



Membangun Model Machine Learning Untuk Meninjau Layanan Indosat Ooredoo Dari Twitter Menggunakan Naive Bayes Classifier

Febri Astiko¹, Achmad Khodar²

^{1,2}Program Studi Teknik Informatika, Fakultas Ilmu Komputer, Universitas Mercu Buana Jakarta
¹41515120130@student.mercubuana.ac.id, ²achmad.kodar@mercubuana.ac.id

Abstract

This study aims to design a machine learning model of sentiment analysis on Indosat Ooredoo service reviews on social media twitter using the Naive Bayes algorithm as a classifier of positive and negative labels. This sentiment analysis uses machine learning to get patterns an model that can be used again to predict new data.

Keywords: sentiment analysis, indosat ooredoo, twitter, naive bayes classifier, positive negative.

Abstrak

Penelitian ini bertujuan untuk membuat rancangan model machine learning tentang sentimen analisis pada peninjauan layanan Indosat Ooredoo dalam media sosial twitter menggunakan algoritma Naive Bayes sebagai classifier dari label positif dan negatif. Sentimen analisis ini menggunakan machine learning agar mendapatkan pola dan model yang dapat digunakan lagi untuk memprediksi data baru.

Kata kunci: sentimen analisis, indosat ooredoo, twitter, naive bayes classifier, positif negatif.

1. Pendahuluan

Dalam era digital sekarang ini semua orang membutuhkan smartphone dan koneksi internet. Tersedia beberapa jasa provider internet yang bisa digunakan, dan kali ini kami akan mengangkat topik mengenai kualitas layanan internet dari provider Indosat yang dirasa kurang memuaskan bagi para pelanggannya. Contohnya saat pulsa terpotong padahal sudah membeli paket data internet secara terpisah, dan kecepatan internet yang low dan tidak stabil. Pada penelitian ini akan menggunakan tanggapan dari para pelanggan Indosat yang pada media sosial twitter untuk mendapatkan review apakah itu positif atau negatif. Karena itu penelitian ini dibuat agar memberikan umpan balik kepada pihak provider Indosat agar dapat meningkatkan kualitas layanan internet mereka agar saling menguntungkan diantara pelanggan dan provider Indosat sendiri. Hal ini hampir sama kondisinya dengan topik jurnal yang di tulis oleh Irvan Wahyudi, Warih Maharani dan Tjokorda Agung Budi Wirayuda dengan judul “Analisis Sentimen Terhadap Poduk di Twitter Menggunakan Metode Support Vektor Machine” [1]. Kebutuhan akan analisis sentimen melalui twitter menarik untuk diteliti agar memberikan sentimen yang lebih efektif dari masyarakat umum, terutama dalam kasus layanan internet dimana orang senang dan tidak senang dengan mudah dan sering membawa pandangan

mereka ke twitter. Jumlah data yang ada di twitter sangat banyak dan mudah ditemukan, karena itu akan menjadi topik pada penelitian ini. Tweet yang dikirim oleh publik massa dapat diperlakukan sebagai opini dan pendapat mereka perspektif tentang situasi menjadikannya sumber informasi berharga bagi organisasi mana pun yang berusaha menarik lebih banyak orang dengan meningkatkan dan menghadiri kebutuhan pelanggan. Sentimen analisis pada twitter dapat digunakan sebagai instrumen yang bagus untuk ulasan pelanggan dan umpan balik terutama untuk yang baru produk dirilis di pasar [2]. Sebagian besar metode klasifikasi sentimen media sosial saat ini menilai polaritas sentimen terutama berdasarkan konten tekstual, dan model yang digunakan pada penelitian ini adalah Naive Bayes Classifier yang menunjukkan bahwa melampaui konten pada tweet adalah bermanfaat dalam klasifikasi sentimen yang menghasilkan opini positif dan negatif dari para pelanggan yang menggunakan layanan internet dari provider Indosat Ooredoo, karena dapat memberikan hasil classifier dengan pemahaman yang mendalam tentang kualitas dari produk yang telah digunakan dari provider tersebut [3].

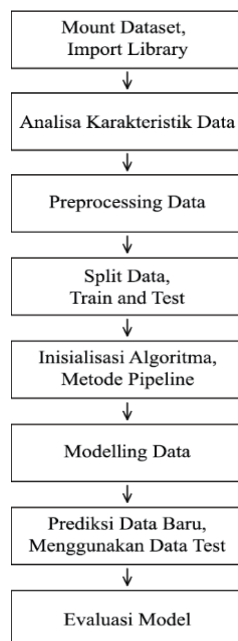
Analisis sentimen media sosial dapat digunakan untuk meningkatkan layanan pelanggan dan pemasaran dan juga berfungsi sebagai ukuran kinerja media sosial. Data analisis sentimen berasal dari postingan twitter yang



memunculkan kata produk yang dimaksudkan kemudian mengekstraksi tweet tersebut. Proses pengolahan data mentah hingga didapatkan kumpulan informasi yang diinginkan melalui tahapan-tahapan antara lain seleksi data, representasi data, algoritma pembelajaran, dan beberapa kombinasi masukan yang dapat meningkatkan performans, hasilnya dapat mengklasifikasikan mereka ke berbagai polaritas yaitu positif dan negatif. Digunakan sebagai pertimbangan produk di masyarakat apakah diterima atau tidak untuk menentukan strategi untuk meningkatkan kualitas dari produk tersebut [4].

2. Metode Penelitian

Metode penelitian yang di rancang adalah sebagai berikut :



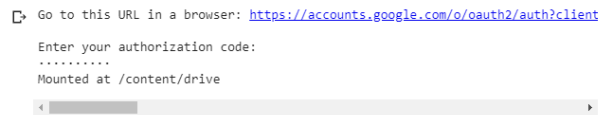
Gambar 1. Tahapan penelitian pembuatan model

a. Mount Dataset dan Import Library yang dibutuhkan

Data latih yang digunakan dalam penelitian ini diambil dari twitter yang berdasarkan hasil pencarian tentang pendapat para pengguna twitter tentang pengalaman menggunakan layanan provider Indosat Ooredoo. Data latih ini bertujuan untuk membuat model machine learning agar mendapatkan pola yang dapat digunakan lagi untuk memprediksi data baru. Data yang di ambil dari twitter untuk data latih adalah lebih dari 200 data posting dan komentar tweet. Dimana data yang di ambil adalah data tweet yang mengandung sentimen terhadap umpan balik dari para konsumen. Data latih akan dikategorikan secara manual yang dilakukan oleh user dan memilih sentimen yang terkandung di dalam tweet tersebut dan menandai tweet tersebut menjadi 2 kategori sentimen yaitu tweet yang mengandung sentimen positif dan negatif [5]. Data latih yang sudah diperoleh ini akan disimpan ke dalam model machine learning yang nantinya akan digunakan sebagai data

training untuk data yang baru. Dataset disimpan dalam google drive, kemudian di akses dengan login gmail untuk mendapatkan key dari google. Dan untuk melakukan mount dataset, langkah pertama adalah

Mount data From Google Drive



mengimport library yang dibutuhkan, yaitu pandas, numpy, regex, NLTK dan matplotlib.

Gambar 2. Login dengan gmail untuk mendapatkan key dari google

b. Analisa Karakteristik Data

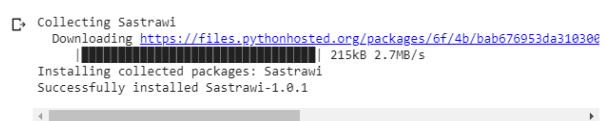
Tahapan ini kita harus mengetahui karakteristik data sebelum diproses, apa saja proses cleaning yang dibutuhkan, berapa baris dan kolom, seperti apa contoh dataset dan labelnya, berapa jenis label, dan rasio setiap label [6]. Caranya adalah dengan load data dan menampilkan 10 data teratas, menampilkan banyak baris dan kolom, melihat jumlah data pada setiap label, dan menampilkan data pie gram dari kolom sentimen.

	KASUS	OUTPUT
0	KUOTA MASIH ADA TAPI PULSA MALAH DIAMBIL	Negatif
1	Kadang juga kuota nya gak bisa di pakai	Negatif
2		Lemot Negatif
3		Pulsa terpotong Negatif
4		Pulsa hilang Negatif
5		Indosat payah Negatif
6		Indosat jelek Negatif
7		Berhenti langganan Negatif
8		IM3 Gangguan Negatif
9		Murah tapi lambat Negatif

Gambar 3. Load data untuk menampilkan 10 data teratas

c. Preprocessing Data

Pada tahap Preprocessing data kita akan menyiapkan komponen apa saja yang akan kita hilangkan untuk mendukung performa akurasi yang baik, dan menyiapkan function untuk dipanggil pada tahap modeling. Tahapan ini dimulai dari cleaning text, dengan pembuatan function clean data yang berfungsi membersihkan data. Kita dapat menggunakan stopwords dari sastrawi, yaitu stopwords bahasa Indonesia untuk menghilangkan stopwords dari data yang kita miliki [7]. Sastrawi adalah library machine learning dalam bahasa pemrograman python. Setelah itu kita menggunakan wordPunctTokenizer untuk menghilangkan data berupa karakter seperti "!,?,@,\$." agar data yang masuk menjadi lebih mudah untuk di proses.



Gambar 4. Install library Sastrawi

Program Jurnal

```

import string
from Sastrawi.StopwordRemover.
    StopwordRemoverFactory
import StopwordRemoverFactory
from StopwordRemoverFactory from tqdm import
    tqdm_notebook as tqdm

factory = StopwordRemoverFactory()
stopword = factory.create_stop_word_remover()

def clean_text(list_of_text)
    output_text = []
    for text in tqdm(list_of_text):
        text = text.translate(str.maketrans(
            ',',',',string.punctuation)).lower()
        text = re.sub(r'[a-zA-Z0-9]', '', str(text))
        text = stopword.remove(str(text))
        output_text.append(str(text))
    return output_text

```

Function `clean_text` data terdiri dari:

`Punct` untuk menghilangkan punctuation seperti simbol, titik, koma, tanda tanya dll dan menjadikan semua huruf menjadi kecil.

`text` memastikan object, only ASCII A-Za-z0-9.

`text` untuk menghilangkan stopwords berbahasa indonesia menggunakan `sastrawi`.

Hasil `text` yang sudah dibersihkan akan dimasukkan kedalam variable `output_text`.

`tqdm` adalah library yang dapat menampilkan progress bar dan waktu proses.

d. Split Data lalu Create Train and Test

Tahapan ini untuk memisahkan data train dan data test menggunakan `scikitlearn train_test_split` untuk inialisasi variable, lalu memisahkan `x_train`, `y_train`, `x_test`, `y_test` yang digunakan untuk membuat model [8]. Sedangkan `train_test_split` digunakan untuk memisahkan variable (`x,y` dengan pembagi data 80%(train:20%(test) sehingga nilai `test_size` adalah 0.2)). Pada tahap ini kita menggunakan library python `scikitlearn train_test_split` dan menginstall `AutoML tables client library`.

Program Jurnal

```

#Inialisasi Variable x dan y
from sklearn.model_selection
import train_test_split
x = data.KASUS
y = data.OUTPUT
x.head(10)

0    KUOTA MASIH ADA TAPI PULSA MALAH DIAMBIL
1    Kadang juga kuota nya gak bisa di pakai
2                                     Lemot
3                                     Pulsa terpotong
4                                     Pulsa hilang
5                                     Indosat payah
6                                     Indosat jelek
7    Berhenti langganan
8    IM3 Gangguan
9    Murah tapi lambat
Name: KASUS, dtype: object

```

Gambar 5. Install AutoML table client library

e. Inialisasi Algoritma dan Metode Pipeline

Algoritma yang digunakan dalam penelitian ini adalah Naive Bayes Classifier (BernoulliNB) yang memiliki asumsi sangat kuat akan independensi dari masing-masing kondisi atau suatu kejadian [9].

$$P(C|A) = \frac{P(A|C)P(C)}{P(A)} \quad (1)$$

Keterangan :

A : Data dengan kelas yang belum diketahui

C : Representasi data C merupakan kelas yang spesifik

$P(C|A)$: Probabilitas C berdasarkan kondisi A

$P(C)$: Probabilitas

$P(A|C)$: Probabilitas A terhadap kondisi C

$P(A)$: Probabilitas A

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Gambar 6. Prediksi Value Algoritma Naive Bayes Classifier

Program Jurnal

```

from sklearn.feature_extraction.text
import TfidfVectorizer
from sklearn.naive_bayes
import BernoulliNB
from sklearn.pipeline

[ ] hasil_bnb = model_bnb.predict(x_test)

[ ] #Test Kata
model_bnb.predict(["Harga Paketnya murah, tapi kecepatan internetnya lambat"])

array(['Negatif'], dtype='<U7')

import Pipeline

tfvec = TfidfVectorizer()
clf_bnb = BernoulliNB()
model_bnb = Pipeline([('tfidf', tfvec),
                       ('classifier', clf_bnb)])

```

Pada tahap ini `BernoulliNB` digunakan untuk membuat pipeline, dan di extract dengan feature dari library `Sklearn` dan `TfidfVectorizer`.

f. Modelling Data

Tahapan modelling data adalah untuk membuat suatu model machine learning yang akan menjadi fungsi utama untuk melatih data-data baru yang akan digunakan kedepannya [10]. Proses ini adalah fungsi utama supaya program yang dibuat bisa berjalan untuk membuat model dapat berjalan dengan lancar. Di dalam modelling data diharapkan membuat program yang simple tapi efektif dan juga mudah dikembangkan agar ke depannya dapat di modifikasi sesuai dengan

kebutuhan atau untuk melakukan komparasi dengan metode algoritma machine learning yang lain. Untuk listing program algoritma Naive Bayes dalam bahasa pemrograman python untuk membuat train data X dan Y menggunakan BernoulliNB dalam model ini, contoh script modelling datanya sebagai berikut:

Program Jurnal

```
Model_bnb.fit(x_train, y_train)

Pipeline(memory=None, steps=[('tfidf',
    TfidfVectorizer(analyzer='word',
    binary=False, decode_error=
    'strict', dtype=<class 'numpy.float64'>,
    encoding='utf-8', input='content',
    lowercase=True, max_df=1.0,
    max_features=None, min_df=1,
    ngram_range=(1, 1), norm='l2',
    preprocessor=None, smooth_idf=True,
    stop_words=None, strip_accents=None,
    sublinear_tf=False, token_pattern=
    '(?u)\b\w+\b', tokenizer=None,
    use_idf=True, vocabulary=None))
    ('classifier',
    BernoulliNB(alpha=1.0, binarize=0.0,
    class_prior=None, fit_prior=True))],
    verbose=False)
```

g. Prediksi Data baru menggunakan Data Test

Tahap ini adalah dengan menggunakan data baru yang akan di diprediksi dengan data test yang telah dibuat, yaitu dengan model machine learning yang sudah berhasil memprediksi data yang pertama di training. Data baru yang bisa digunakan jumlahnya tidak terbatas, karena model yang dibuat selalu fleksibel dan semakin pintar dengan semakin bertambahnya bermacam data baru yang di proses.

```
[ ] #Accuracy Score
print("""
    Hasil Skor Akurasi Menggunakan Naive Bayes
    """)
print(accuracy_score(hasil_bnb,y_test))
```

```
Hasil Skor Akurasi Menggunakan Naive Bayes

0.7317073170731707
```

Gambar 7. Prediksi dan test dengan contoh kalimat komentar

h. Evaluasi Model

Ini adalah tahapan terakhir dalam metode penelitian ini, yang bertujuan untuk mengevaluasi hasil akhir dari model yang telah dibuat. Model ini akan menghasilkan output akurasi menggunakan Naive Bayes Classifier, dengan melalui proses confusion matrix, classification report, accuracy score. Untuk contoh hasil dari eksekusi programnya adalah sebagai berikut:

```
[ ] #Classification Report
print("""
    Hasil Laporan Classification Menggunakan Naive Bayes
    """)
print(classification_report(y_test,hasil_bnb))
```

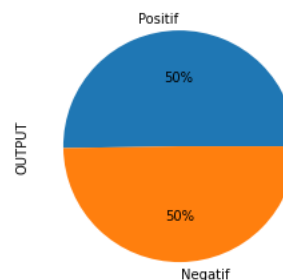
```
Hasil Laporan Classification Menggunakan Naive Bayes
```

	precision	recall	f1-score	support
Negatif	0.68	0.90	0.78	21
Positif	0.85	0.55	0.67	20
accuracy			0.73	41
macro avg	0.76	0.73	0.72	41
weighted avg	0.76	0.73	0.72	41

Gambar 8. Hasil eksekusi Confusion Matrix

3. Hasil dan Pembahasan

Pengujian model machine learning ini menggunakan lebih dari 200 data latih, yang terdiri dari 100 kelas positif dan 100 kelas negatif. Fitur yang digunakan meliputi tweet yang mengandung hashtag, retweet, comment. Fitur-fitur yang terdiri dari jumlah kata benda, sifat, kerja, keterangan, seru dan persentase kata benda, sifat, kerja, keterangan, seru belum dapat menentukan suatu data dapat masuk ke kelas positif atau negatif dengan tepat. Banyaknya jumlah kata suatu data yang memiliki banyak kata benda tidak selalu akan masuk ke kelas negatif atau positif. Begitupun dengan kata sifat, kerja, keterangan, seru, dan persentase jumlah dari kata-kata tersebut. Kata-kata yang tidak termasuk ke dalam daftar kata bersentimen akan tetap tercatat sebagai bag of words features dan dihitung jumlah frekuensinya dari setiap dokumen. Fitur kata yang sering muncul di data latih dengan kelas negatif akan memiliki nilai frekuensi yang tinggi untuk kelas negatif juga, sehingga fitur kata tersebut dapat menjadi petunjuk dalam proses klasifikasi untuk menunjukkan mana data uji yang akan masuk kelas negatif dan sebaliknya, seperti kata “jelek” yang sering muncul pada data kelas negatif. Semakin banyak jumlah datanya maka fiturnya akan semakin beragam dan akan menghasilkan banyak keluaran klasifikasi yang tepat. Hasil pengujian model ini adalah sebagai berikut:



Gambar 10. Menampilkan Data Pie Gram dari kolom sentimen

Akurasi dari pengujian Naive Bayes Classifier menunjukkan hasil yang cukup signifikan, dengan skor 0.8536585365853658 menunjukkan bahwa model ini bekerja dengan hasil yang sangat baik, akurasi yang akurat dan dapat digunakan untuk menganalisis data baru yang jumlahnya lebih banyak lagi. Model ini menunjukkan bahwa semakin banyak jumlah data latihan yang di olah, maka akan semakin banyak juga variasi kata-kata baru maupun singkatan kata dari tweet yang berupa hashtag, retweet, comment yang dipelajari dan disimpan oleh model machine learning ini.

Keterangan Library Python yang digunakan:

Pandas berfungsi untuk membersihkan data mentah ke dalam sebuah bentuk tabel untuk analisis, melakukan perbandingan dan penggabungan set data dan juga penanganan data yang hilang.

Numpy, berfungsi untuk melakukan operasi vektor dan matriks dengan mengolah array dan array multidimensi, digunakan untuk kebutuhan dalam menganalisis data.

Regex, adalah deretan karakter yang digunakan untuk pencarian string atau teks dengan menggunakan pola

NLTK, berfungsi untuk proses pengolahan teks bahasa natural seperti classification, tokenization, stemming, tagging, parsing.

Matplotlib, adalah library python 2D yang dapat menghasilkan plot, digunakan untuk membuat dan menampilkan grafik hasil dari data yang di proses.

Sklearn, berfungsi untuk melakukan processing data ataupun melakukan training data sebagai kebutuhan untuk pembuatan model program machine-learning ini.

Sastrawi, adalah library python sederhana yang berfungsi untuk mengubah kata berimbuhan bahasa Indonesia yang tidak baku menjadi bentuk dasarnya.

BernoulliNB, adalah varian dari algoritma Naive Bayes digunakan sebagai fitur yang menentukan data latihan yang digunakan dan di proses bernilai negatif atau positif.

Untuk menyelesaikan penelitian ini, beberapa tools yang digunakan dengan spesifikasi sebagai berikut :

Tools	Software	Hardware
Colaboratory	Google	PC/Laptop
	Chrome	Windows 10
Excel	Microsoft	PC/Laptop
	Office	Windows 10

4. Kesimpulan

Perbedaan literature review dari jurnal referensi yang di citasi adalah tentang sumber data yang digunakan, serta pokok permasalahan yang melibatkan pengalaman sendiri maupun pengalaman publik yang dituangkan lewat media sosial twitter, berbagai kesan positif dan negatif yang terdapat dalam setiap postingan di twitter

membuktikan kualitas layanan dari provider Indosat Ooredoo memang perlu untuk ditingkatkan. Dari pengalaman pelanggan, sering mengeluh tentang banyaknya sms spam masuk, pulsa yang hilang tanpa digunakan, serta kecepatan internet yang lambat. Dari segi ini sangat perlu diperbaiki kualitas layanan dengan merespon setiap umpan balik yang masuk terutama di dalam media sosial, karena itu sangat mudah di akses oleh semua orang.

Variabel yang digunakan mengandung 2 opini yaitu positif dan negatif, semua kata dan kalimat yang digunakan dalam data latihan dikelompokkan berdasarkan tipe data dan opini masing-masing variabel tersebut. Setiap variabel yang memiliki imbuhan, contohnya kata 'jeleknya' akan di sortir dengan library Sastrawi yang digunakan untuk mendapatkan kata dasarnya, yang berarti kata 'jelek' masuk dalam variabel opini negatif.

Variabel yang digunakan dalam penentuan penelitian adalah berupa tabel sebagai berikut:

Variabel	Kata (Sample)
(x) Positif	Internet Murah, Internet Cepat, Indosat Baik, Banyak Promo, Banyak Bonus, Terima kasih Indosat, hemat, Senang, Puas.
(y) Negatif	Internet Lambat, Pulsa Hilang, Jaringan Jelek, Operator Bodoh, Indosat Parah, Indosat goblok, Kecewa, lemot, Indosat Maling Pulsa, Sinyal Kacau, Pulsa Ludes, Sinyal Error, Rugi Pulsa.

Klasifikasi analisis konten twitter untuk meninjau kualitas layanan provider Indosat Ooredoo dengan menggunakan sentimen positif dan negatif dilakukan dengan baik dengan metode algoritma Naive Bayes Classifier khususnya BernoulliNB. Namun ada beberapa aspek penting di dalam model yang perlu ditingkatkan seperti penggunaan library, keberagaman stemmer dan stopword list. Tujuannya agar meningkatkan skor akurasi dan prediksi pengklasifikasian dengan berbagai jenis pola penulisan berbahasa Inggris atau metode klasifikasi lain yang dapat diterapkan pada penelitian serupa, dengan dataset yang sama atau dataset yang berbeda. Dan juga bisa melakukan komparasi penelitian dengan menggunakan metode algoritma machine learning yang lain untuk menemukan hasil perbandingan terbaik yang lebih efektif dan lebih efisien. Untuk itu kami sarankan menggunakan dataset dengan ruang lingkup yang lebih banyak dan yang lebih beragam dari kategori jenisnya.

Daftar Rujukan

- [1] I. Wahyudi, W. Maharani, and A. B. Tjokorda, "Analisis Sentimen terhadap Produk di Twitter menggunakan Metode Support Vektor Machine," pp. 1-6, 2012.

- [2] D. Dutta Das, S. Sharma, S. Natani, N. Khare, and B. Singh, "Sentimental Analysis for Airline Twitter data," *IOP Conf. Ser. Mater. Sci. Eng.*, vol. 263, no. 4, 2017.
- [3] F. Hemmatian and M. K. Sohrabi, "A survey on classification techniques for opinion mining and sentiment analysis," *Artif. Intell. Rev.*, vol. 52, no. 3, pp. 1495–1545, 2019.
- [4] S. Vashishtha and S. Susan, "Fuzzy rule based unsupervised sentiment analysis from social media posts," *Expert Syst. Appl.*, vol. 138, 2019.
- [5] A. Pandhu and W. Diki, "Analisa sentimen dan Klasifikasi Komentar Positif Pada Twitter dengan Naïve Bayes Classification Sentiment Analysis and Classification of Positive Comments on Twitter with Naive Bayes Classification," vol. 1, no. 2, 2020.
- [6] S. Fransiska, "Seri Sains dan Teknologi ANALISIS SENTIMEN TWITTER UNTUK REVIEW FILM MENGGUNAKAN ALGORITMA NAIVE BAYES CLASSIFIER (NBC) PADA SENTIMEN R Jurnal Siliwangi Vol . 5 . No . 2 , 2019 Seri Sains dan Teknologi P-ISSN 2477-3891 E-ISSN 2615-4765," vol. 5, no. 2, 2019.
- [7] D. Ramayanti and U. Salamah, "Complaint Classification Using Support Vector Machine for Indonesian Text Dataset," vol. 3, no. 7, pp. 179–184, 2018.
- [8] W. Sharif *et al.*, "An empirical approach for extreme behavior identification through tweets using machine learning," *Appl. Sci.*, vol. 9, no. 18, 2019.
- [9] A. S. M. Alharbi and E. de Doncker, "Twitter sentiment analysis with a deep neural network: An enhanced approach using user behavioral information," *Cogn. Syst. Res.*, vol. 54, pp. 50–61, 2019.
- [10] A. C. E. S. Lima and L. N. De Castro, "Tecla: A temperament and psychological type prediction framework from Twitter data," *PLoS One*, vol. 14, no. 3, pp. 1–18, 2019.