

ANALISIS SENTIMEN MEDIA SOSIAL TWITTER TERHADAP PRODUK THE BODY SHOP INDONESIA MENGGUNAKAN METODE NAÏVE BAYES

(Analysis of Twitter Social Media Sentiment Towards the Body Shop Indonesia Products Using Naïve Bayes Method)

Nurul Nadiyah Sholihah.^[1], I Gde Putu Wirarama WW^[2], Ariyan Zubaidi^[3]

Program Studi Teknik Informatika, Fakultas Teknik, Universitas Mataram
Jl. Majapahit 62, Mataram, Lombok NTB, INDONESIA

Email: nurulnadya19@gmail.com, [wirarama, zubaidi13]@unram.ac.id

Abstract

The Body Shop is one of the largest cosmetics franchises in the world that has 3000 stores spread across more than 70 countries including Indonesia. The Body Shop has around 144 stores throughout Indonesia with more than one million members. Seeing the large scale of business from The Body Shop Indonesia, it is necessary to improve the quality and innovation in terms of business. This can be done in several ways, one of which is by knowing what opinions are developing in the community about The Body Shop Indonesia. Twitter, for example, is a platform where a person can send, read and share a tweet and interact with many people from different countries or communities. The choice of Twitter social media is due to the rapid flow of information that develops on this platform so that the data obtained will be varied and diverse. From the many opinions given by the public towards The Body Shop, it will take a long time to manually separate several types of comments such as positive, negative and neutral comments. So, because the information that develops on Twitter occurs exponentially, the Twitter API will be used to automate data collection. Then pre-processing will be done on the collected data and then the data will be put into several models for data training purposes. Each extracted tweet will be classified with Naïve Bayes Classifier based on its sentiment whether positive, negative or neutral. The results will show the percentage of each tweet in the positive, neutral and negative classes. After the model was successfully created and implemented, the accuracy result of classification with Naïve Bayes was 62% for grouping into three classes, namely negative, neutral and positive classes. The amount of data used consists of 2564 data with the composition of positive sentiment data as much as 1524 data, 711 neutral data and 329 negative data. Meanwhile, grouping into two classes, namely positive and negative only, resulted in an accuracy of 85%.

Keywords: sentiment analysis, naïve bayes, preprocessing, twitter, twitter api, classification.

1. PENDAHULUAN

Industri kosmetik seperti produk kecantikan merupakan salah satu aspek yang semakin populer di masyarakat, hal ini diketahui dari pernyataan Kementerian Perindustrian RI yang mengatakan industri kosmetik mengalami pertumbuhan sebesar 7% di sepanjang 2019 dan diprediksi mencapai 9% pada 2020. Penjualan produk kecantikan yang semakin meningkat ditambah lagi dengan bermunculannya berbagai tren kecantikan di sosial media menjadikan pasar industri kosmetik semakin luas [1].

Banyaknya pilihan produk kecantikan yang dapat digunakan harus diikuti dengan pengetahuan dan kesadaran konsumen dalam memilih produk yang tepat dan aman baik bagi diri sendiri maupun bagi lingkungan. Di Indonesia, salah satu perusahaan yang giat mempromosikan produk kecantikan aman dan ramah lingkungan adalah The Body Shop Indonesia.

Seiring dengan berkembangnya teknologi termasuk dalam dunia bisnis, pendekatan berbasis data diidentifikasi sebagai cara untuk meningkatkan pertumbuhan bisnis [3]. Hal ini dapat dilakukan dengan analisis sentimen yang merupakan suatu proses untuk menganalisa pendapat orang dari sepotong teks untuk menentukan apakah sentimen itu positif, negatif atau netral [4]. Analisis sentimen dilakukan untuk melihat pendapat atau kecenderungan opini terhadap sebuah masalah, apakah cenderung berpandangan negatif atau positif [5].

Pemilihan media sosial Twitter dikarenakan kembali maraknya penggunaan Twitter di Indonesia. Menurut data yang dirilis oleh Statista pada 29 Juni 2021, jumlah pengguna Twitter di Indonesia menempati posisi ke-6 di dunia dengan jumlah sebanyak 15.1 juta pengguna hingga April 2021. Jumlah ini bertambah dari tahun sebelumnya dimana Indonesia berada di urutan ke-7 dengan jumlah

pengguna Twitter sebanyak 13.2 juta pengguna[6]. Banyaknya pengguna Twitter akan berakibat terhadap banyaknya data cuitan yang akan dikumpulkan.

Terdapat beberapa penelitian sebelumnya yang menggunakan data yang diperoleh dari Twitter sebagai bahan acuan untuk melakukan penelitian. Misalnya, pada penelitian yang dilakukan oleh Connor Gallagher, Kevin Curran dan Eoghan Furey I tahun 2019 [9] mengenai analisis sentimen dan text mining yang dilakukan untuk mengetahui bagaimana pendapat pelanggan terhadap produk smartphone berdasarkan ulasan pelanggan dari situs Amazon yang kemudian hasilnya akan digunakan untuk menyusun strategi dan model bisnis demi meningkatkan pelayanan bagi konsumen.

2. TINJAUAN PUSTAKA

Penelitian tentang analisis sentimen sebelumnya sudah sering dilakukan seperti pada tahun 2017 lalu, Buntoro [10] melakukan penelitian pemilihan Gubernur DKI Jakarta 2017 dengan menggunakan dataset sebanyak 300 cuitan yang didapatkan dari Twitter dengan menggunakan 3 kata kunci yaitu "AHY", "Ahok" dan "Anies" dari masing-masing kata kunci diambil 100 cuitan. Pada proses klasifikasi digunakan metode *Naïve Bayes Classifier* (NBC) dan Support Vector Machine (SVM). Hasil akurasi tertinggi didapatkan dengan menggunakan metode Naïve Bayes yang memiliki akurasi sebesar 95% untuk dataset "AHY" sedangkan akurasi dengan menggunakan SVM menghasilkan nilai akurasi tertinggi sebesar 90% untuk dataset "AHY".

Pada penelitian yang dilakukan oleh Saputra dkk [14] mengenai analisis sentimen data Presiden Jokowi, data yang digunakan diambil menggunakan *search technique* dengan penggunaan data yang berasal dari Twitter, Facebook serta blog politik. *Preprocessing* yang dilakukan pada penelitian ini melalui beberapa tahapan yaitu perubahan kata tidak baku menjadi kata baku, stemming atau pencarian akar kata, tokenisasi dan *stopword removal*.

Penelitian serupa mengenai sentimen analisis juga dilakukan oleh Fiarni dkk [17] untuk *review* toko ritel online Indonesia dengan menggunakan metode Naïve Bayes. Penelitian ini bertujuan untuk membantu industri usaha menengah di Indonesia untuk memahami kebutuhan target pasarnya. Sentimen akan dikelompokkan ke dalam tiga kelas yaitu positif, netral dan negatif. Sistem yang dihasilkan pada penelitian ini mampu mengklasifikasikan pendapat masyarakat dengan nilai presisi sebesar 97.25% dan akurasi

sebesar 89.21%, hal ini menunjukkan bahwa sistem yang diusulkan memiliki keandalan yang baik.

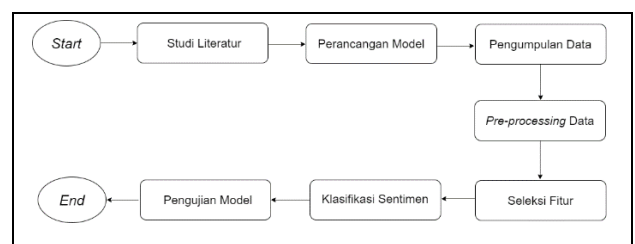
Penelitian mengenai analisis sentimen terhadap transportasi umum MRT Jakarta dengan Naïve Bayes Classifier dilakukan oleh Ikasari dkk [20] pada tahun 2020. Data didapatkan melalui proses pengunduhan data dengan Twitter API dari media sosial Twitter. Sentimen yang telah melewati tahap prapemrosesan lalu diklasifikasikan menggunakan metode Naïve Bayes. penelitian ini menghasilkan tingkat akurasi sebesar 95.88%.

Berdasarkan penelitian-penelitian terkait yang sudah dilakukan sebelumnya, dapat dilihat bahwa Naïve Bayes menghasilkan akurasi yang baik untuk mengklasifikasikan sentimen. Sehingga penulis ingin melakukan analisis sentimen pengguna Twitter terhadap brand The Body Shop Indonesia dengan menggunakan metode Naïve Bayes untuk mengategorikan sentimen-sentimennya ke dalam kelas positif, netral dan negatif. Pertimbangan penulis dalam pemilihan Naïve Bayes sebagai metode untuk melakukan klasifikasi adalah karena prosesnya yang sederhana, cepat dan bisa menghemat banyak waktu, tidak sensitif terhadap fitur yang tidak relevan.

3. METODE PENELITIAN

3.1. Alur Penelitian

Tahapan jalannya penelitian tugas akhir ini diilustrasikan pada diagram alir di Gambar 1.



Gambar 1. Alur Penelitian

Berdasarkan Gambar 1 diuraikan tahapan alur penelitian dimulai dengan studi literatur, lalu melakukan perancangan model, pengumpulan data yang dilakukan secara otomatis dengan bantuan Twitter API, kemudian data yang telah dikumpulkan akan melalui tahap *preprocessing*, selanjutnya data yang sudah bersih kemudian akan menjadi bahan untuk melakukan pembobotan dengan TF-IDF untuk memilih fitur-fitur yang akan digunakan pada tahap klasifikasi. Data yang sudah melalui tahap preprocessing dan pembobotan dengan TF-IDF akan menjadi *input* pada tahap pembangunan model klasifikasi sentimen dengan Naïve Bayes. Dimana

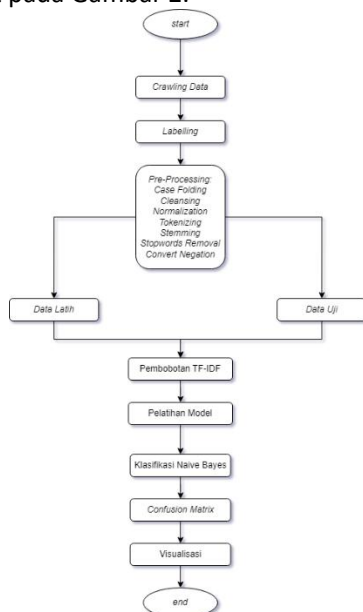
hasilnya akan melalui tahap evaluasi atau pengujian untuk melihat kelayakan model yang dibuat. Dari hasil klasifikasi dan pengujian ini akan dilakukan penarikan kesimpulan agar dapat diketahui informasi yang dibutuhkan yaitu sentimen pengguna Twitter terhadap brand The Body Shop. Hasil akhir yang diharapkan dari penelitian ini adalah pengelompokan polaritas sentimen ke dalam kelas positif, netral dan negatif.

3.1.1. Studi literatur

Studi literatur dalam penelitian ini dilakukan dengan mempelajari referensi-referensi seperti jurnal penelitian dan prosiding seminar yang berkaitan dengan penelitian. Materi yang menjadi fokus untuk dipelajari dalam studi literatur adalah konsep Twitter API, preprocessing data, pembobotan kata dengan TF-IDF dan konsep klasifikasi analisis sentimen dengan metode Naive Bayes.

3.1.2. Perancangan Model

Perancangan model mencakup desain mengenai bagaimana model akan bekerja untuk memproses data dan melakukan klasifikasi sentiment. Rancangan model dapat dilihat pada Gambar 2.



Gambar 2. Rancangan Model

3.1.3. Pengumpulan data

Proses pengunduhan data yang berupa cuitan dilakukan dengan menggunakan *library* Python Tweepy dan pengunduhan cuitan dilakukan berdasarkan kata kunci untuk mendapatkan data yang spesifik sesuai kebutuhan. Cuitan yang diambil merupakan cuitan yang mengandung kata kunci “the body shop” dan “tbs”. Data yang diunduh berupa tweet, nama pengguna dan waktu tweet dibuat.

Cuitan-cuitan yang dikumpulkan lalu akan diseleksi terlebih dahulu kesesuaiannya untuk dipakai sebagai dataset dalam pelatihan model klasifikasi sentimen.

3.1.4. Pelabelan data

Data yang telah dikumpulkan selanjutnya akan dibagi menjadi data latih dan data uji dengan menggunakan fungsi `train_test_split` dari *library* sklearn. Data akan diberi label dengan label positif, label netral dan label negatif dengan mengacu pada daftar kata opini positif dan negatif yang mulanya berasal dari penelitian Liu dkk [44] yang kemudian diterjemahkan dan dimodifikasi oleh Wahid dan Azhari [45] ke dalam Bahasa Indonesia.

3.1.5. Pre-processing data

Data yang didapatkan dari hasil pengunduhan dengan Twitter API masih berupa data mentah sehingga perlu melalui tahapan *preprocessing* terlebih dahulu. Langkah ini dibutuhkan untuk membuat data tidak terstruktur yang didapatkan ke dalam bentuk yang siap diolah. Setelah semua tahap *preprocessing* selesai, maka hasilnya akan disimpan dan dijadikan sebagai dataset untuk melakukan proses klasifikasi analisis sentimen.

3.1.5.1 Case Folding

Pada tahap *case folding* akan dilakukan penyeragaman bentuk huruf dari cuitan yang telah dikumpulkan dengan menjadikan bentuk huruf dalam cuitan ke dalam huruf kecil dengan menggunakan method `.lower()`.

3.1.5.2 Cleansing

Tahap ini akan menghapus karakter-karakter yang tidak penting dan tidak memiliki peran dalam menganalisis sentimen. Data akan dibersihkan dari tanda baca, tautan, symbol, mention dan hastag.

3.1.5.3 Normalisasi

Pada proses normalisasi akan dilakukan penanganan bahasa alay atau gaul dan penghapusan karakter berulang. Tahap ini dilakukan dengan bantuan kamus bahasa gaul dimana dilakukan pemeriksaan apakah suatu kata terdapat dalam kamus bahasa gaul atau tidak, jika iya maka kata akan tersebut akan diganti menjadi kata baku sesuai dengan data pada kamus.

3.1.5.4 Tokenization

Tahap ini akan memecah teks menjadi kata atau token berdasarkan karakter pemisah yang memisahkannya yaitu *whitespace*. Selain itu semua karakter yang bukan huruf seperti angka dan tanda baca serta delimiter lainnya akan dihapus. Pada proses ini digunakan method `word_tokenize()` yang disediakan NLTK.

3.1.5.5 Stemming

Tahap ini akan mengubah kata yang berimbuhan menjadi kata dasar. Dalam melakukan proses stemming digunakan library Sastrawi yang menyediakan stemmer untuk Bahasa Indonesia.

3.1.5.6 Stopwords Removal

Pada tahap ini akan dilakukan penghapusan *stopwords* atau kata-kata umum yang memiliki frekuensi kemunculan tinggi tapi tidak memiliki arti yang bermakna dalam penganalisisan sentimen. Daftar *stopwords* Bahasa Indonesia didapatkan dari NLTK dan daftar *stopwords* yang disusun oleh Tala, serta terdapat beberapa *stopwords* yang ditambahkan secara manual dengan fungsi `extend`. Selain penambahan *stopwords* secara manual, dilakukan juga pengurangan *stopwords* secara manual, karena beberapa kata yang terdapat pada daftar *stopwords* merupakan kata negasi yang hendak digunakan untuk proses konversi negasi pada tahap selanjutnya.

3.1.5.7 Negation Handling

Pada tahap ini akan dilakukan penggabungan kata-kata negasi yang terdiri dari "tidak", "tak", "bukan", "tanpa", "jangan", "belum" dengan kata yang mengikutinya. Proses ini dilakukan dengan menggunakan fungsi `Series.str.replace` dari modul *regular expression*.

3.1.6. Seleksi fitur dengan TF-IDF

Analisis sentimen pada penelitian ini akan menggunakan metode pembobotan *Term Frequency - Inverse Document Frequency* (TF-IDF). Setiap *term* (kata) yang telah diekstrak akan dihitung besar bobotnya terhadap suatu data (dokumen / *tweet*). Metode TF-IDF merupakan kombinasi dari model TF (*Term Frequency*) dan IDF (*Inverse Document Frequency*). TF merupakan frekuensi kemunculan sebuah *term* dalam satu data dimana semakin sering suatu kata muncul maka bobotnya akan semakin besar, sedangkan IDF berperan untuk mengurangi dominasi kata yang sering muncul di berbagai data sehingga akan ditafsirkan sebagai *common term* (*term* yang

umum ditemukan) dan nilainya dianggap tidak memiliki pengaruh [39][40].

Berikut merupakan persamaan yang digunakan untuk mencari bobot suatu *term* (*Term Frequency*):

$$tf = \frac{\text{frekuensi kemunculan term } t \text{ di dokumen } d}{\text{jumlah total terms pada dokumen } d}$$

Perhitungan nilai *Inverse Document Frequency* dilakukan dengan menggunakan persamaan berikut:

$$idf_t = \log_{10} \times \frac{N}{df_t}$$

Dimana :

idf_t	=	nilai IDF <i>term t</i>
N	=	jumlah dokumen
df_t	=	jumlah dokumen yang mengandung <i>term t</i>

Sedangkan untuk menghitung bobot masing-masing dokumen terhadap *term* digunakan persamaan TF- IDF sebagai berikut:

$$W_{dt} = tf_{dt} \times idf_t$$

Dimana :

W_{dt}	=	nilai bobot TF-IDF
tf_{dt}	=	frekuensi kemunculan <i>term t</i> pada dokumen d
idf_t	=	nilai IDF <i>term t</i>

3.1.7. Klasifikasi dengan Naïve Bayes

Naïve Bayes Classifier merupakan *classifier probabilistic* yang mengaplikasikan properti dari teorema Bayes dengan asumsi terdapat independensi yang kuat antara fitur-fiturnya (*naïf*). Kelebihan dari penggunaan pengklasifikasi ini adalah tidak banyak dibutuhkannya data latih untuk memperhitungkan parameter-parameter yang akan digunakan dalam melakukan prediksi. Karena independensi fitur ini, alih-alih menghitung *matrix covariance* secara lengkap, *Naïve Bayes Classifier* hanya akan memperhitungkan varian dari fitur yang dibutuhkan [41]. Misalnya terdapat *tweet "d"* dan himpunan kelas "*c*" (positif dan negatif), maka untuk menentukan kelas dari *tweet "d"* digunakan persamaan sebagai berikut:

$$P(c|d) = \frac{P(d|c) \times P(c)}{P(d)}$$

Dimana :

- d = data dengan kelas yang tidak diketahui
- c = representasi data “d” merupakan kelas yang spesifik
- $P(c|d)$ = probabilitas *tweet* dengan atribut “d” terdapat pada kelas “c” (*posterior probability*)
- $P(c)$ = probabilitas awal kelas “c” (*prior probability*)
- $P(d)$ = probabilitas “d”
- $P(d|c)$ = probabilitas independen kelas “c” dari semua fitur dalam vektor “d” (*Likelihood*)

3.1.8. Pengujian model

Confusion Matrix digunakan untuk mengevaluasi kinerja proses klasifikasi. Matriks ini menunjukkan hubungan antara data yang diklasifikasikan dengan benar dan salah. Pada *confusion matrix*, *True Positive* (TP) mewakili jumlah sentimen yang diklasifikasikan dengan benar, sedangkan *False Positive* (FP) memberikan jumlah banyaknya sentimen negatif yang diklasifikasikan sebagai sentimen positif oleh *classifier*. Sama halnya, *True Negative* (TN) adalah jumlah sentimen negatif yang diklasifikasikan dengan benar dan *False Negative* (FN) adalah jumlah sentimen positif yang diklasifikasikan sebagai sentimen negatif oleh *classifier* [41].

Tabel 1 Tabel *Confusion Matrix*

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

Dari *confusion matrix* ini akan dihitung parameter-parameter berbeda seperti *precision*, *recall*, *F-1 score* dan *accuracy* untuk mengevaluasi kinerja proses klasifikasi [41].

Berikut merupakan rumus dari *precision* yang akan memberikan nilai ketepatan dari *classifier* dengan menghitung perbandingan dari jumlah sentimen positif yang diklasifikasikan dengan benar terhadap jumlah total sentimen yang diklasifikasikan positif [41].

$$Precision = \frac{TP}{TP + FP}$$

Di bawah ini merupakan cara perhitungan untuk mendapatkan nilai *recall* yang menggambarkan keberhasilan sistem dalam menemukan kembali

informasi kelengkapan dari *classifier* dengan menghitung perbandingan antara jumlah sentimen positif yang diklasifikasikan dengan benar terhadap jumlah sebenarnya dari sentimen positif yang ada pada data [41].

$$Recall = \frac{TP}{TP + FN}$$

Rumus berikut adalah perhitungan yang digunakan untuk menghitung nilai *F-1 Score* yang merupakan rasio rata-rata *precision* dan *recall* yang dibobotkan dimana nilainya berkisar antara 0 dan 1 dengan 1 sebagai nilai yang terbaik [42].

$$F - 1 Score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Accuracy dihitung sebagai perbandingan antara jumlah sentimen yang diklasifikasikan dengan benar terhadap jumlah total sentimen yang ada pada data [41]. Rumus yang digunakan untuk menghitung *accuracy* adalah sebagai berikut:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

4. HASIL DAN PEMBAHASAN

4.1 Pengumpulan data

Data yang berhasil dikumpulkan berjumlah 2564 data dengan rincian sebanyak 1524 data merupakan data dengan sentiment positif, 711 data bersentimen netral dan sebanyak 329 data bersentimen negative. Data ini terkumpul dalam kurun waktu 4 bulan, mulai dari tanggal 1 April 2021 sampai dengan 30 Agustus 2021. Data hasil *scraping* dapat dilihat pada Gambar 3.

	Username	Created_at	Tweets
0	lorenzoizza	2021-04-01 20:32	The Body Shop'in- Shea (8/10) : berasa powdery ...
1	narapawaka	2021-04-01 08:55	Dengan produk yang dikenal sebagai produk yang...
2	dindamiko	2021-04-01 07:50	@beautything The body shop charcoal mask bagoo...
3	_peacheje	2021-04-01 02:34	Mau tanya body mist nya the body shop yg enak ...
4	lorenzoizza	2021-04-01 20:35	Klo tbs yg basic kyk strawberry, mango, vanill...
...
2559	charmingyul	2021-08-30 23:24	@ohmybeautybank Black musk tbs!!! Sumpah ini t...
2560	TXT_HVENINGKAI	2021-08-30 15:06	@dixneuf @sugarwaxinc tea tree oilnya tbs bag...
2561	idekadoshopee	2021-08-30 12:20	@bogorfess_ Sepengalamanku tbs tahan lama kaa
2562	quokkasky	2021-08-30 10:03	Pengen beli tbs black musk tapi agak pricey
2563	annidaferani	2021-08-30 09:13	@shdrina Lagi banyak yg diskon nih tbs, bener2...

2564 rows x 3 columns

Gambar 3. Data Hasil *Scraping*.

4.2 Pelabelan Data

Pada pelabelan digunakan angka sebagai label, angka 0 diberikan untuk data dengan sentiment negatif, angka 1 diberikan untuk data dengan sentiment netral dan angka 2 diberikan untuk data dengan sentiment positif. Data yang sudah diberi label dapat dilihat pada Gambar 4.

	Label	Tweets
0	2	The Body Shop\n- Shea (8/10) : berasa powdery ...
1	2	Dengan produk yang dikenal sebagai produk yang...
2	2	@beaughtingy The body shop charcoal mask bagoos...
3	1	Mau tanya body mist nya the body shop yg enak ...
4	1	Klo tbs yg basic kyk strawberry, mango, vanill...
...
2559	2	@ohmybeautybank Black musk tbs!!! Sumpah ini t...
2560	2	@dixnneuf @sugarwaxinc tea tree oilnya tbs bag...
2561	2	@bogorfess_ Sepengalamanku tbs tahan lama kaa
2562	0	Pengen beli tbs black musk tapi agak pricey
2563	1	@shdrina Lagi banyak yg diskon nih tbs, bener2...

2564 rows × 2 columns

Gambar 4. Data dengan Label

4.3 Text Pre-Processing

Dari hasil scraping didapatkan data sejumlah 2564 data yang kemudian disimpan dalam satu file dengan format csv. Selanjutnya data akan melalui tahapan pra-pemrosesan yang terdiri dari beberapa tahapan seperti berikut:

4.3.1 Case Folding

Pada tahap ini akan dilakukan konversi data dari huruf kapital menjadi huruf kecil dengan menggunakan method `.lower()`. Hasil case folding dapat dilihat pada Gambar 5

```
[1] the body shop\n- shea (8/10) : berasa powdery klo agak lama dipake, ada hint vanillanya, agak nutty, elegan, soft, bias
nya aku pake klo abis mandi, harganya under 200rb\n- black musk (8,5/10) : musk vanilla gitu, lumayan susah dijelasin t
p ini cocok bgt dipake klo malem hari',
'dengan produk yang dikenal sebagai produk yang natural & mendukung global warming, itu saya tertarik untuk mencoba
produk the body shop & langsung merasa cocok sejak pandemi ini. setelah saya research, ternyata brand the body shop a
dih produk yg paling sering dicari sejak q3 2020 https://t.co/T5skunttdt',
'@beaughtingy the body shop charcoal mask bagoosooooo bgtb' 'vaid',
'sau tanya body mist nya the body shop yg enak rasa apa',
'klo tbs yg basic kyk strawberry, mango, vanilla, gitu2 gausah dijelasin laya krn baunya sesuai namanyaaa',
'@beaughtingy tergantung produknya ga sil aku pake tbs tea tree engga bermiyak kl kebanyakan',
'grocnya itu kalo bentuknya masih jerawat gt coba pake tea tree oil nya tbs kak. kalo sy sejauh ini cocok sm itu :)',
'@ohmybeautybank tbs mah paling enak yg spritz sweet love',
'berkali kali ikutin rekomen mb mb sociolla lah tbs lah bbw lah ga ada yg bener ngasi rekomennyaab' 'vaid perkara buruz be
lanja tp di saktuinnya aye',
'@ohmybeautybank aku gatau harumnya sih neder tapi tbs ada juga di shopee kok',
'@ohmybeautybank tbs mango lebih seger, vanilla lebih sweet kalo kata aku',
'@ber_uangbanyak gokil si body butter w pernah pake herboris yang coklat wanginya enak banget, sehari-an. apalagi tbs, du
rea enak .dah keluar sabuk direm.'
```

Gambar 5 Hasil Case Folding

4.3.2 Cleansing

Pada tahap ini akan dilakukan pembersihan data dari noise seperti tanda baca, mention, hastag, dsb. Hasil case folding dapat dilihat pada Gambar 6

```
[1] the body shop\n- shea (8/10) : berasa powdery klo agak lama dipake, ada hint vanillanya, agak nutty, elegan, soft, bias
nya aku pake klo abis mandi, harganya under 200rb\n- black musk (8,5/10) : musk vanilla gitu, lumayan susah dijelasin t
p ini cocok bgt dipake klo malem hari',
'dengan produk yang dikenal sebagai produk yang natural & mendukung global warming, itu saya tertarik untuk mencoba
produk the body shop & langsung merasa cocok sejak pandemi ini. setelah saya research, ternyata brand the body shop a
dih produk yg paling sering dicari sejak q3 2020 https://t.co/T5skunttdt',
'@beaughtingy the body shop charcoal mask bagoosooooo bgtb' 'vaid',
'sau tanya body mist nya the body shop yg enak rasa apa',
'klo tbs yg basic kyk strawberry, mango, vanilla, gitu2 gausah dijelasin laya krn baunya sesuai namanyaaa',
'@beaughtingy tergantung produknya ga sil aku pake tbs tea tree engga bermiyak kl kebanyakan',
'grocnya itu kalo bentuknya masih jerawat gt coba pake tea tree oil nya tbs kak. kalo sy sejauh ini cocok sm itu :)',
'@ohmybeautybank tbs mah paling enak yg spritz sweet love',
'berkali kali ikutin rekomen mb mb sociolla lah tbs lah bbw lah ga ada yg bener ngasi rekomennyaab' 'vaid perkara buruz be
lanja tp di saktuinnya aye',
'@ohmybeautybank aku gatau harumnya sih neder tapi tbs ada juga di shopee kok',
'@ohmybeautybank tbs mango lebih seger, vanilla lebih sweet kalo kata aku',
'@ber_uangbanyak gokil si body butter w pernah pake herboris yang coklat wanginya enak banget, sehari-an. apalagi tbs, du
rea enak .dah keluar sabuk direm.'
```

Gambar 6 Hasil cleansing

4.3.3 Normalisasi

Pada proses normalisasi dilakukan proses penggantian kata dari kata tidak baku menjadi kata

baku dengan mengacu pada kamus singkatan. Hasil data yang sudah di normalisasi dapat dilihat pada Gambar 7

	Label	Tweets	Normalized
0	2	the body shop shea berasa powdery klo agak lam...	the body shop shea berasa powdery kalau agak l...
1	2	dengan produk yang dikenal sebagai produk yang...	dengan produk yang dikenal sebagai produk yang...
2	2	the body shop charcoal mask bagoos bgt	the body shop charcoal mask bagoos sangat...
3	1	mau tanya body mist nya the body shop yg enak ...	mau tanya body mist nya the body shop yang ena...
4	1	klo tbs yg basic kyk strawberry mango vanilla ...	kalau tbs yang basic kayak strawberry mango va...
...
2559	2	black musk tbs sumpah ini tuh makin lama dipa...	black musk tbs sumpah ini tuh makin lama dipak...
2560	2	tea tree oilnya tbs bagus semalem jerewi kempes	tea tree oilnya tbs bagus semalam jerewi kempes
2561	2	sepengalamanku tbs tahan lama kaa	sepengalamanku tbs tahan lama kaa
2562	0	pengen beli tbs black musk tapi agak pricey	ingin beli tbs black musk tapi agak pricey
2563	1	lagi banyak yg diskon nih tbs bener menggoda...	lagi banyak yang diskon nih tbs benar menggoda...

2564 rows × 3 columns

Gambar 7 Hasil Normalisasi

4.3.4 Tokenization

Tahap ini akan memecah teks menjadi kata atau token berdasarkan karakter pemisah yang memisahkannya yaitu *whitespace*. Selain itu semua karakter yang bukan huruf seperti angka dan tanda baca serta *delimiter* lainnya akan dihapus. Pada proses ini digunakan *method* `word_tokenize()` yang disediakan NLTK. Hasil *tokenization* dapat dilihat pada Gambar 8

```
0 [the, body, shop, shea, berasa, powdery, kalau...
1 [dengan, produk, yang, dikenal, sebagai, produ...
2 [the, body, shop, charcoal, mask, bagoos, sangat]
3 [mau, tanya, body, mist, nya, the, body, shop,...
4 [kalau, tbs, yang, basic, kayak, strawberry, m...
...
2559 [black, musk, tbs, sumpah, ini, tuh, makin, la...
2560 [tea, tree, oilnya, tbs, bagus, semalam, jerew...
2561 [sepengalamanku, tbs, tahan, lama, kaa]
2562 [ingin, beli, tbs, black, musk, tapi, agak, pr...
2563 [lagi, banyak, yang, diskon, nih, tbs, benar, ...
Name: Tokenized, Length: 2564, dtype: object
```

Gambar 8 Hasil Tokenization

4.3.5 Stemming

Tahap ini akan mengubah kata yang berimbuhan menjadi kata dasar. Dalam melakukan proses *stemming* digunakan *library* Sastrawi yang menyediakan *stemmer* untuk Bahasa Indonesia. Hasil *stemming* dapat dilihat pada Gambar 9.

```
0 the body shop shea asa powdery pakai hint vani...
1 produk kenal produk natural dukung global warm...
2 the body shop charcoal mask bagoos
3 body mist the body shop enak
4 tbs basic kayak strawberry mango vanilla dijel...
...
2559 black musk tbs sumpah pakai enak wangi tahan p...
2560 tea tree oilnya tbs bagus malam jerewi kempes
2561 alam tbs tahan kaa
2562 beli tbs black musk pricey
2563 diskon tbs goda iman dompet
Name: Ulasan_clean, Length: 2564, dtype: object
```

Gambar 9 Hasil Stemming

4.3.6 Stopwords Removal

Pada tahap ini akan dilakukan penghapusan stopwords atau kata-kata umum yang memiliki frekuensi kemunculan tinggi tapi tidak memiliki arti yang bermakna dalam penganalisisan sentimen. Daftar stopwords Bahasa Indonesia didapatkan dari NLTK dan daftar stopwords yang disusun oleh Tala, serta terdapat beberapa stopwords yang ditambahkan secara manual dengan fungsi `extend`. Selain penambahan stopwords secara manual, dilakukan juga pengurangan stopwords secara manual, karena beberapa kata yang terdapat pada daftar stopwords merupakan kata negasi yang hendak digunakan untuk proses konversi negasi pada tahap selanjutnya. Data hasil penghapusan stopwords dapat dilihat pada Gambar 10.

```
0 [the, body, shop, shea, berasa, powdery, dipak...
1 [produk, dikenal, produk, natural, mendukung, ...
2 [the, body, shop, charcoal, mask, bagoos]
3 [body, mist, the, body, shop, enak]
4 [tbs, basic, kayak, strawberry, mango, vanilla...
...
2559 [black, musk, tbs, sumpah, dipakai, enak, wang...
2560 [tea, tree, oilnya, tbs, bagus, semalam, jerew...
2561 [sepengalamanku, tbs, tahan, kaa]
2562 [beli, tbs, black, musk, pricey]
2563 [diskon, tbs, menggoda, iman, dompet]
Name: no_stopwords, Length: 2564, dtype: object
```

Gambar 10 Hasil *Stopwords Removal*

4.3.7 Negation handling

Pada tahap ini akan dilakukan penggabungan kata-kata negasi yang terdiri dari "tidak", "tak", "bukan", "tanpa", "jangan", "belum" dengan kata yang mengikutinya. Proses ini dilakukan dengan menggunakan fungsi `Series.str.replace` dari modul *regular expression*. Data hasil *negation handling* dapat dilihat pada Gambar 11.

	Label	Negation Handling
0	2	the body shop shea asa powdery pakai hint vani...
1	2	produk kenal produk natural dukung global warm...
2	2	the body shop charcoal mask bagoos
3	1	body mist the body shop enak
4	1	tbs basic kayak strawberry mango vanilla tidak...
...
2559	2	black musk tbs sumpah pakai enak wangi tahan p...
2560	2	tea tree oilnya tbs bagus malam jerewi kempes
2561	2	alam tbs tahan kaa
2562	0	beli tbs black musk pricey
2563	1	diskon tbs goda iman dompet

2564 rows × 2 columns

Gambar 11 Hasil *Negation Handling*

4.4 TF-IDF

Pada tahap ini dilakukan pembobotan *terms* dengan metode TF-IDF untuk menyeleksi fitur yang akan digunakan pada tahap klasifikasi. Hasil pembobotan *terms* dengan metode TF-IDF dapat dilihat pada gambar 12 berikut.

Gambar 12 Hasil Pembobotan TF-IDF

4.5 Data training dan testing

Proses pemilahan data uji dan data latih dilakukan dengan proses split dengan bantuan dari *library* `sklearn` yang memiliki fungsi `model_selection.train_test_split`. Dataset dibagi menjadi 30% data uji dan 70% data latih. Jumlah data yang digunakan adalah sebanyak 2564 data dengan komposisi 770 data uji dan 1794 data latih. Ditetapkan *random state* dengan nilai 42, fungsi dari *random state* ini adalah untuk mengontrol pengacakan data uji dan data latih.

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    tfidf, data_label['Label'], test_size=0.3, random_state=42)
```

Gambar 13 Pembagian Data

4.6 Klasifikasi naïve bayes

Pada tahap ini akan dilakukan proses klasifikasi dengan menggunakan algoritma Naïve Bayes. Hasil yang diperoleh dari pengklasifikasian dengan Naïve Bayes adalah sebagai berikut.

```
from sklearn.naive_bayes import MultinomialNB
from sklearn import metrics
clf = MultinomialNB().fit(X_train, y_train)
predicted= clf.predict(X_test)
print("Accuracy:",metrics.accuracy_score(y_test, predicted)*100,"%")
```

Accuracy: 62.077922077922075 %

Gambar 14 Akurasi klasifikasi dengan Naïve Bayes

Berdasarkan hasil klasifikasi yang dilakukan dengan metode Naïve Bayes dengan pembobotan *tfidf*, didapatkan nilai akurasi sebesar 62%. Nilai akurasi yang didapatkan termasuk rendah, hal ini dikarenakan terdapat ambiguitas pada data yang bersentimen netral.

Dikarenakan banyaknya ambiguitas pada data dengan sentiment netral maka akan dicoba untuk melakukan klasifikasi naïve bayes hanya dengan 2 kelas yaitu kelas positif dan kelas negative sebagai

perbandingan. Hasil yang diperoleh dari pengklasifikasian dengan Naïve Bayes dengan 2 kelas adalah sebagai berikut

```
from sklearn.naive_bayes import MultinomialNB
from sklearn import metrics
clf = MultinomialNB().fit(X_train, y_train)
predicted= clf.predict(X_test)
print("Accuracy:",metrics.accuracy_score(y_test, predicted)*100,"%")
Accuracy: 85.02994011976048 %
```

Gambar 15 Akurasi Klasifikasi Naïve Bayes dengan Dua Kelas

Berdasarkan hasil klasifikasi yang dilakukan dengan metode Naïve Bayes dengan pembobotan fitur TF-IDF dengan hanya menggunakan kelas positif dan negatif, didapatkan nilai akurasi sebesar 85%. Nilai akurasi yang didapatkan meningkat jauh dibandingkan dengan akurasi menggunakan 3 kelas.

4.7 Evaluasi

Evaluasi dilakukan dengan menggunakan mencetak *Classification Report* seperti yang terlihat pada Tabel 2.

Tabel 2. *Classification Report Naïve Bayes dengan Tiga Kelas*

	precision	recall	F1-score	support
Negative	1.00	0.01	0.02	96
Netral	0.50	0.06	0.11	208
Positive	0.62	1.00	0.77	466
Accuracy			0.62	770
Macro avg	0.71	0.36	0.30	770
Weighted avg	0.64	0.62	0.50	770

Pada Tabel 2 dapat dilihat bahwa didapatkan akurasi sebesar 0.62 atau 62% dalam persen. Akurasi ini menggambarkan persentase dari total sentiment yang berhasil diklasifikasikan secara benar, angka 62% menunjukkan bahwa dari 770 data uji yang digunakan, sebanyak 477.4 data berhasil diklasifikasikan ke dalam kelas yang sebenarnya. Dari 770 data uji, komposisinya adalah sebanyak 96 data berlabel negative, sebanyak 208 data berlabel netral dan 466 data berlabel positif.

Evaluasi untuk klasifikasi dengan menggunakan 2 kelas dapat dilihat pada Tabel 3.

Tabel 3. *Classification Report Naïve Bayes dengan Dua Kelas*

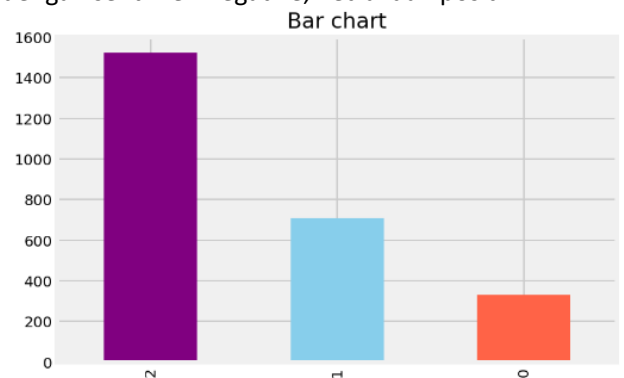
	precision	recall	F1-score	support
Negative	0.00	0.00	0.00	25
Positive	0.85	1.00	0.92	142
Accuracy			0.85	167
Macro avg	0.43	0.50	0.46	167

Weighted avg	0.72	0.85	0.78	167
--------------	------	------	------	-----

Pada Tabel 3 dapat dilihat bahwa didapatkan akurasi sebesar 0.85 atau 85% dalam persen. Akurasi ini menggambarkan persentase dari total sentiment yang berhasil diklasifikasikan secara benar, angka 85% menunjukkan bahwa dari 167 data uji yang digunakan, sebanyak 141.95 data berhasil diklasifikasikan ke dalam kelas yang sebenarnya. Dari 167 data uji, komposisinya adalah sebanyak 25 data memiliki sentimen negative dan sebanyak 142 data memiliki sentimen positif.

4.8 Visualisasi

Dibawah ini merupakan gambar persebaran data dengan sentimen negative, netral dan positif:



Gambar 16 Persebaran Data

Berdasarkan Gambar 16 didapatkan informasi mengenai jumlah data pada masing-masing kelas sentiment. Ditemukan bahwa data dengan sentiment positif berjumlah paling banyak yaitu sebanyak 1524 data, diikuti dengan data bersentimen netral sebanyak 711 data sedangkan data dengan sentiment negative berjumlah paling sedikit yaitu sebanyak 329 data. Dapat diambil kesimpulan bahwa pelanggan yang menggunakan produk the body shop merasa puas dengan pengalaman yang didapatkan.

Berikut merupakan gambar *wordcloud* dari data cuitan.

Observer : Analyzing and Comparing Opinions on the Web," pp. 342–351.

- [12] A. Tripathy, A. Agrawal, and S. K. Rath, "Classification of Sentimental Reviews Using Machine Learning Techniques," *Procedia Comput. Sci.*, vol. 57, pp. 821–829, 2015, doi:

10.1016/j.procs.2015.07.523.

- [13] E. Haddi, X. Liu, and Y. Shi, "The role of text pre-processing in sentiment analysis," *Procedia Comput. Sci.*, vol. 17, pp. 26–32, 2013, doi: 10.1016/j.procs.2013.05.005.