

Sentimen Pada Komentar Youtube Tentang Pencegahan Dan Penanganan Kekerasan Seksual Pada Permendikbud Berbasis Naïve Bayes Dan Support Vector Machine

Muhammad Derryl Qinanda¹, Agung Nilogiri², Taufiq Timur W³

^{1,2,3} Fakultas Teknik, Universitas Muhammadiyah Jember

Email : ¹derrylqinanda@gmail.com, ²agungnilogiri@unmuhjember.ac.id,

³taufiqtimur@unmuhjember.ac.id

(Naskah masuk: 16 Mei 2022, diterima untuk diterbitkan: 3 Agustus 2022, terbit: 28 Agustus 2022)

ABSTRAK

Perkembangan teknologi informasi membawakan perubahan modern untuk masyarakat. Beberapa penggunaan media sosial, diantaranya adalah untuk mendapatkan dan menyampaikan informasi kepada masyarakat atau kerabat. Salah satu dari beberapa media sosial yang saat ini sering digunakan oleh masyarakat yaitu situs *Youtube*. Sentimen Analisis merupakan sebuah teknik dimana untuk mengekstrak sebuah data yang berbentuk teks yang digunakan untuk memperoleh sebuah informasi tentang sentimen bernilai positif dan negatif. Ruang lingkup penelitian ini dilakukan hanya pada komentar masyarakat terhadap permasalahan pencegahan dan penanganan kekerasan seksual di lingkungan perguruan tinggi pada *channel* MataNajwa. Selanjutnya data tersebut akan dilakukan pengolahan data menggunakan metode *Naïve Bayes* dan *Support Vector Machine* menggunakan ekstraksi fitur *TF-IDF*. Hasil akurasi yang didapatkan pada penelitian ini ialah 64% dengan menggunakan *Naïve Bayes* dan dilakukan pengujian dengan menggunakan data yang dipilih menggunakan *K-Fold Cross Validation* lalu menghasilkan akurasi sebesar 51,7%. Sedangkan nilai akurasi yang didapatkan saat menggunakan *Support Vector Machine* ialah sebesar 92% dengan dilakukan pengujian menggunakan *unseen data test* yang dipilih random menghasilkan nilai akurasi sebesar 62%.

Kata kunci: analisis sentimen, youtube, support vector machine, naive bayes, akurasi

ABSTRACT

The development of information technology brings modern changes to society. Some of the uses of social media include obtaining and conveying information to the public or relatives, one of the several social media that is currently often used by the public, namely the *Youtube* site. Sentiment Analysis is a technique in which to extract data in the form of text which is used to obtain information about positive and negative sentiments. The scope of this research was carried out only on public comments on the problem of preventing and handling sexual violence in tertiary institutions on the *MataNajwa* channel. Furthermore, the data will be processed using the *Naïve Bayes* method and *Support Vector Machine* using *TF-IDF* feature extraction. The accuracy results obtained in this study were 64% using *Naïve Bayes* and testing was carried out using selected data using *K-Fold Cross Validation* which resulted in an accuracy of 51.7%. While the accuracy value obtained when using the *Support Vector Machine* is 92% by testing using an *unseen data test* that is randomly selected to produce an accuracy value of 62%.

Keywords: Sentiment analysis, youtube, support vector machine, naive bayes, accuration.

1. PENDAHULUAN

Sentimen Analisis merupakan Teknik yang digunakan untuk mengekstrak sebuah data yang berbentuk teks yang digunakan untuk memperoleh sebuah informasi yang bernilai positif dan negatif. Sentimen Analisis tersebut bisa didapatkan dari seseorang yang menggunakan internet untuk menyampaikan suatu komentar yang bersifat menilai atau opini pribadi (Sari, 2019). Melalui penelitian ini, akan dilakukan sebuah penelitian yang bertujuan untuk melakukan analisis sentimen terhadap komentar pengguna *youtube*.

Ruang lingkup penelitian ini dilakukan hanya pada komentar masyarakat terhadap permasalahan tindakan preventif dan penanganan kekerasan seksual pada lingkup perguruan tinggi. Pada tahun 2020, sebanyak 962 laporan yang diterima tentang kasus kekerasan seksual pada perguruan tinggi, selanjutnya pemerintah bersepakat untuk membuat peraturan yang nantinya peraturan tersebut bisa berfungsi untuk melakukan pencegahan dan bisa melakukan penanganan pada kekerasan seksual di perguruan tinggi (tegas.co, 2021).

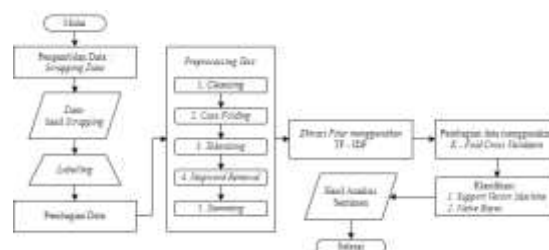
Selanjutnya data tersebut akan dilakukan pengolahan data menggunakan dua metode yaitu metode *Naïve Bayes* dan *Support Vector Machine* dengan ekstraksi fitur TF-IDF. Perlunya melakukan perbandingan antara hasil kerja metode *Naïve Bayes* dan *Support Vector Machine* dikarenakan terdapat beberapa artikel dengan judul diantaranya yaitu “Perbandingan Metode *Naïve Bayes* dan *Support Vector Machine* pada Analisis Sentimen Twitter” dimana mengatakan bahwa hasil kerja metode *Naïve Bayes* didapatkan beberapa nilai yaitu nilai

akurasi, presisi, dan recall dengan nilai yang lebih tinggi dibandingkan menggunakan metode *Support Vector Machine* (Fikri, 2020). Selain itu, pada artikel dengan judul “Komparasi Algoritma *Naive Bayes* Dan *Support Vector Machine* Untuk Analisa Sentimen Review Film” dimana jurnal tersebut mengatakan bahwa nilai akurasi yang didapatkan oleh metode *Naïve Bayes* sebesar 84.50%, dan *Support Vector Machine* sebesar 90% (Indrayuni, 2020).

Dengan perbedaan pendapat dari jurnal yang didapat maka akan dilakukan penelitian hasil kinerja antara metode *Support Vector Machine* dan *Naïve Bayes* guna membuktikan adanya kemungkinan perbedaan hasil nilai rata – rata akurasi. Pada percobaan menggunakan kedua metode tersebut diharapkan mampu mendapatkan nilai yang lebih baik dibandingkan dengan penelitian yang sebelumnya, selain itu melalui penelitian ini dapat diketahui perbandingan kinerja hasil metode mana yang lebih baik antara algoritma *Support Vector Machine* dan *Naïve Bayes* berdasarkan nilai akurasi yang diperoleh.

2. METODOLOGI PENELITIAN

Tahapan penelitian dari proses awal hingga akhir ditunjukkan pada Gambar 1 di bawah ini



Gambar 1. Tahapan penelitian

A. Pengambilan Data

Data

yang digunakan pada penelitian ini

yaitu data komentar Youtube. Dataset tersebut didapatkan dari channel Najwa Shihab dengan judul “Pro Kontra Permendikbud soal Kekerasan Seksual – (Part 4)”. Proses pengumpulan data tersebut dilakukan dengan menggunakan teknik scrapping pada bulan November tahun 2021.

B. Preprocessing Text

Dilakukannya *pre-processing* pada data yang digunakan melalui tahapan (Tuhuteru, 2020) seperti berikut:

1. Cleansing

Pada tahapan ini, semua karakter yang bukan alfabet akan dihapus sehingga mengurangi karakter yang tidak dikehendaki dan tidak memiliki arti dalam analisis sentimen.

2. Case Folding

Pada tahap ini, setiap huruf yang terdapat dalam kalimat akan diubah menjadi *lowercase* atau menjadi huruf kecil semua.

3. Tokenizing

Tahapan ini, komentar yang telah melalui tahapan *cleansing* dan *case folding* akan dipisahkan dari kalimatnya menjadi per kata atau token.

4. Stopword Removal

Tahapan ini ialah proses memilih kata yang penting dari hasil token yaitu kata-kata apa saja yang digunakan dalam dokumen atau disebut dengan *filtering*.

5. Stemming

Proses pada tahapan ini yaitu mengembalikan kata menjadi kata dasar dengan cara membuang awalan, akhiran atau sisipan atau bisa disebut dengan imbuhan. Proses ini dilakukan agar hanya kata dasar yang tersimpan dalam index database dalam sebuah dokumen.

C. TF-IDF

TF-IDF atau singkatan dari *Term Frequency-Inverse Document Frequency* adalah suatu proses untuk melakukan transformasi data dari data tekstual ke dalam data numerik untuk dilakukan pembobotan pada tiap kata atau fitur. Metode ini terkenal efisien, mudan dan akurat (Ria, 2018). Rumus perhitungan TF adalah sebagai berikut (Ria, 2018)

$$tf_{t,d} = \frac{n_{t,d}}{N} \quad (1)$$

Keterangan:

$n_{t,d}$ = nilai istilah yang muncul

N = total dokumen dalam aset

Tf = frekuensi kemunculan kata pada sebuah dokumen

Rumus perhitungan IDF adalah sebagai berikut (Ria, 2018)

$$idf_d = \log\left(\frac{N}{df}\right) \quad (2)$$

Keterangan:

N = total semua dokumen

df = banyak dokumen yang mengandung term tersebut.

Rumus TF-IDF adalah sebagai berikut (Ria, 2018)

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

Keterangan:

TF-IDF = *Term Frequency-Inverse Document Frequency*

TF = nilai TF

IDF = nilai IDF

D. K – Fold Cross Validation

K – Fold Cross-validasi adalah suatu teknik validasi untuk menilai bagaimana hasil statistik analisis akan menggeneralisasi kumpulan data independen. Teknik ini digunakan untuk melakukan prediksi model dan

memperkirakan seberapa akurat sebuah model prediktif ketika dijalankan dalam praktiknya. Salah satu teknik dari validasi silang adalah *k-fold cross validation*, dimana Teknik ini memecah data menjadi k bagian set data dengan ukuran yang sama. Penggunaan *k-fold cross validation* untuk menghilangkan bias pada data. Pelatihan dan pengujian dilakukan sebanyak k kali. (Tempola, 2018).

E. Naïve Bayes

Naive Bayes adalah sebuah algoritma yang terdapat pada *data mining* ialah penggunaan metode *Naïve Bayes* untuk melakukan pemrosesan data memiliki waktu yang cepat, selain pemrosesan yang cepat, metode ini mudah untuk diimplementasikan dengan struktur yang cukup sederhana dan memiliki tingkat efektifitas yang tinggi (Fahriza, 2021). Berikut adalah persamaan teorema bayes sebagai berikut (Raharja, 2018):

$$P(X_i = x_i | Y = y_i) = \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-\frac{(x_i - \mu_{ij})^2}{2\sigma_{ij}^2}} \quad (4)$$

Keterangan:

- P = Peluang
- Xi = Atribut ke-i
- x_i = nilai x ke-i
- Y = kelas yang dicari
- Yi = sub-kelas yang dicari
- σ = standar deviasi
- μ = nilai rata-rata hitung (mean)

F. Support Vector Machine

Support Vector Machine (SVM) merupakan metode yang bisa digunakan untuk menyelesaikan permasalahan secara linier maupun permasalahan non-linier. Dalam menyelesaikan permasalahan non-linier digunakan konsep kernel pada ruang

kerja berdimensi tinggi, dengan mencari hyperplane yang dapat memaksimalkan margin antar kelas data. *Hyperplane* berfungsi dalam memisahkan 2 kelompok *class +1* dan *class -1* dimana setiap *class* memiliki *pattern* masing-masing (Rian, 2020).

Berikut adalah langkah-langkah metode *Support Vector Machine* (Rian, 2020):

a. Menentukan kata pada dokumen dimana kata tersebut sering muncul atau komentar yang digunakan pada *youtube*.

b. Menentukan inialisasi awal dimana nilai α=0.5, C=1, λ=0.5, gamma=0.5 dan epsilon=0.001.

c. Melakukan perhitungan matriks dengan rumus (2.4):

$$D_{ij} = y_i y_j (K(\vec{x}_i \cdot \vec{x}_j) + \lambda^2) \quad (5)$$

Keterangan:

- D_{ij} = anggota matriks data ke-ij
- y_i = label atau kelas data ke-i
- y_j = label atau kelas data ke-j
- λ = turunan batas teoritis
- K($\vec{x}_i \cdot \vec{x}_j$) = fungsi kernel

d. Bagi data ke n= 1,2,3,... n gunakan persamaan (2.5)(2.6)(2.7) berikut ini :

$$E_i = \sum_{j=1}^i a_j D_{ij} \quad (6)$$

$$\delta\alpha_i = \min \{ \max [\gamma (1 - E_i), -\alpha_i], C - \alpha_i \} \quad (7)$$

$$\alpha_i = \alpha_i + \delta\alpha_i \quad (8)$$

Keterangan:

- E_i = nilai error data ke-i
- Γ = tingkat pembelajaran
- max_(i)D_{ij} = nilai maksimum diagonal matriks hessian

e. Melakukan pencari nilai bias (b) menggunakan persamaan (2.8)

$$b = -\frac{1}{2} [w \cdot x^+ + w \cdot x^-] \quad (9)$$

- f. Menguji pada dokumen yang akan diuji
- g. Perhitungan keputusan (10)

$$h(x) = \begin{cases} +1, & \text{if } w \cdot x + b \geq 0 \\ -1, & \text{if } w \cdot x + b < 0 \end{cases}$$

Jika hasil $sign\ h(x)$

dari perhitungan keputusan lebih dari sama dengan 0 maka hasil dari nilai tersebut adalah +1 dan termasuk ke dalam kelas positif dan sebaliknya, jika hasil $sign\ h(x)$ dari perhitungan keputusan kurang dari 0 maka hasil dari nilai -1 dan termasuk ke dalam kelas negatif. Penggunaan persamaan untuk mencari nilai keputusan seperti pada persamaan (11).

$$h(x) = w \cdot x + b \quad (11)$$

atau

$$h(x) = \sum_{i=1}^m a_i y_i K(x, x_i) + b \quad (12)$$

3. HASIL dan PEMBAHASAN

Berikut adalah hasil nilai akurasi yang dihasilkan pada masing-masing metode.

a. Naive Bayes

Berikut adalah nilai akurasi yang didapatkan menggunakan metode *Naive Bayes* yang bisa dilihat pada tabel 1.

Tabel 1 Nilai Akurasi Metode *Naive Bayes*

K-fold cross	Langkah Uji	Naive Bayes
2-fold	Langkah Uji 1	0,488
	Langkah Uji 2	0,512
5-fold	Langkah Uji 1	0,52
	Langkah Uji 2	0,52
	Langkah Uji 3	0,56
	Langkah Uji 4	0,54
	Langkah Uji 5	0,46
10-fold	Langkah Uji 1	0,4
	Langkah Uji 2	0,56
	Langkah Uji 3	0,64
	Langkah Uji 4	0,44
	Langkah Uji 5	0,48
	Langkah Uji 6	0,6
	Langkah Uji 7	0,56

Langkah Uji 8	0,6
Langkah Uji 9	0,52
Langkah Uji 10	0,44

b. Support Vector Machine

Berikut adalah nilai akurasi yang didapatkan menggunakan metode *Support Vector Machine* yang bisa dilihat pada tabel 2.

Tabel 2 Nilai Akurasi Metode *Support Vector Machine*

K-fold cross	Langkah Uji	Support Vector Machine
2-fold	Langkah Uji 1	0,672
	Langkah Uji 2	0,68
5-fold	Langkah Uji 1	0,64
	Langkah Uji 2	0,82
	Langkah Uji 3	0,72
	Langkah Uji 4	0,74
	Langkah Uji 5	0,62
10-fold	Langkah Uji 1	0,64
	Langkah Uji 2	0,6
	Langkah Uji 3	0,92
	Langkah Uji 4	0,72
	Langkah Uji 5	0,68
	Langkah Uji 6	0,84
	Langkah Uji 7	0,72
	Langkah Uji 8	0,72
	Langkah Uji 9	0,6
	Langkah Uji 10	0,6

Berikutnya adalah pengujian data dimana pengujian data ini, masing-masing langkah uji pada setiap fold akan diuji menggunakan 29 *unseen data test* dimana data tersebut telah diambil secara random setelah data dikelompokkan menggunakan metode *K-Fold Cross Validation* dimana komposisi data tersebut adalah 13 positif dan 16 negatif. Berikut adalah tabel dari hasil pengujian data menggunakan *unseen data* baik metode *Naive Bayes* dan *Support Vector Machine*.

Tabel 3 Pengujian *Unseen Data Test Naive Bayes*

K-fold cross	Langkah Uji	Naïve Bayes	
		Data Training	Data Unseen
2-fold	Langkah Uji 1	0,488	0,448
	Langkah Uji 2	0,512	0,551
5-fold	Langkah Uji 1	0,52	0,482
	Langkah Uji 2	0,52	0,551
	Langkah Uji 3	0,56	0,551
	Langkah Uji 4	0,54	0,586
	Langkah Uji 5	0,46	0,551
10-fold	Langkah Uji 1	0,4	0,517
	Langkah Uji 2	0,56	0,517
	Langkah Uji 3	0,64	0,517
	Langkah Uji 4	0,44	0,62
	Langkah Uji 5	0,48	0,551
	Langkah Uji 6	0,6	0,551
	Langkah Uji 7	0,56	0,551
	Langkah Uji 8	0,6	0,551
	Langkah Uji 9	0,52	0,482
	Langkah Uji 10	0,44	0,62

Tabel 4 Pengujian *Unseen Data Test Support Vector Machine*

K-fold cross	Langkah Uji	Support Vector Machine	
		Data Training	Data Unseen
2-fold	Langkah Uji 1	0,672	0,551
	Langkah Uji 2	0,68	0,551
5-fold	Langkah Uji 1	0,64	0,551
	Langkah Uji 2	0,82	0,551
	Langkah Uji 3	0,72	0,62
	Langkah Uji 4	0,74	0,586
	Langkah Uji 5	0,62	0,551
10-fold	Langkah Uji 1	0,64	0,62
	Langkah Uji 2	0,6	0,551
	Langkah Uji 3	0,92	0,62
	Langkah Uji 4	0,72	0,586
	Langkah Uji 5	0,68	0,62

Langkah Uji 6	0,84	0,586
Langkah Uji 7	0,72	0,551
Langkah Uji 8	0,72	0,586
Langkah Uji 9	0,6	0,551
Langkah Uji 10	0,6	0,586

Berdasarkan langkah pengujian data yang telah dilakukan, diperoleh rekapitulasi data hasil akurasi yang disajikan pada tabel berikut:

Tabel 5 Hasil Perhitungan pada metode *Naïve Bayes dan Support Vector Machine*

K-fold cross	Langkah Uji	Naïve Bayes		Support Vector Machine	
		Data Training	Data Unseen	Data Training	Data Unseen
2-fold	Langkah Uji 1	0,488	0,448	0,672	0,551
	Langkah Uji 2	0,512	0,551	0,68	0,551
5-fold	Langkah Uji 1	0,52	0,482	0,64	0,551
	Langkah Uji 2	0,52	0,551	0,82	0,551
	Langkah Uji 3	0,56	0,551	0,72	0,62
	Langkah Uji 4	0,54	0,586	0,74	0,586
	Langkah Uji 5	0,46	0,551	0,62	0,551
10-fold	Langkah Uji 1	0,4	0,517	0,64	0,62
	Langkah Uji 2	0,56	0,517	0,6	0,551
	Langkah Uji 3	0,64	0,517	0,92	0,62
	Langkah Uji 4	0,44	0,62	0,72	0,586
	Langkah Uji 5	0,48	0,551	0,68	0,62
	Langkah Uji 6	0,6	0,551	0,84	0,586
	Langkah Uji 7	0,56	0,551	0,72	0,551
	Langkah Uji 8	0,6	0,551	0,72	0,586
	Langkah Uji 9	0,52	0,482	0,6	0,551
	Langkah Uji 10	0,44	0,62	0,6	0,586

Tabel nilai akurasi diatas menunjukkan perolehan nilai akurasi tertinggi metode naïve bayes yaitu 64% dengan k=10 pada langkah uji ketiga, sedangkan untuk nilai akurasi tertinggi yang dihasilkan oleh metode support vector machine yaitu 92% dengan k=10 pada langkah uji ketiga. Dari hasil akurasi tersebut menunjukkan bahwa nilai yang dihasilkan oleh Support Vector Machine

yaitu berupa nilai akurasi menghasilkan nilai yang lebih besar dari pada nilai akurasi yang dihasilkan oleh metode naive bayes.

4. KESIMPULAN

Pada penelitian ini yakni pada proses *K-fold cross validation* nilai akurasi tertinggi yang dihasilkan oleh metode *naïve bayes* berdasarkan langkah uji ketiga pada $K = 10$ yaitu 64%, sedangkan nilai akurasi tertinggi yang dihasilkan oleh metode *Support Vector Machine* (SVM) berdasarkan langkah uji ketiga pada $K = 10$ yaitu 92%. Untuk masing-masing metode dilakukan pengujian dengan menggunakan *29 unseen data test* menghasilkan nilai akurasi sebesar 51,7% dan 62% pada langkah uji ketiga dengan $K=10$.

Metode *Support Vector Machine* (SVM) pada penelitian ini mampu menghasilkan nilai akurasi yang lebih besar daripada menggunakan metode *Naïve Bayes*, dengan rata-rata nilai akurasi yang diperoleh ialah 70,18% dibandingkan dengan 52%.

Sehingga kedepannya dimungkinkan untuk menggunakan metode yang berbeda agar bisa menghasilkan nilai akurasi yang lebih akurat.

DAFTAR PUSTAKA

- Buntoro, Ghulam Asrofi. 2017. Analisis Sentimen Calon Gubernur DKI Jakarta 2017 Di Twitter. Universitas Muhammadiyah Ponorogo.
- Firman, T., Miftah, M., Amal, K. 2018. Perbandingan Klasifikasi Antara Knn Dan Naive Bayes Pada Penentuan Status Gunung Berapi Dengan K-Fold Cross Validation. Universitas Khairun Ternate
- Hakim, Sulton Nur. 2021. Analisis Sentimen Persepsi Pengguna Myindihome Menggunakan Metode Support Vector Machine (Svm) Dan Naïve Bayes Classifier (Nbc). Universitas Islam Indonesia
- Cahyono, A G. 2016. Pengaruh Media Sosial Terhadap Perubahan Sosial Masyarakat Di Indonesia.
- Sari. 2019. Analisis Sentimen Pelanggan Toko Online Jd.Id Menggunakan Metode Naïve Bayes Classifier Berbasis Konversi Ikon Emosi. Bandung
- Sunoto, Y. 2017. Analisis Testimonial Wisatawan Menggunakan Text Mining Dengan Metode Naive Bayes Dan Decision Tree, Studi Kasus Pada Hotel – Hotel Di Jakarta. Institut Bisnis Dan Informatika Kwik Kian Gie
- Tuhuteru. 2020. Analisis Sentimen Masyarakat Terhadap Pembatasan Sosial Berskala Besar Menggunakan Algoritma Support Vector Machine. Universitas Kristen Indonesia Maluku
- Luqyana, 2018. Analisis Sentimen Cyberbullying Pada Komentar Instagram Dengan Metode Klasifikasi Support Vector Machine. Universitas Brawijaya
- Ria, M. 2018. Penerapan Metode Term Frequency Inverse Document Frequency (Tf-Idf) Dan Cosine Similarity Pada Sistem Temu Kembali Informasi Untuk Mengetahui Syarah Hadits Berbasis Web (Studi Kasus: Syarah Umdatil Ahkam), Universitas Islam Negeri Syarif Hidayatullah Jakarta

- Rian, T. 2020. Analisis Sentimen Terhadap Layanan Indihome Berdasarkan Twitter Dengan Metode Klasifikasi Support Vector Machine (Svm). Universitas Nasional, Jakarta
- Fahriza, F. 2021. Klasifikasi Sentimen Terhadap Gubernur Dki Jakarta Di Media Sosial Twitter Menggunakan Metode Naive Bayes Classifier. Universitas Negeri Sultan Syarif Qasim Riau Pekanbaru
- Fikri, M. 2020. Perbandingan Metode Naive Bayes Dan Support Vector Machine Pada Analisis Sentimen Twitter. Universitas Muhammadiyah Malang