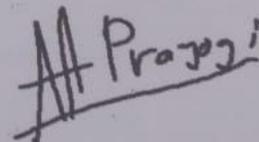
Pembimbing I



Dr. Amil Ahmad Ilham, S.T, M.IT.

NIP. 19731010 199802 1 001

Pembimbing II



A. Ais Prayogi A., S.T, M.Eng.

NIP. 19830510 201404 1 001

Diterima dan disahkan oleh:

Ketua Program Studi S1 Teknik Informatika



Dr. Amil Ahmad Ilham, S.T, M.IT.

NIP. 19731010 199802 1 001

## PERNYATAAN KEASLIAN

Saya yang bertanda tangan dibawah ini,

Nama : Safina  
NIM : D42116011  
Jurusan/program studi : Teknik Informatika

dengan ini menyatakan dengan sebenar-benarnya bahwa skripsi yang berjudul:

Implementasi Metode Analisis Sentimen dan Deteksi Sarkasme Menggunakan Data Media Sosial

Adalah karya ilmiah saya sendiri dan sepanjang pengetahuan saya di dalam naskah skripsi ini tidak terdapat karya ilmiah yang pernah diajukan oleh orang lain untuk memperoleh gelar akademik di suatu perguruan tinggi, dan tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain kecuali yang secara tertulis dikutip dalam naskah ini dan disebutkan dalam sumber kutipan dan daftar pustaka

Apabila di kemudian hari ternyata di dalam naskah skripsi ini dapat dibuktikan terdapat unsur-unsur jiplakan, saya bersedia menerima sanksi atas perbuatan tersebut dan diproses sesuai dengan peraturan perundang-undangan yang berlaku (UU No. 20 Tahun 2003, pasal 25 ayat 2 dan pasal 70)

Makassar 30 November 2020

Yang membuat pernyataan



Safina

## KATA PENGANTAR

Segala puji dan syukur penulis panjatkan kepada Tuhan Yang Maha Esa karena atas nikmat dan rahmat-Nya sehingga penulis dapat menyelesaikan skripsi ini. Penulisan skripsi ini dilakukan dalam rangka memenuhi syarat untuk memperoleh gelar Sarjana Strata Satu pada Departemen Teknik Informatika Fakultas Teknik Universitas Hasanuddin. Dalam proses penyelesaian skripsi ini, penulis memperoleh banyak bantuan dan dukungan dari berbagai pihak, oleh karena itu pada kesempatan ini penulis ingin berterima kasih kepada:

1. Kedua orang tua penulis, Bapak Saifuddin dan Ibu Zumana yang selalu mendoakan dan memberikan dukungan kepada penulis dalam menyelesaikan tugas akhir ini.
2. Bapak Dr. Amil Ahmad Ilham S.T, M.IT, selaku dosen pembimbing I yang telah memberikan bimbingan dan masukan dalam penyusunan tugas akhir ini
3. Bapak A.Ais Prayogi Alimuddin, S.T., M.Eng, selaku dosen pembimbing II yang telah memberikan bimbingan dan masukan dalam penyusunan tugas akhir ini
4. Bapak Amil Ahmad Ilham S.T, M.IT, Ph.D selaku ketua Departemen Teknik Informatika Fakultas Teknik Universitas Hasanuddin
5. Seluruh dosen Departemen Fakultas Teknik Informatika Fakultas Teknik Universitas Hasanuddin
6. Seluruh staf Departemen Teknik Informatika Fakultas Teknik Universitas Hasanuddin
7. Teman-teman Teknik Informatika Angkatan 2016 selaku rekan belajar selama masa perkuliahan
8. Seluruh pihak yang tak sempat penulis sebutkan satu persatu yang telah membantu dalam penyusunan tugas akhir ini

Akhir kata, penulis menyadari bahwa masih banyak kekurangan dalam penulisan skripsi ini. Oleh karena itu, penulis mengharapkan kritik dan saran yang bersifat membangun untuk kesempurnaan skripsi ini. Semoga skripsi ini dapat bermanfaat bagi pembaca.

Makassar, 25 November 2020

Safina

## ABSTRAK

Sebagai salah satu sumber data terbesar, percakapan di media sosial dapat dimanfaatkan untuk berbagai keperluan, salah satunya adalah analisis sentimen. Melalui analisis sentimen kalimat dapat dikelompokkan menjadi sentimen positif, negatif, atau netral. Kalimat yang mengandung sarkasme mengakibatkan hasil klasifikasi sentimen menjadi kurang akurat. Pada penelitian ini, pengaruh deteksi sarkasme terhadap akurasi analisis sentimen diuji. Penelitian ini menggunakan data dari Twitter berbahasa Indonesia. Fitur yang digunakan untuk deteksi sarkasme adalah sentiment score awal dan akhir kalimat serta kata interjeksi. Algoritma yang digunakan untuk deteksi sarkasme adalah *Random Forest*. Analisis sentimen dilakukan menggunakan fitur TF-IDF dan algoritma *Naïve bayes*. Hasil pengujian menunjukkan bahwa deteksi sarkasme tidak meningkatkan akurasi klasifikasi sentimen. Rata-rata akurasi klasifikasi sentimen adalah 74.3 % .Rata-rata akurasi analisis sentimen dengan deteksi sarkasme adalah 72.9%. Rata-rata akurasi deteksi sarkasme adalah 71.7%.

**Kata Kunci:** Analisis sentimen, Deteksi Sarkasme, Random Forest, Naïve Bayes

# DAFTAR ISI

HALAMAN PENGESAHAN.....	iii
PERNYATAAN KEASLIAN.....	iv
KATA PENGANTAR .....	iv
ABSTRAK .....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	x
DAFTAR GAMBAR .....	xi
DAFTAR LAMPIRAN.....	xii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang .....	1
1.2 Rumusan Masalah .....	2
1.3 Tujuan Penelitian.....	2
1.4 Manfaat Penelitian.....	2
1.5 Batasan Masalah.....	2
1.6 Sistematika Penulisan.....	2
BAB II TINJAUAN PUSTAKA.....	3
2.1 Penelitian Terkait .....	3
2.2 Twitter .....	3
2.3 Analisis Sentimen.....	4
2.4 Sarkasme .....	4
2.5 TF-IDF.....	5
2.6 Interjeksi.....	6
2.7 Sentiwordnet.....	7
2.8 Naïve Bayes.....	7
2.9 Random Forest .....	8
2.10 CART .....	9
2.11 Confussion Matrix .....	10
2.12 K-Fold Cross Validation .....	10

BAB III METODOLOGI PENELITIAN.....	12
3.1    Deksripsi Umum.....	12
3.2    Pengumpulan Data .....	15
3.3    Preprocessing Data .....	15
3.4    Ekstraksi Fitur Deteksi Sarkasme .....	18
3.5    Ekstraksi Fitur Analisis sentimen.....	20
3.6    Analisis Sentimen.....	21
3.7    Deteksi Sarkasme .....	23
3.8    Sentiment Reversal.....	24
3.9    Pengujian .....	25
BAB IV HASIL PENELITIAN DAN PEMBAHASAN.....	26
4.1    Dataset .....	26
4.2    Analisis Sentimen.....	26
4.3    Deteksi Sarkasme .....	33
4.4    Analisis Sentimen dengan Deteksi Sarkasme.....	43
4.5    Perbandingan Akurasi Sistem dengan Library.....	45
4.6    Pengembangan Website Sentimen Analisis .....	46
BAB V KESIMPULAN DAN SARAN.....	48
5.1    Kesimpulan.....	48
5.2    Saran.....	48
DAFTAR PUSTAKA .....	49
LAMPIRAN.....	52

## DAFTAR TABEL

Tabel 3. 1 Contoh Preprocessing Data.....	17
Tabel 4. 1 Rincian Dataset .....	26
Tabel 4. 2 Contoh Data Train.....	27
Tabel 4. 3 Jumlah Kemunculan Kata pada Tweet Pertama .....	28
Tabel 4. 4 Contoh Perhitungan Nilai TF pada Tweet Pertama .....	28
Tabel 4. 5 Bobot TF Data.....	28
Tabel 4. 6 Bobot IDF .....	29
Tabel 4. 7 Contoh Perhitungan TF-IDF pada Tweet Pertama .....	29
Tabel 4. 8 Bobot TF-IDF .....	29
Tabel 4. 9 Data Input.....	30
Tabel 4. 10 Mean Data .....	30
Tabel 4. 11 Variance Data.....	31
Tabel 4. 12 Data Test .....	31
Tabel 4. 13 Bobot TF-IDF Data Test .....	31
Tabel 4. 14 Probabilitas Atribut pada Setiap Kelas .....	32
Tabel 4. 15 Contoh Data Train.....	36
Tabel 4. 16 Data Bootsrap.....	38
Tabel 4. 17 Contoh Data Test .....	41
Tabel 4. 18 Hasil Klasifikasi Trees dan Random Forest.....	42
Tabel 4. 19 Rata-Rata Akurasi Deteksi Sarkasme .....	43
Tabel 4. 20 Akurasi Analisis Sentimen dengan Deteksi Sarkasme .....	43
Tabel 4. 21 Perbandingan Akurasi Sistem dengan Library .....	45

## DAFTAR GAMBAR

Gambar 2. 1 Pembagian Data Train dan Data Test dalam 5-Fold Cross Validation ..	11
Gambar 3. 1 Flowchart Deteksi Sarkasme .....	12
Gambar 3. 2 Flowchart Analisis sentimen .....	14
Gambar 3. 4 Flowchart Menghitung Sentimen Awal dan Akhir Tweet .....	19
Gambar 3. 3 Flowchart TF-IDF .....	20
Gambar 3. 5 FLOWchart Algoritma Naive Bayes.....	22
Gambar 3. 6 Flowchart Algoritma Random Forest.....	23
Gambar 3. 7 Flowchart Sentiment Reversal .....	25
Gambar 4. 1 Pengaruh Jumlah Kata terhadap Akurasi Analisis Sentimen .....	33
Gambar 4. 2 Program Translate Data.....	33
Gambar 4. 3 Program untuk Membagi Data .....	34
Gambar 4. 4 Program Menghitung Sentiment Score Data.....	35
Gambar 4. 5 Program Hitung Interjeksi .....	36
Gambar 4. 6 Program Bootstrap Sampling .....	37
Gambar 4. 7 Program Random Forest .....	38
Gambar 4. 8 Program Split Tree .....	39
Gambar 4. 9 Program Menghitung Nilai Gini .....	39
Gambar 4. 10 Program Cari Split Point Terbaik.....	40
Gambar 4. 11 Decision Tree .....	40
Gambar 4. 12 Program Prediksi Random Forest .....	41
Gambar 4. 13 Program Prediksi Data pada Tree.....	42
Gambar 4. 14 Hasil Prediksi Tree .....	42
Gambar 4. 15 Halaman Awal.....	46
Gambar 4. 16 Halaman Hasil .....	47

## DAFTAR LAMPIRAN

<b>Lampiran 1</b> Tabel Keyword dan Tanggal Crawling Data .....	53
<b>Lampiran 2.</b> Class DataPreparation.....	54
<b>Lampiran 3.</b> Kelas TF-IDF .....	57
<b>Lampiran 4.</b> Kelas Naïve Bayes .....	58
<b>Lampiran 5.</b> Kelas Sentiment Score .....	60
<b>Lampiran 6.</b> Kelas Random Forest.....	61
<b>Lampiran 7.</b> Kelas CART.....	62

# BAB I

## PENDAHULUAN

### 1.1 Latar Belakang

Di era globalisasi ini, jumlah data yang dihasilkan setiap harinya sangat besar. Salah satu yang menjadi sumber data terbesar adalah sosial media, seperti Instagram, Facebook, Twitter dan lain-lain. Sebagai salah satu raksasa sosial media, rata-rata 500 juta *tweet* diunggah ke Twitter setiap harinya [1]. *Tweet* yang diunggah sering kali berisi aspirasi dari pengguna. Hal ini dapat dimanfaatkan untuk berbagai tujuan, salah satunya adalah untuk memperoleh gambaran umum dari opini publik terhadap topik-topik tertentu, misalnya bisnis dan politik. Untuk mengetahui hal tersebut, dapat dilakukan analisis sentimen.

Percakapan dari Twitter dapat diklasifikasikan berdasarkan sentimennya. Melalui analisis sentimen, data dapat dikelompokkan menjadi sentimen positif, negatif, dan netral. Ada beberapa hal yang mengakibatkan klasifikasi sentimen menjadi kurang akurat, salah satunya adalah sarkasme. Sarkasme adalah bentuk mengekspresikan perasaan negatif menggunakan kata-kata positif. Kalimat sarkasme biasa ditemukan pada topik pemerintahan, merek, dan politik. Berdasarkan pengamatan yang dilakukan oleh Lunando dan Purwarianti [2], dari 100 *tweet* yang bertopik pemerintahan, merek dan politik, terdapat 18 *tweet* sarkasme. Sarkasme mengakibatkan polaritas kalimat berubah dari positif ke negatif dan sebaliknya. Oleh karena itu, penelitian ini menguji pengaruh deteksi sarkasme terhadap akurasi klasifikasi sentimen.

Menurut Rajadesingan dkk.[3] adanya sentimen yang saling kontras pada suatu kalimat merupakan salah satu ciri sarkasme. Ciri lain dari kalimat sarkasme adalah penggunaan kata interjeksi, seperti yang disebutkan dalam penelitian Bouazizi dan Ohtsuki[4]. Oleh karena itu, pada penelitian ini *sentiment score* awal dan akhir kalimat serta kata interjeksi digunakan sebagai fitur untuk mendeteksi sarkasme yang terdapat pada suatu kalimat.

## **1.2 Rumusan Masalah**

1. Bagaimana melakukan analisis sentimen pada komentar Twitter?
2. Bagaimana mendeteksi sarkasme pada komentar Twitter?
3. Bagaimana pengaruh deteksi sarkasme terhadap akurasi analisis sentimen?

## **1.3 Tujuan Penelitian**

1. Untuk mendeteksi sarkasme yang terdapat pada komentar Twitter
2. Untuk mengetahui pengaruh deteksi sarkasme terhadap akurasi analisis sentimen

## **1.4 Manfaat Penelitian**

Manfaat dari penelitian ini adalah adanya gambaran mengenai pengaruh deteksi sarkasme terhadap akurasi analisis sentimen pada data sosial media.

## **1.5 Batasan Masalah**

Data yang digunakan adalah data dari Twitter berbahasa Indonesia

## **1.6 Sistematika Penulisan**

BAB I PENDAHULUAN berisi latar belakang, tujuan penelitian, manfaat penelitian, batasan masalah dan sistematika penulisan.

BAB II TINJAUAN PUSTAKA berisi penelitian terkait dan dasar teori yang digunakan untuk menyusun tugas akhir ini.

BAB III METODOLOGI PENELITIAN berisi tentang langkah-langkah penelitian yang dilakukan.

BAB IV HASIL PENELITIAN DAN PEMBAHASAN berisi hasil dan pembahasan dari penelitian yang dilakukan.

BAB V KESIMPULAN DAN SARAN berisi kesimpulan dari penelitian dan saran untuk pengembangan penelitian.

## **BAB II**

### **TINJAUAN PUSTAKA**

#### **2.1 Penelitian Terkait**

Penelitian yang dilakukan oleh Prasad, dkk (2017): *Sentiment Analysis for Sarcasm Detection on Streaming Short Text Data* [5]. Penelitian ini membandingkan performa beberapa algoritma klasifikasi, yaitu *Random Forest, Gradient Boosting, Decision Tree, Adaptive Boost, Logistic Regression and Gaussian Naïve Bayes* dalam mendeteksi sarkasme pada komentar. Pada penelitian ini, kamus emoji dan bahasa slang digunakan digunakan dalam *preprocessing data*. Fitur *blob polarity, blob subjectivity, capitalization, positive sentiment, negative sentiment* digunakan untuk mendeteksi sarkasme pada komentar.

Penelitian yang dilakukan oleh Septiani dan Sibaroni (2019): *Sentiment Analysis Terhadap Tweet Bernada Sarkasme Berbahasa Indonesia* [6]. Penelitian ini menggunakan fitur unigram dan interjeksi untuk mendeteksi sarkasme. Klasifikasi dilakukan menggunakan algoritma *SVM* dan *Naïve Bayes*. Hasil penelitian menunjukkan penggunaan fitur interjeksi dapat meningkatkan akurasi klasifikasi. Hasil akurasi tertinggi diperoleh menggunakan algoritma *Naïve Bayes* dengan fitur ekstraksi *unigram* dan interjeksi.

Penelitian yang dilakukan oleh Yunitasari, Musdholifah dan Sari (2019): *Sarcasm Detection For Sentiment Analysis in Indonesian Tweets* [7]. Penelitian ini menggunakan unigram dan 4 set fitur boazizi untuk deteksi sarkasme dan TF-IDF sebagai fitur untuk analisis sentimen. Klasifikasi sentimen dilakukan menggunakan algoritma *Naïve Bayes*.

#### **2.2 Twitter**

Twitter adalah layanan sosial media yang memungkinkan penggunaanya untuk mengirim pesan yang biasa disebut dengan *tweet*. Pesan yang dikirimkan

dapat berupa teks, gambar, video dan audio. Twitter pertama kali didirikan pada tahun 2006 oleh Jack Dorsey.

Twitter menyediakan akses data kepada perusahaan, pengembang dan pengguna melalui API (*Application User Interface*). Salah satu API yang disediakan oleh twitter adalah *search API*. Search API memungkinkan perusahaan, pengembang atau pengguna untuk memperoleh *tweet* yang telah diunggah ke twitter [8].

### **2.3 Analisis Sentimen**

Analisis sentimen adalah bidang ilmu yang mempelajari tentang opini publik terhadap suatu topik, baik berupa produk, individu, masalah, dan sebagainya[9]. Analisis sentimen merupakan salah satu aplikasi dari *Natural Language Processing*, yang digunakan untuk mengidentifikasi dan mengklasifikasi pendapat dari suatu sumber. Secara garis besar, analisis sentimen bertujuan untuk menentukan sikap dari penulis terhadap suatu topik atau polaritas kontekstual keseluruhan dokumen. Sikap dari penulis dapat berupa penilaian, keadaan emosi penulis saat menulis, atau emosi yang ingin disampaikan penulis kepada pembaca [10].

Tujuan analisis sentimen adalah untuk mengidentifikasi dan mengklasifikasi apakah suatu dokumen atau kalimat mengungkapkan pendapat dan apakah pendapat tersebut bersifat positif, negatif atau netral. Dengan kata lain, tujuan analisis sentimen yaitu mengidentifikasi apakah suatu teks mengekspresikan pendapat atau tidak, dan menentukan orientasi teks yang mengekspresikan pendapat [10].

### **2.4 Sarkasme**

Sarkasme berasal dari bahasa Yunani '*sarkasmos*' yang merupakan turunan dari kata kerja '*sarkein*' yang berarti 'merobek-robek daging seperti anjing', 'menggigit bibir karena marah', atau berbicara dengan kepahitan. Majas sarkasme termasuk majas pertentangan. Majas Sarkasme merupakan turunan dari majas ironi (Lee dan Katz, 1998)[11]. Majas ironi adalah sindiran yang

mengungkapkan suatu makna berlainan dengan rangkaian katanya[12]. Majas sarkasme bersifat biasanya bersifat menyindir dan lebih kasar dari majas ironi [11].

Seseorang yang menggunakan majas sarkasme biasanya mengekspresikan kata-kata yang bersifat terlalu positif, namun bermaksud sebaliknya. Ucapan sarkasme memiliki makna literal positif, tetapi memiliki maksud negatif. Sarkasme merupakan bentuk dari majas ironi, yang muncul dalam konteks komunikasi interpersonal, baik secara vokal, atau tertulis, dan sering kali ditukukan untuk mengejek orang tertentu [13].

Menurut Joana Maren Hjele Oslen, yang dikutip dari [11], ada beberapa penanda majas sarkasme, sebagai berikut:

1. *Self-contradiction*

Penanda *self-Contradiction* dapat berupa interjeksi, kalimat emotif dan evaluatif, kata-kata berlebihan dan super relatif, dan pasangan kata kerja dan kata sifat positif

2. *Hyperbolic Combination*

Ungkapan Hiperbola juga dapat digunakan sebagai penanda majas sarkasme. Kata kata yang biasa digunakan untuk majas hiperbola adalah *excessive adjective*, seperti *excellent, lovely, gorgeous, brilliant, terrible, horrible* dan sebagainya

3. *Manner Violation*

*Manner Violation* sering muncul berupa ekspresi repetisi, terutama yang menggunakan sub maksim '*be brief*'.

## 2.5 TF-IDF

TF-IDF (*Tem Frequency - Inverse Document Frequency* ) adalah ukuran statistik yang digunakan untuk mengetahui seberapa penting suatu kata dalam sebuah dokumen [14]. TF-IDF dapat digunakan untuk pembobotan kata sebagai fitur untuk klasifikasi sentimen kalimat.

TF (*Term Frequency*) adalah bobot dari suatu kata dalam suatu dokumen. Bobot kata dihitung berdasarkan frekuensi kemunculannya dalam dokumen.

Misalnya jika kata bagus muncul 3 kali dalam suatu dokumen dan jumlah kata pada dokumen adalah 9, maka nilai TF dari kata tersebut adalah  $3 / 9 = 0.3$ . TF menggunakan konsep *bag of word* dimana urutan kata tidak diperhatikan [15].

DF (*Document Frequency*) adalah banyaknya dokumen yang mengandung suatu kata. Konsep IDF (*Inverse Document Frequency*) digunakan untuk mengatasi efek dari nilai TF yang terlalu tinggi. Semakin banyak dokumen yang mengandung suatu kata, maka bobot IDF pada kata tersebut semakin kecil [15]. Untuk menghitung nilai IDF, digunakan persamaan 2.1

$$\text{Idf}_t = \log \frac{N}{df_t} \quad (2.1)$$

Untuk menghitung TF-IDF, digunakan persamaan 2.2:

$$\text{Tf.idf}_{t,d} = \text{tf}_{t,d} \times \text{Idf}_t \quad (2.2)$$

Dimana:

$\text{Idf}_t$  = inversi frekuensi dokumen dari kata t

N = banyaknya dokumen

$df_t$  = banyaknya dokumen yang mengandung kata t

$\text{Tf.idf}_{t,d}$  = nilai bobot kata t pada dokumen d

$\text{tf}_{t,d}$  = frekuensi kemunculan term t pada dokumen d

## 2.6 Interjeksi

Kata interjeksi adalah kata yang digunakan untuk mengungkapkan perasaan seseorang. Kata ini digunakan untuk memperkuat perasaan. Perasaan yang diungkapkan dapat berupa kejjikan, kekaguman, keheranan, kekesalan, kesyukuran, kekagetan, ajakan, panggilan, simpulan dan sebagainya. Oleh karena itu kata interjeksi termasuk kata yang bernilai emosi tinggi dan termasuk kata afektif. Kata interjeksi dapat berasal dari bahasa asing atau bahasa daerah. Kata interjeksi dapat di bagi menjadi dua, yaitu kata seru singkat yang mengungkapkan perasaan batin pembicara, seperti wah, wow, lho, dan kata seru yang terdiri dari kata-kata biasa seperti syukur, astaga, dll, yang mengungkapkan perasaan [16].

## 2.7 Sentiwordnet

Sentiwordnet adalah basis data/sumber daya leksikal yang dapat digunakan untuk *opinion mining*. Sentiwordnet berisi istilah-istilah yang diekstrak dari basis data WordNet. WordNet merupakan kamus bahasa Inggris yang dikembangkan oleh Princeton University, yang berisi makna dari suatu kata. Sentiwordnet berisi informasi mengenai nilai positif dan negatif dari suatu kata. Setiap kata pada Sentiwordnet dilengkapi dengan nilai  $Pos(s)$ , dan  $Neg(s)$  yang menunjukkan seberapa positif atau negatif kata tersebut[17].

## 2.8 Naïve Bayes

*Naïve Bayes* adalah metode klasifikasi yang berdasar peluang. Algoritma ini termasuk dalam algoritma klasifikasi statistik yang dapat memprediksi probabilitas keanggotaan kelas. *Naïve Bayes* pertama kali dikemukakan oleh Thomas Bayes. Algoritma ini disebut *naïve* karena kondisi antar atribut diasumsikan saling bebas[18].

Algoritma *Naïve Bayes* menghitung peluang setiap atribut diklasifikasikan ke dalam suatu kelas. Selanjutnya, data akan diklasifikasikan ke kelas yang memiliki peluang tertinggi.

Kelebihan dari algoritma *Naïve Bayes* diantaranya adalah membutuhkan jumlah data latih yang kecil untuk proses pengklasifikasian[19]. Persamaan algoritma *Naïve Bayes* dapat dilihat pada persamaan 2.3

$$P(C|X) = \frac{P(X|C)P(C)}{P(X)} \quad (2.3)$$

Dimana:

$P(C|X)$  = peluang data dengan atribut X diklasifikasikan kelas C.

$P(X|C)$  = peluang atribut X terjadi di dalam kelas C pada data latih.

$P(C)$  = peluang kelas C pada data latih.

$P(X)$  = peluang atribut X pada data latih.

Jika variabel prediktor merupakan variabel kontinu, maka peluang  $P(X|C)$  akan sangat kecil, sehingga hasil akurasi tidak akurat. Oleh karena itu jika

variabel prediktor merupakan variabel kontinu digunakan rumus *denitas gauss* seperti peramaan 2.4

$$P(X_i = x_i | Y = y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}}} e^{-\frac{(x_i - \mu_{ij})^2}{2\sigma^2}} \quad (2.4)$$

Dimana:

P = Peluang atribut xi di klasifikasikan ke kelas yj

X<sub>i</sub> = Data ke i

x<sub>i</sub> = nilai data ke i

Y = nilai kelas yang dicari

y<sub>i</sub> = sub kelas Y yang dicari

μ = rata-rata dari setiap atribut

σ = Standar deviasi dari setiap atribut

Pada proses pengklasifikasian dokumen, *Naïve Bayes* akan memilih kategori yang memiliki probabilitas paling tinggi, seperti pada persamaan 2.5

$$P(C|T = \{X_i, \dots X_n\}) = p(C_i) \text{Argmax} \prod_{i=1}^n P(X_i|C) \quad (2.5)$$

Dimana:

P(X|C) = peluang atribut X terjadi di dalam kelas C pada data latih.

P(C) = peluang kelas C pada data latih

P(X) = peluang atribut X pada data latih

P(C|X) = peluang data dengan atribut X diklasifikasikan kelas C.

## 2.9 Random Forest

*Random Forest* adalah algoritma klasifikasi yang terdiri dari kumpulan *decision tree*, dimana *decision tree* tersebut digunakan untuk mengklasifikasikan data ke suatu kelas. *Random Forest* termasuk metode *ensemble*, yaitu metode untuk meningkatkan akurasi klasifikasi dengan menggabungkan beberapa metode klasifikasi[20]. Algoritma *Random Forest* dapat digunakan untuk klasifikasi dan regresi. *Random Forest* juga dapat digunakan untuk variabel prediktor yang bersifat kategorikal, maupun kontinu[21].

Sesuai dengan namanya, algoritma *Random Forest* membangun hutan (*forest*) dari kumpulan pohon. Setiap pohon dibangun menggunakan data sampel dari *training set* yang diambil menggunakan teknik *bootstrap sampling*. *Test set* kemudian diklasifikasi berdasarkan pohon-pohon yang telah dibangun. Setiap pohon mengklasifikasikan *data test* ke suatu kelas. Data akan di klasifikasikan ke kelas dengan suara terbanyak, atau biasa disebut *majority vote* [20].

## 2.10 CART

CART (*Classification and Regression Trees*) adalah salah satu jenis dari algoritma *decision tree*. CART adalah algoritma prediksi yang dibangun dengan membagi set data secara rekursif dan menyesuaikan model untuk setiap bagian. CART Termasuk metode *nonparametric*. Algoritma CART dapat digunakan untuk klasifikasi dan regresi [22].

*Root node* adalah *node* paling awal dari sebuah *tree*. *Root node* tidak memiliki *input* dan memiliki dua *output*. *Terminal node* atau biasa disebut *leaf* adalah *node* yang paling akhir dan tidak bisa terbagi lagi. *Node* ini memiliki 1 *input* dan tidak memiliki *output*. Setiap *node* yang bukan *node terminal* terpecah menjadi dua turunan, sesuai dengan nilai dari salah satu variabel prediktor. Untuk variabel prediktor kontinu, pemisahan di tentukan oleh *split-point*. Titik yang nilai prediktornya lebih kecil dari *split-point* ke kiri, dan sisanya ke kanan [21].

Ide dasar untuk membentuk pohon didasarkan pada bagaimana memilih satu *split point* pada sebuah *node*, sehingga memperoleh *child node* yang paling *pure*. Pemilihan *split point* mempertimbangkan setiap kemungkinan pemisahan pada setiap variabel prediktor dan memilih yang terbaik sesuai dengan beberapa kriteria [21]. Pada algoritma CART, pemilihan *split* terbaik dipilih berdasarkan indeks *gini* seperti pada persamaan 2.6:

$$gini = \sum_{i=1}^n (1 - \sum_{j=1}^m p_j^2) \frac{t_n}{t} \quad (2.6)$$

Dimana

P = proporsi kelas j pada partisi n

n = banyaknya partisi

m = banyaknya label kelas

t<sub>n</sub> = jumlah data di partisi n

T = jumlah total partisi

## 2.11 Confussion Matrix

*Confussion Matrix* digunakan untuk mengetahui tingkat akurasi hasil klasifikasi. Nilai Akurasi dapat dihitung dengan persamaan 2.7:

$$\frac{TP+TN}{TP+TN+FP+FN} \quad (2.7)$$

Dimana:

TP = *True Positive*, yaitu jumlah kelas positif yang diklasifikasi sebagai kelas positif

TN = *True Negative*, yaitu jumlah kelas negative yang di klasifikasi sebagai kelas negatif

FP = *False Positive*, yaitu jumlah kelas negative yang di klasifikasi sebagai kelas positif

NF = *False Negative*, yaitu jumlah kelas positif yang diklasifikasi sebagai kelas negative

## 2.12 K-Fold Cross Validation

*K-fold cross validation* adalah metode yang digunakan untuk mengevaluasi kinerja algoritma. Pada metode ini, data dipisahkan menjadi dua subset, yaitu data training dan data testing [6]. Data training dan data testing pada disilangkan secara terus menerus sehingga setiap titik memiliki peluang untuk divalidasi [20].

Dalam k-fold cross validation, data di partisi menjadi k bagian yang sama besar. Selanjutnya data di iterasi sehingga setiap partisi data dapat divalidasi sementara k-1 partisi lainnya menjadi data training [20]. Gambar 2.1 mengilustrasikan cara kerja

k-fold cross validation, dengan  $k = 5$ . Bagian yang berwarna gelap menggambarkan data testing, dan bagian yang berwarna putih menggambarkan data training.

Fold				
1	2	3	4	5

Gambar 2. 1 Pembagian Data Train dan Data Test dalam 5-Fold Cross Validation