

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA

Denni Kurniawan¹ dan Muhammad Yasir²

^{1,2}Magister Ilmu Komputer, Fakultas Teknologi Informasi,
Universitas Budi Luhur, Jakarta, Indonesia

E-mail: denni.kurniawan@budiluhur.ac.id, muhammad90yasir@gmail.com

Abstract

Freedom of expression on social media like Twitter will not always has positive effects, sometimes it contains negative things such as fake news, hate speech, and racism, where these kinds of tweet can be categorized as an act of Cyberbullying. This cyberbullying tends to increase every time. The aim of this study is to use the Naïve Bayes method in classifying types of sentiment on Twitter. The keyword used is Saipul Jamil, and the tweet was taken in September 2021. A total of 18,067 tweets were collected and then they will be labelled with a positive or negative value. This study also uses the CRIPS-DM method which is consist of Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment stages. The results of this study obtained the value of Accuracy (85.6%), Negative Recall (82.1%), Positive Recall (90.23%), and Negative Precision (91.76%) Positive Precision (79.18%).

Keywords: Cyberbullying, CRIPS-DM, Naïve Bayes, Sentiment Analysis, Twitter

Abstrak

Kebebasan mengeluarkan pendapat di media sosial Twitter tidak selamanya bernilai positif, karena terkadang mengandung hal negatif seperti cuitan yang mengandung berita bohong, penyebaran ujaran kebencian, dan rasisme. Ramai yang berpendapat bahwa jenis cuitan ini dapat dikategorikan sebagai *cyberbulling*. Dimana tindakan *cyberbulling* ini cenderung meningkat di setiap waktunya. Penelitian ini membahas penggunaan metode Naïve Bayes dalam melakukan klasifikasi jenis sentimen pada media sosial Twitter. Adapun kata kunci yang digunakan adalah Saipul Jamil, dan cuitan tersebut diambil pada bulan September 2021. Sebanyak 18.067 cuitan berhasil dikumpulkan dan dilanjutkan dengan memberikan label positif atau negatif pada data. Penelitian ini juga menggunakan metode CRIPS-DM dengan tahapan *Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, dan Deployment*. Hasil penelitian ini mendapatkan nilai Accuracy (85,6%), Recall Negatif (82,1%), Recall Positif (90,23%), dan Precision Negatif (91,76%) Precision Positif (79,18%).

Kata Kunci: Cyberbullying, CRIPS-DM, Naïve Bayes, Analisis Sentimen; Twitter.

1. Pendahuluan

Media sosial merupakan sebuah sarana komunikasi yang paling efektif, transparan dan efisien dimasa sekarang ini. Media sosial juga turut berperan penting dalam segala jenis perubahan dan pembaharuan yang terjadi [1]. Kebebasan berekspresi di media sosial memiliki tantangan yang sangat besar, dimana penyebaran berita yang bersifat negatif menjadi sangat mudah untuk disebar. Beberapa tindakan negatif yang sering dilakukan pengguna media sosial Twitter diantaranya adalah penyebaran berita bohong, penyebaran ujaran kebencian, dan rasisme. Tindakan ini merupakan perilaku negatif yang memiliki unsur tindakan *cyberbullying*. Tindakan yang mengandung *cyberbullying* di lingkungan pendidikan maupun di sosial media memiliki kecenderungan untuk terus naik [2].

Penelitian dengan topik *cyberbullying* merupakan salah satu penelitian yang sangat menarik untuk dilakukan, dimana opini masyarakat atau disebut juga dengan sentimen analisis di media sosial dapat dilihat dengan menggunakan beberapa metode populer dalam klasifikasi. Sebagai contoh penelitian yang dilakukan oleh Maulana dan Ernawati [3], melakukan klasifikasi *cyberbullying* di media sosial Twitter terhadap akun publik figur politik di Indonesia. Metode Naïve Bayes yang ditunakan berhasil melakukan klasifikasi dengan tingkat akurasi sebesar 76%. Penelitian lainnya yang dilakukan oleh Saputril dan Zuhri [4] menggunakan metode Naïve Bayes dalam melakukan klasifikasi terhadap ujaran kebencian pada periode Pemilihan Presiden 2019. Penggunaan metode Naïve Bayes yang digunakan berhasil mendapatkan tingkat akurasi sebesar 71%. Penelitian lain yang dilakukan oleh Sulastri [5], menggunakan metode Naïve Bayes pada analisis sentimen kasus penolakan dibukanya larangan ekspor benih Lobster. Penggunaan Naïve Bayes pada penelitian ini berhasil mendapatkan nilai akurasi sebesar 72,5%.

Dari beberapa penelitian di atas, dapat dilihat bahwa dengan menggunakan metode Naïve Bayes, tingkat akurasi yang dihasilkan masih belum terlalu tinggi. Sehingga tingkat akurasi masih bisa ditingkatkan lagi dengan menggunakan fitur tambahan. Adapun tujuan dari penelitian ini adalah untuk melakukan analisis sentimen *cyberbullying* dengan menggunakan metode CRIPS-DM dan Naïve Bayes dengan tingkat akurasi yang lebih tinggi dibandingkan penelitian sebelumnya.

2. Metode Penelitian

Pada bagian ini, akan dijelaskan teori-teori dasar yang digunakan pada penelitian ini, yaitu penjelasan *cyberbullying*, metode Naïve Bayes, dan model CRISP-DM.

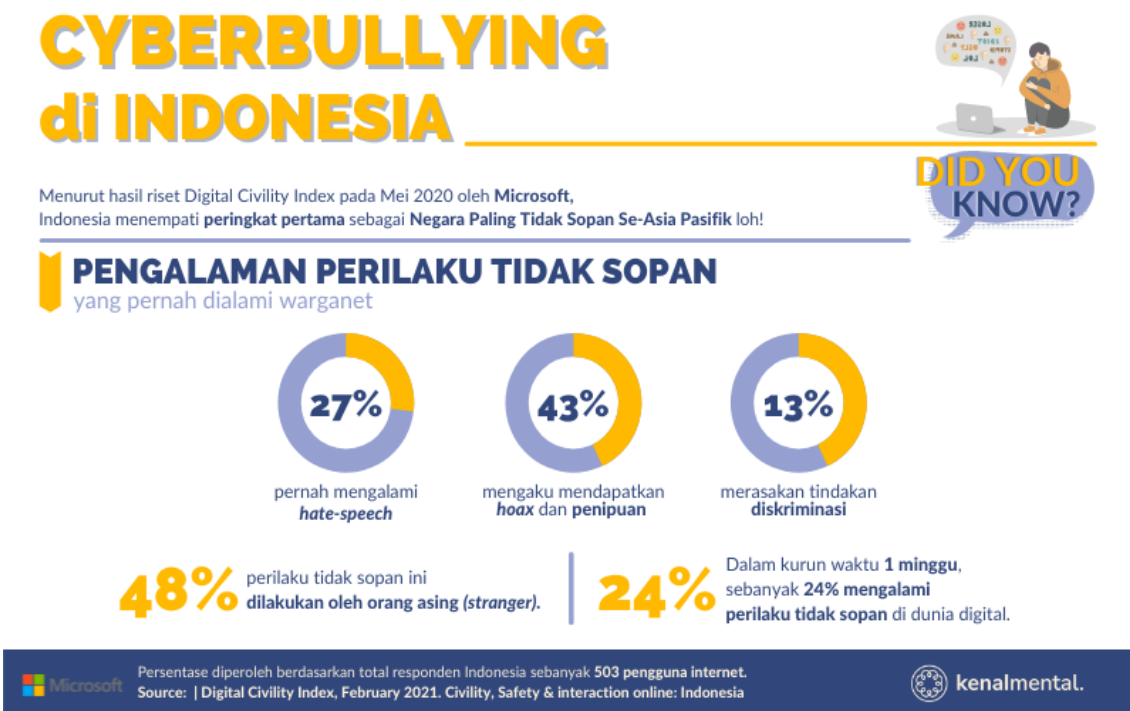
Cyberbullying

Ada beberapa contoh tindakan yang dilakukan pengguna media sosial dengan menggunakan perangkat teknologi yang dapat dikategorikan sebagai perbuatan *cyberbullying*, diantaranya adalah intimidasi, ujaran kebencian, memermalukan seseorang, rasisme, atau body shaming yang dilakukan oleh pengguna media sosial menggunakan perangkat teknologi [6]. Menurut Alanzi dan Alves-foss [7], terdapat juga beberapa jenis perilaku lainnya yang dianggap sebagai tindakan *cyberbullying*, diantaranya adalah:

1. *Flooding* yaitu melibatkan pelaku intimidasi berulang kali mengirimkan komentar/postingan yang tidak penting sehingga tidak memungkinkan korban target untuk berpartisipasi dalam percakapan.
2. *Masquerade* yaitu di mana pelaku intimidasi berpura-pura meniru atau menyamar sebagai korban.

3. *Flaming/Bashing* melibatkan perkelahian online di mana pelaku intimidasi mengirim atau memposting konten yang menghina, menyakitkan, dan vulgar kepada korban yang ditargetkan secara pribadi atau publik dalam grup online.
4. *Trolling* melibatkan dengan sengaja menerbitkan komentar yang tidak sesuai dengan komentar lain untuk memicu argumen atau emosi negatif meskipun komentar itu sendiri mungkin tidak vulgar atau menyakitkan.
5. Pelecehan adalah jenis percakapan di mana pelaku intimidasi sering mengirimkan pesan yang menghina dan kasar kepada korban secara pribadi.
6. Penghinaan, terjadi ketika pelaku intimidasi mengirimkan atau mempublikasikan gosip atau pernyataan palsu tentang korban untuk merusak persahabatan/reputasi korban
7. Outing terjadi ketika seorang pengganggu memposting atau mempublikasikan informasi pribadi atau memalukan di ruang obrolan atau forum publik. Jenis *cyberbullying* ini mirip dengan fitnah.
8. Pengucilan, yaitu dengan sengaja mengecualikan seseorang dari grup online (mengucilkan orang lain). Jenis *cyberbullying* ini terjadi di kalangan remaja dan remaja lebih menonjol.

Seperti yang dikemukakan oleh Abdussalam, bahwa tindakan yang mengandung *cyberbullying* di lingkungan pendidikan maupun di sosial media memiliki kecenderungan untuk terus naik [2]. Menurut penelitian yang dilakukan oleh Microsoft, pengalaman perilaku tidak sopan yang dialami oleh pengguna sosial media juga tidak dapat dianggap enteng. Seperti terlihat pada Gambar 1, terlihat sebanyak 27% responden pernah menjadi korban ujaran kebencian, 43% responden mendapatkan berita bohong, dan 13% responden pernah merasakan tindakan diskriminasi [8].



Gambar 1. Riset Microsoft tentang *cyberbullying* di Indonesia [8].

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA

Metode Naïve Bayes

Penggunaan Metode Naive Bayes sering diimplementasikan pada klasifikasi di berbagai bidang. Menurut Khairati [9], bahwa hasil klasifikasi yang dilakukan dalam implementasinya menunjukkan hasil yang beragam. Meskipun pada beberapa penelitian yang menggunakan optimalisasi Naïve Bayes menunjukkan kebenaran, namun terkadang tingkat perkiraan probabilitasnya tidak terlalu tinggi. Sehingga masih memerlukan pemahaman yang lebih dalam terhadap karakteristik data yang mempengaruhi kinerja metode ini. Metode Naïve Bayes juga menyediakan mekanisme untuk menggunakan informasi dalam sampel data untuk memperkirakan probabilitas posterior $P(H|X)$ dari setiap kelas H yang diberikan objek X . Setelah memiliki perkiraan seperti itu, baru dapat menggunakannya untuk klasifikasi atau aplikasi pendukung keputusan lainnya [10]. Menurut pemaparan Bustami [11], formula yang digunakan dalam metode Naïve Bayes adalah seperti diperlihatkan pada persamaan (1) di bawah:

$$P(H|X) = \frac{P(H|X) \cdot P(H)}{P(X)} \quad (1)$$

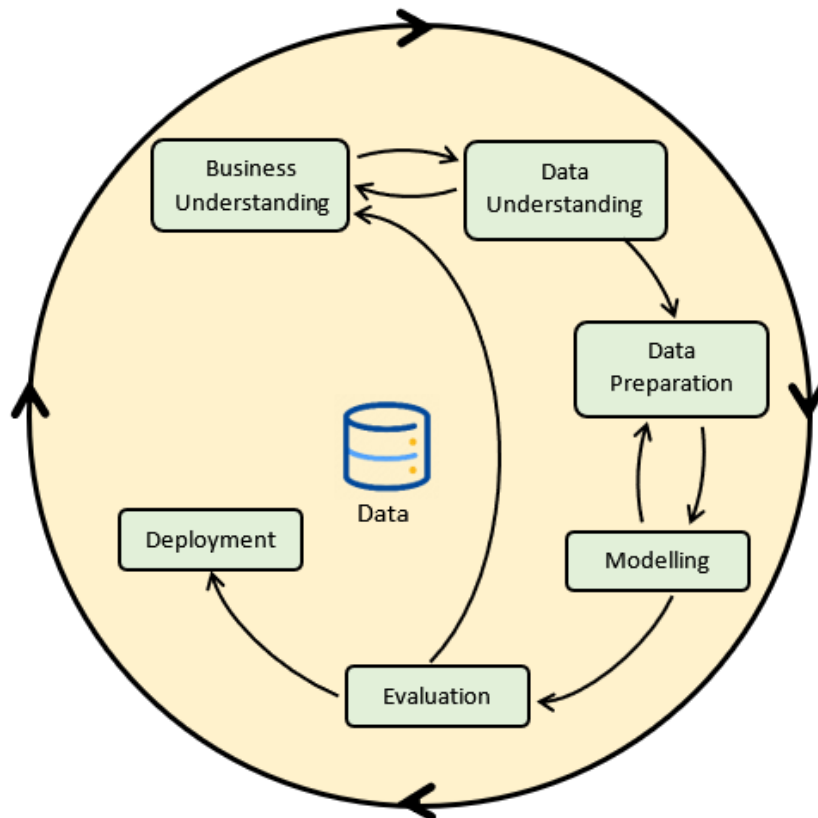
Dimana X adalah data dengan class yang belum diketahui; H adalah hipotesis data X meruapakan suatu class spesifik; $P(H|X)$ adalah probabilitas hipotesis H berdasarkan kondisi X ; $P(H)$ adalah probabilitas hipotesis H ; $P(X|H)$ adalah probabilitas X berdasarkan kondisi pada hipotesis H ; dan $P(X)$ adalah probabilitas X .

Model CRISP-DM

Cross Industry Standard Process for Data Mining (CRISP-DM) adalah model proses industri-independen yang lazim digunakan untuk keperluan *data mining*. Metode ini terdiri dari enam fase berulang seperti terlihat pada Gambar 2. Fase pertama dimulai dari *Business Understanding*, *Data Understanding*, *Data Preparation*, *Modelling*, *Evaluation*, dan *Deployment* [12, 13]. Dan dalam siklus pengembangannya CRISP-DM dianggap sebagai metodologi data mining terlengkap dalam hal pemenuhan kebutuhan proyek industri, dan telah menjadi yang paling luas penggunaannya dalam proyek analitik, *data mining*, serta *data science* [14]. Tahapan penelitian menggunakan CRISP-DM terlihat pada Gambar 2, yang merupakan aliran proses pada model CRISP-DM yang biasa digunakan untuk *data mining*. Setiap proses atau fase memiliki tugas dan fungsi masing-masing. Panah pada gambar menunjukkan keterikatan antar fase, namun ada kalanya urutan fase ini tidak kaku dan juga tergantung pada hasil di setiap fase. Lingkaran luar dari setiap fase juga melambangkan siklus *data mining* itu sendiri. Model CRISP-DM juga bertujuan untuk membuat proyek penambangan data besar, lebih murah, lebih andal, lebih dapat diulang, lebih mudah dikelola, dan lebih cepat.

Fase pertama dari model CRISP-DM adalah *Business Understanding*, dimana fase ini berfokus untuk memahami tujuan dan persyaratan proyek dari perspektif bisnis. Pengetahuan ini kemudian diubah menjadi definisi masalah penambangan data dan rencana proyek awal yang dirancang untuk mencapai tujuan. Kemudian dilanjutkan dengan fase *Data Understanding* atau pemahaman data, dimana fase ini dimulai dengan proses pengumpulan data dan dilanjutkan dengan kegiatan mengidentifikasi data, mengidentifikasi masalah kualitas data, menemukan wawasan pertama ke dalam data, atau mendeteksi subset yang menarik untuk membentuk hipotesis untuk informasi

tersembunyi. Di tahap ini penulis melakukan pelabelan *dataset* manual dengan bantuan anator, ada baiknya dalam proses pelabelan ini melibatkan ahli atau pakar bahasa sesuai bahasa yang digunakan. Tujuannya adalah ahli atau pakar bahasa tersebut lebih memahami pemaknaan dari setiap teks yang tertulis [15].



Gambar 2. Ilustrasi tahapan pelaksanaan CRISP-DM [12]

Dilanjutkan dengan fase *Data Preparation* atau persiapan data, yang mencakup semua kegiatan untuk membangun dataset akhir dari data mentah awal. Tugas persiapan data kemungkinan akan dilakukan beberapa kali, tetapi tidak dalam urutan yang ditentukan. Tugas seperti tabel, catatan, dan pemilihan atribut, *preprocessing* data, konstruksi atribut baru, dan transformasi data untuk pemodelan. Tahap *preprocessing* atau pembersihan data Twitter adalah skenario standar *preprocessing* data teks. Menurut penelitian yang dilakukan oleh Giachanou dan Crestani [16], terdapat beberapa karakteristik unik tweet yang memerlukan langkah yang berbeda untuk mengatasinya. Berdasarkan penelitian sebelumnya yang dilakukan oleh Kane *et al.*, terdapat beberapa langkah yang digunakan untuk *preprocessing* Twitter yaitu *Tokenize*, *Stopword*, *Case Folding*, *Remove Punctuation*, dan *Stemming* [17]. Pada tahap tahapan *preprocessing* Twitter sebagai berikut:

1. *Tokenize*: Tahap ini merupakan proses pembagian kata. Kalimat akan dibagi menjadi beberapa bagian yang disebut token. Token dapat berupa kata, frasa, atau elemen makna lainnya. Kita bisa menggunakan *library nltk.tokenize* Kita bisa menggunakan *library nltk.tokenize* untuk membagi kata menjadi kelompok huruf.
2. *Stopword*: Tahap ini merupakan proses untuk menghilangkan kata-kata umum dan sering yang tidak memiliki pengaruh signifikan dalam sebuah kalimat. sehingga data

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA

dapat diproses lebih efisien pada tahap selanjutnya. Mengimpor daftar *stopwords* bisa menggunakan dari *library nltk.corpus*.

3. *Case Folding*: Tahap ini merupakan proses mengubah kata menjadi bentuk yang sama. Pada langkah ini, mengonversi semua kata menjadi huruf kecil menggunakan metode *lower case* pada Python.
4. *Remove Punctuation*: Tahap ini merupakan proses menghapus karakter, angka, *string ASCII*, dan tanda baca. Pesan *twitter* biasanya berisi simbol, angka, dan tanda baca. Semua ini dihapus menggunakan sintaks ekspresi reguler
5. *Stemming*: Tahap ini adalah proses mendapatkan dasar atau akar kata dengan menghilangkan imbuhan dan sufiks. Penelitian ini menggunakan *library python* sastrawi untuk menghilangkan kata-kata imbuhan dalam bahasa Indonesia ke bentuk dasarnya.

Setelah fase ketiga selesai dilakukan, maka dilanjutkan dengan fase berikutnya yaitu fase *Modeling*. Dimana pada fase ini terdapat berbagai teknik pemodelan dipilih dan diterapkan. Beberapa metode memerlukan format data tertentu. Pada penelitian ini, tahap *Modeling* menggunakan metode Naïve Bayes untuk mendapatkan nilai optimal. Tahapan berikutnya adalah *Evaluation*, dimana pada tahapan ini satu atau lebih model dibangun untuk menghasilkan analisis data yang memiliki kualitas tinggi. Sebelum melanjutkan ke implementasi akhir model, penting untuk mengevaluasi model secara lebih menyeluruh, dan meninjau langkah-langkah yang diambil untuk membangun model, untuk memastikan bahwa model mencapai tujuan bisnis dengan benar. Tujuan utamanya adalah untuk menentukan apakah ada beberapa masalah yang belum dipertimbangkan secara memadai. Pada akhir fase ini, keputusan harus dibuat tentang bagaimana menggunakan hasil data mining. Fase terakhir adalah *Deployment*, dimana pada fase akhir ini pemodelan bukanlah akhir dari proyek. Secara umum, pengetahuan yang diperoleh perlu diatur dan disajikan sedemikian rupa sehingga pengguna dapat menggunakannya. Bergantung pada apa yang dibutuhkan pengguna, fase *Deployment* bisa dibuat sesederhana membuat laporan atau serumit menerapkan proses data mining berulang.

3. Hasil dan Pembahasan

Penelitian ini dilakukan dengan tahapan yang sesuai dengan pengembangan dari tahapan CRIPS-DM yang telah dijelaskan sebelumnya, yaitu terdiri dari enam tahapan yaitu *Business Understanding*, *Data Understanding*, *Data Preparation*, *Modeling*, *Evaluation*, dan *Deployment*. Penjelasan secara lebih terperinci terkait dengan tahapan-tahapan tersebut adalah seperti berikut:

a) *Business Understanding*

Tahapan *Business Understanding* berfokus pada pemahaman tujuan yaitu melakukan *experiment research* bagaimana melakukan klasifikasi sentimen yang mengandung unsur *cyberbullying* pada media sosial Twitter. Kemudian data Twitter dikumpulkan dengan kata kunci “Saipul Jamil” yang diambil selama delapan hari, dimulai pada tanggal 2 sampai tanggal 9 September 2021 menggunakan API Twitter dengan menggunakan Python Google Colab.

b) *Data Understanding*

Tahapan selanjutnya yaitu tahapan *Data Understanding* yaitu proses pendataan awal terkait tweet *cyberbullying* penyanyi dangdut Saipul Jamil. Dari proses pengumpulan data

selama delapan hari ini berhasil mendapatkan sebanyak 18.0067 tweet. Di tahap ini penulis melakukan pelabelan *dataset* yang dilakukan secara manual. Proses ini juga melibatkan seorang ahli bahasa atau pakar bahasa sesuai dengan bahasa yang digunakan, yang bertujuan agar pakar bahasa tersebut dapat memberikan pelabelan tanpa kesalahan dan memiliki kemampuan pemahaman yang lebih tinggi terhadap setiap teks yang tertulis pada penelitian.

c) Data Preparation

Tahap *Data Preparation* merupakan tahap persiapan data yang mencakup semua kegiatan untuk membangun dataset akhir. Pada tahap ini, dilakukan beberapa preprocessing text yaitu tahap *Tokenize*, *Case Folding & Remove Punctuation*, *Stopword*, dan *Stemming*. Pada tahapan awal yang disebut *Tokenize*, adalah proses yang berfungsi untuk memisahkan teks menjadi potongan-potongan kata yang menyusunnya. Proses pemisahan ini dapat dilakukan dengan menggunakan `nlk.word_tokenize` pada Python. Perbedaan hasil antara sebelum dan sesudah tahapan *Tokenize* dilakukan dapat dilihat pada Tabel 1.

TABEL 1. IMPLEMENTASI *TOKENIZE*

No	Sebelum	Sesudah
1	Bebas dari Penjara, Saipul Jamil Ngaku Trauma\n@saipuljamil\n https://t.co/sgRJ4nXieH	['Bebas', 'dari', 'Penjara', ',', 'Saipul', 'Jamil', 'Ngaku', 'Trauma', '@', 'saipuljamil', 'https', ':', '//t.co/sgRJ4nXieH']
2	@saipuljamil Tuhan Maha Pemaaf dan menerima Taubat.Taubat artinya tdk melakukan perbuatan tercela yg sama lagi. Ini hanya Peringatan Bang Saipul Jamil ,Akibat perbuatan anda trauma seumur hidup diterima korban. Tdk bisa perbaiki kerusakan yg dibuat. Jadi Jgn sampai berbuat lagi ! https://t.co/iayMFyx0TT	['@', 'saipuljamil', 'Tuhan', 'Maha', 'Pemaaf', 'dan', 'menerima', 'Taubat.Taubat', 'artinya', 'tdk', 'melakukan', 'perbuatan', 'tercela', 'yg', 'sama', 'lagi', ',', 'Ini', 'hanya', 'Peringatan', 'Bang', 'Saipul', 'Jamil', ',', 'Akibat', 'perbuatan', 'anda', 'trauma', 'seumur', 'hidup', 'diterima', 'korban', ',', 'Tdk', 'bisa', 'perbaiki', 'kerusakan', 'yg', 'dibuat', ',', 'Jadi', 'Jgn', 'sampai', 'berbuat', 'lagi', '!', 'https', ':', '//t.co/iayMFyx0TT']

Tahapan selanjutnya, yaitu Tahap *Case folding* dimulai dengan proses merubah teks menjadi huruf kecil. Sedangkan *Remove Punctuation* adalah proses menghilangkan tanda baca atau simbol yang ada dalam dataset. Perbedaan hasil antara sebelum dan sesudah tahapan *Tokenize* dilakukan dapat dilihat pada Tabel 2.

TABEL 2. IMPLEMENTASI *CASE FOLDING* DAN *REMOVE PUNCTUATION*

No	Sebelum	Sesudah
1	Bebas dari Penjara, Saipul Jamil Ngaku Trauma\n@saipuljamil\n https://t.co/sgRJ4nXieH	bebas dari penjara saipul jamil ngaku trauma
2	@saipuljamil Tuhan Maha Pemaaf dan menerima Taubat.Taubat artinya tdk melakukan perbuatan tercela yg sama lagi. Ini hanya Peringatan Bang Saipul Jamil ,Akibat perbuatan anda trauma seumur hidup diterima korban. Tdk bisa perbaiki kerusakan yg dibuat. Jadi Jgn sampai berbuat lagi ! https://t.co/iayMFyx0TT	tuhan maha pemaaf dan menerima taubattaubat artinya tdk melakukan perbuatan tercela yg sama lagi ini hanya peringatan bang saipul jamil akibat perbuatan anda trauma seumur hidup diterima korban tdk bisa perbaiki kerusakan yg dibuat jadi jgn sampai berbuat lagi

Tahapan akhir yang disebut tahap *Stopword* merupakan suatu proses yang berfungsi untuk menghilangkan kata-kata umum dan sering yang tidak memiliki pengaruh

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA

signifikan dalam sebuah kalimat. sehingga data dapat diproses lebih efisien pada tahap selanjutnya. Proses *import* daftar *Stopwords* dapat dilakukan dengan menggunakan *library* *nlTK.corpus*. Perbedaan hasil antara sebelum dan sesudah tahapan *Tokenize* dilakukan dapat dilihat pada Tabel 3.

TABEL 3. IMPLEMENTASI *STOPWORD*

No	Sebelum	Sesudah
1	bebas dari penjara saipul jamil ngaku trauma	bebas penjara ngaku trauma
2	tuhan maha pemaaf dan menerima taubattaubat artinya tdk melakukan perbuatan tercela yg sama lagi ini hanya peringatan bang saipul jamil akibat perbuatan anda trauma seumur hidup diterima korban tdk bisa perbaiki kerusakan yg dibuat jadi jgn sampai berbuat lagi	tuhan maha pemaaf menerima taubattaubat tdk perbuatan tercela yg peringatan bang akibat perbuatan trauma seumur hidup diterima korban tdk perbaiki kerusakan yg jgn berbuat

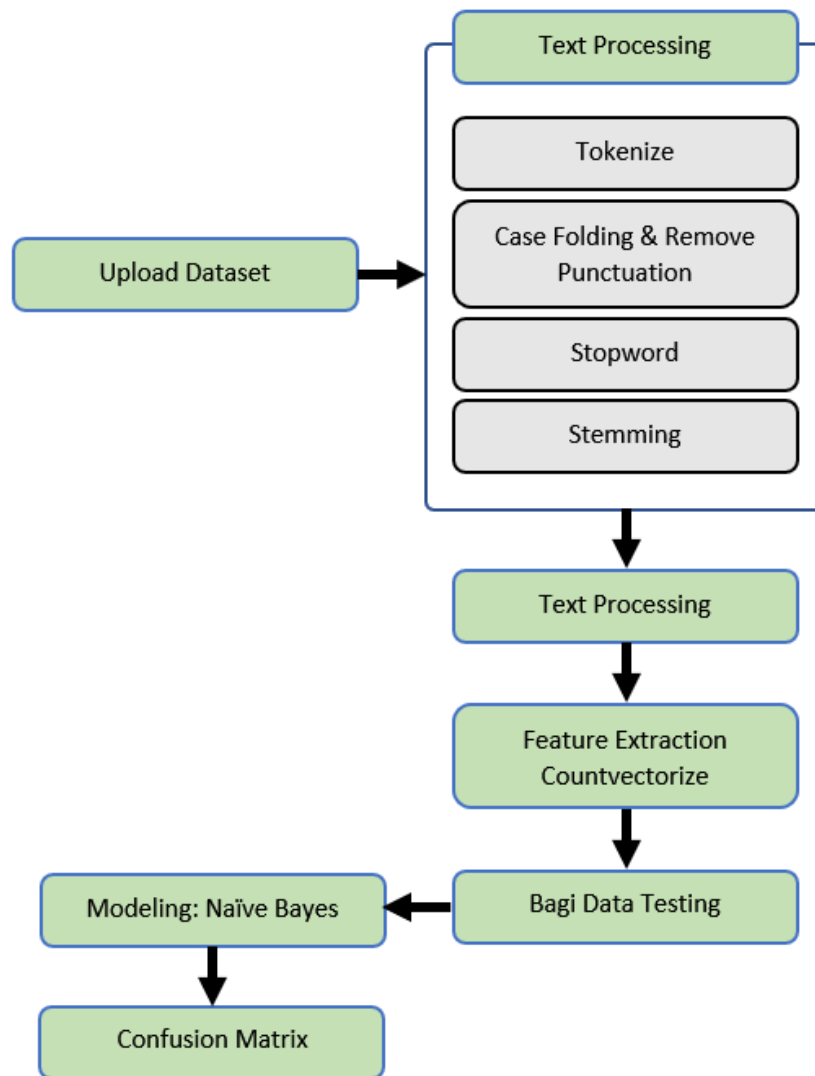
Proses mengubah kata berimbuhan menjadi kata dasar dapat dilakukan pada Tahap *Stemming*. Pada tahap ini menggunakan *library* Python Sastrawi untuk menghilangkan kata-kata imbuhan dalam bahasa Indonesia ke bentuk dasarnya. Perbedaan hasil antara sebelum dan sesudah tahapan *Tokenize* dilakukan dapat dilihat pada Tabel 4.

TABEL 4. IMPLEMENTASI *STEMMING*

No	Sebelum	Sesudah
1	bebas penjara ngaku trauma	bebas penjara ngaku trauma
2	tuhan maha pemaaf menerima taubattaubat tdk perbuatan tercela yg peringatan bang akibat perbuatan trauma seumur hidup diterima korban tdk perbaiki kerusakan yg jgn berbuat	tuhan maha maaf terima taubattaubat tdk buat cela yg ingat bang akibat buat trauma umur hidup terima korban tdk baik rusak yg jgn buat

d) Modeling

Pada tahap ini merupakan proses *mining* menggunakan model Naïve Bayes dengan alur seperti terlihat pada Gambar 3. Langkah pertama adalah melakukan *upload dataset* mentah yang sudah diberi label positif atau negatif. Kemudian dilanjutkan dengan melakukan proses *text preprocessing* menggunakan Google Collab. Selanjutnya, dilakukan *text preprocessing* data dengan melakukan *Tokenize*, *Case Folding* dan *Remove Punctuation*, *Stopword*, *Stemming*. Terdapat 18.067 data mentah dan setelah melakukan proses *text preprocessing*, hanya tersisa 6.787 data. Setelah proses pembersihan data dilakukan, maka dilanjutkan dengan melakukan perubahan informasi label tersebut menjadi angka. Fungsi ini untuk mengubah data ke dalam bentuk data integer positif (1) atau negatif (-1). Kemudian dari 6.787 data tersebut, data kembali dibagi menjadi sebanyak 2.500 yang digunakan sebagai sampel data yang didapatkan melalui perhitungan Slovin. Kemudian data ini di ubah ke dalam bentuk vektor menggunakan fitur *extraction counvectorize* untuk digunakan sebagai masukan dalam algoritma pembelajaran mesin. Dari sebanyak 2.500 data yang dimiliki, data kembali dibagi menjadi dua bagian, yang akan digunakan untuk kebutuhan *training*, dan *testing*. Pembagian data ini menggunakan perbandingan 80:20, yang bermaksud menggunakan sebanyak 80% data sebagai data *training* dan 20% data sebagai data *testing*. Setelah itu pada tahap selanjutnya dilakukan modeling menggunakan metode Naïve Bayes dan dilakukan perhitungan untuk melihat akurasi dari model. Setelah proses modeling, tahap selanjutnya adalah melakukan *Confusion Matrix* untuk mengukur performa klasifikasi.

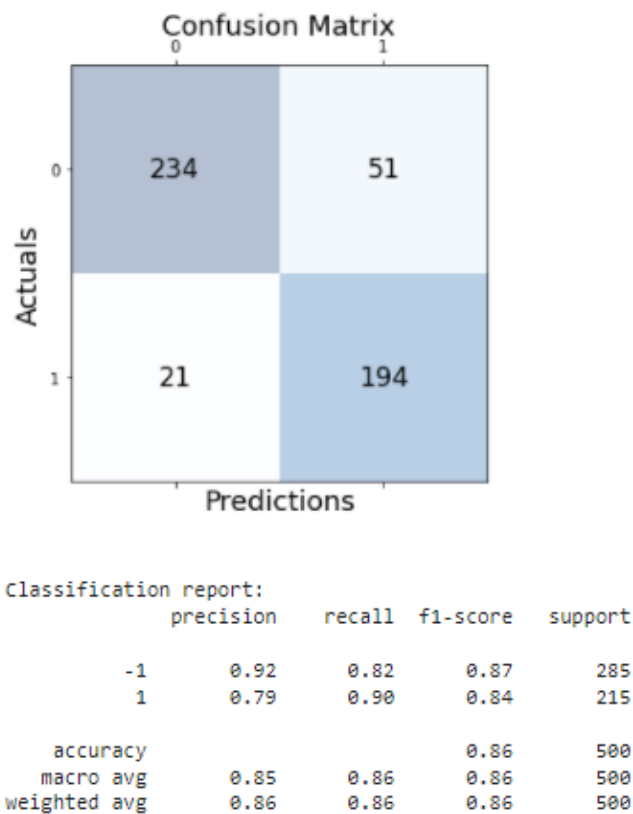


Gambar 3. Alur proses modeling

e) Evaluation

Tahapan *Evaluation* merupakan tahap laporan evaluasi model dengan menggunakan *Confusion Matrix* untuk mengukur performa klasifikasi Naïve Bayes. Perhitungan untuk menentukan *Accuracy*, *Recall Negative*, *Recall Positive*, *Precision Negative*, dan *Precision Positive* dapat dilakukan dengan menggunakan Persamaan (2), (3), (4), (5), dan (6) secara berurutan. Dimana TP adalah *True Positive*, TN adalah *True Negative*, FP adalah *False Positive*, dan FN adalah *True Negative*.

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA



Gambar 4. Hasil Confusion Matrix

Confusion Matrix	Predicted Class		
	Positive	Negative	
Observed Class	Positive	TP = 234	FN = 21
	Negative	FP = 51	TN = 194

Gambar 5. Penjelasan pengambilan nilai TP, TN, FP, dan FN dari *Confusion Matrix*

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (2)$$

$$Recall (neg) = \frac{TP}{TP + FP} \quad (3)$$

$$Recall (pos) = \frac{TN}{TN + FN} \quad (4)$$

$$Precision (neg) = \frac{TP}{TP + FN} \quad (5)$$

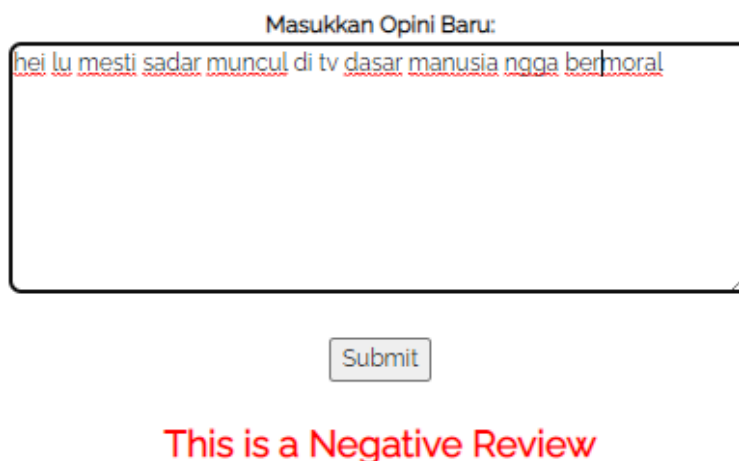
$$Precision (pos) = \frac{TN}{TN + FP} \quad (6)$$

Nilai TP, TN, FP, dan FN dapat diambil dari *Confusion Matrix* seperti pada Gambar 5. Sehingga diperoleh nilai TP, TN, FP, dan FN sebesar 234, 194, 51, dan 21, secara berurutan. Nilai ini kemudian dimasukkan ke dalam *Persamaan* (2), (3), (4), (5), dan (6)

sehingga dari hasil perhitungan diperoleh nilai *Accuracy* 85,6%; *Recall Positive* 90,23%; *Recall Negative* 82,1%; dan *Precision Positif* 79,18%.

f) Deployment

Sebuah prototipe hasil pengembangan model metode Naïve Bayes menggunakan FLASK telah dilakukan pada tahap *Deployment*, seperti diperlihatkan pada Gambar 6. Prototipe yang dibuat memiliki kemampuan untuk mendeteksi kalimat atau tweet yang memiliki unsur *cyberbullying*. Sistem ini juga mampu mengklasifikasikan kalimat opini baru tersebut dengan label *positive* atau *negative*. Prototype ini memiliki kemampuan yang sangat bagus, walaupun adanya keterbatasan penggunaan jumlah karakter yang berlaku di Twitter, yaitu maksimal 140 karakter.



Gambar 6. Implementasi Sistem

4. Kesimpulan

Berdasarkan hasil penelitian yang diperoleh, dapat disimpulkan bahwa algoritma Naïve Bayes yang digunakan untuk mengklasifikasikan tweet yang mengandung unsur *cyberbullying*. Proses pelatihan ini perlu didukung dengan data latih yang baik. Proses pemberian label pada dataset sebaiknya dilakukan dengan bantuan seorang anator atau ahli bahasa agar lebih memahami pemaknaan dari setiap teks yang tertulis. Dari hasil evaluasi menggunakan *Confusion Matrix* didapatkan hasil analisis sentimen *cyberbullying* menggunakan metode Naïve Bayes mendapatkan *Accuracy* sebesar 85,6%, *Recall Positive* sebesar 90,23%; *Recall Negative* sebesar 82,1%; *Precision Negative* sebesar 91,76%; dan *Precision Positive* sebesar 79,18%. Terlihat bahwa hasil pengujian yang dilakukan memiliki tingkat akurasi yang lebih baik dibandingkan hasil penelitian yang serupa.

Referensi

- [1] D. R. Rahadi, "Perilaku Pengguna dan Informasi Hoax di Media Sosial". *Jurnal Manajemen dan Kewirausahaan*, vol.5(1), pp. 58-70. 2017.
- [2] M. S. Abdussalam, Sejumlah kasus bullying sudah warnai catatan masalah anak di awal 2020 begini kata komisioner KPAI, <https://www.kpai.go.id/publikasi/sejumlah-kasus-bullying-sudah-warnai-catatan-masalah-anak-di-awal-2020-begini-kata-komisioner-kpai>, 2020, retrieved February 13, 2020.

OPTIMIZATION SENTIMENT ANALYSIS USING CRISP-DM AND NAÏVE BAYES METHODS IMPLEMENTED ON SOCIAL MEDIA

- [3] F. A. Maulana, and I. Ernawati, “Analisa Sentimen Cyberbullying di Jejaring Sosial Twitter dengan Algoritma Naïve Bayes”, *Seminar Nasional Mahasiswa Ilmu Komputer dan Aplikasinya*, vol. 1(2), pp. 529-538. 2020.
- [4] N. A. O. Saputril, and K. Zuhri, “Analisis sentimen masyarakat terhadap pilpres 2019 berdasarkan opini dari Twitter menggunakan metode Naïve Bayes classifier”. *Jurnal Informanika*, vol. 7(1), pp. 55-62, 2021.
- [5] K. Sulastrri, “Klasifikasi Naïve Bayes pada analisis sentimen atas penolakan dibukanya larangan ekspor benih lobster”, *Jurnal Riset Inovasi Bidang Informatika dan Pendidikan Informatika KERNEL*, vol. 1(2), pp. 68-75. 2020.
- [6] F. A. Maulana, I. Ernawati, P. Labu, and J. Selatan, “Analisa sentimen cyberbullying di jejaring sosial Twitter dengan algoritma Naïve Bayes”, *Seminar Nasional Mahasiswa Ilmu Komputer dan Aplikasi (SENAMIKA)*, pp. 529–538. 2020.
- [7] I. Alanazi, and J. Alves-foss, “Cyber bullying and machine learning: A survey”, *International Journal of Computer Science Security*, vol. 18(10), pp. 1-8, 2020.
- [8] Anon, Civility, Safety & Interaction Online: Indonesia, [shttps://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE4MM81](https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE4MM81), 2021, retrieved November 10, 2021
- [9] A. F. Khairati, A. A. Adlina, G. F. Hertono, and B. D. Handari, “Kajian Indeks Validitas pada Algoritma K-Means Enhanced dan K-Means MMCA”, *Prosiding Seminar Nasional Matematika*, vol. 2. pp. 161-170, 2019.
- [10] G. I. Webb. “Encyclopedia of Machine Learning and Data Mining”, *Encyclopedia of Machine Learning and Data Mining*
- [11] Bustami, “Penerapan Algoritma Naïve Bayes untuk Mengklasifikasi Data Nasabah”, *TECHSI: Jurnal Penelitian Teknik Informatika*, vol. 4, pp. 127-146. 2010.
- [12] W. Rüdiger, and J. Hipp. “CRISP-DM: Towards a Standard Process Model for Data Mining.” *Proceedings of the 4th international conference on the practical applications of knowledge discovery and data mining*, vol. 4, pp. 29–39. 2000.
- [13] P. Chapman, J. Clinton, R. Kerber, T. Khabaza, T. Reinartz, C. Shearer, and R. Wirth. “CRISP-DM 1.0: Step-by-step Data Mining Guide,” 2000.
- [14] C. Schröer, F. Kruse, and J. M. Gómez, “A Systematic Literature Review on Applying CRISP-DM Process Model”, *Procedia Computer Science*, vol. 181 (2019), pp. 526–534. 2021.
- [15] M. S. Hadna, P. I. Santosa, dan W. W. Winarno, “Studi Literatur Tentang Perbandingan Metode Untuk Proses Analisis Sentimen Di Twitter”, *Seminar Nasional Teknologi Informasi dan Komunikasi 2016 (Sentika)*, 2016, pp. 57–64.
- [16] A. Giachanou, dan F. Crestani, “Like it or not: A survey of Twitter sentiment analysis methods”, *ACM Computing Surveys*, vol. 49 (2), 2016. pp. 1-41.
- [17] S. N. Kane, A. Mishra, dan A. K. Dutta, “Preface: International Conference on Recent Trends in Physics (ICRTP 2016)”, *Journal of Physics: Conference Series*, vol. 755 (1), 2016, pp. 3–9.