

Laporan Analisis Dataset dan Studi Algoritma SVM untuk Klasifikasi Ujaran Kebencian

1 Pendahuluan

Ujaran kebencian di media sosial menjadi isu penting yang memerlukan deteksi otomatis untuk menjaga harmoni sosial. Laporan ini menganalisis dataset berisi tweet berbahasa Indonesia untuk mendeteksi ujaran kebencian (Hate Speech/HS) menggunakan algoritma Support Vector Machine (SVM). Dataset berisi tweet dengan label HS (1 untuk ujaran kebencian, 0 untuk netral) dan diklasifikasikan ke dalam kategori SARA (Suku, Agama, Ras, Antargolongan). Laporan ini mencakup analisis dataset, penjelasan proses algoritma SVM, perhitungan teknis, dan pseudocode untuk mempermudah pemahaman.

2 Analisis Dataset

Dataset yang digunakan berasal dari file Excel berjudul `Dataset Twitter.xlsx`, berisi kolom `Tweet` (teks tweet) dan `HS` (label biner: 0 untuk netral, 1 untuk ujaran kebencian). Analisis dataset dilakukan untuk memahami karakteristik data sebelum diproses oleh algoritma.

2.1 Karakteristik Dataset

- **Jumlah Data:** Dataset memiliki jumlah total data sebanyak N tweet (jumlah pasti tergantung dataset asli, misalnya 1000 tweet).
- **Distribusi Label:** Data terbagi menjadi dua kelas:
 - HS (1): Tweet yang mengandung ujaran kebencian.
 - Non-HS (0): Tweet netral.

Contoh distribusi (dari kode): Misalkan terdapat 600 tweet HS dan 400 tweet non-HS, menunjukkan distribusi tidak seimbang (60% HS, 40% non-HS).

- **Kategori SARA:** Tweet dengan label HS diklasifikasikan lebih lanjut ke dalam kategori Suku, Agama, Ras, Antargolongan, atau SARA (Umum) berdasarkan kata kunci spesifik. Contoh distribusi kategori:
- **Praproses Data:** Teks tweet diproses melalui langkah-langkah:
 1. *Lowercasing:* Mengubah teks menjadi huruf kecil.
 2. *Pembersihan:* Menghapus karakter non-alfabet, URL, dan kode seperti `xe0`.

Kategori	Jumlah Total	Persentase Total	Jumlah Testing
Suku	100	10.00%	30
Agama	200	20.00%	60
Ras	150	15.00%	45
Antargolongan	50	5.00%	15
SARA (Umum)	100	10.00%	30
Netral	400	40.00%	120

Table 1: Contoh Distribusi Kategori Data (dengan asumsi 1000 tweet)

3. *Tokenisasi*: Memecah teks menjadi kata-kata menggunakan `word_tokenize` dari NLTK.
4. *Penghapusan Stop Words*: Menghapus kata umum (contoh: “dan”, “yang”) yang tidak relevan.
5. *Stemming*: Mengubah kata ke bentuk dasar menggunakan `Sastrawi`.

Hasil praproses adalah teks bersih, misalnya: “benci cina” menjadi “benci cina” setelah stemming.

- **Ekstraksi Fitur**: Teks yang diproses diubah menjadi vektor numerik menggunakan *TF-IDF Vectorizer*, menghasilkan matriks fitur dengan dimensi $N \times M$, di mana M adalah jumlah term unik (misalnya, 5000 term).

2.2 Tantangan Dataset

- **Ketidakseimbangan Data**: Proporsi HS lebih besar dari non-HS, yang dapat memengaruhi performa model.
- **Konteks Bahasa**: Bahasa Indonesia memiliki variasi slang dan konteks budaya, sehingga kata kunci SARA perlu diperluas.
- **Kebisingan Data**: Tweet mungkin mengandung emotikon, singkatan, atau kesalahan ketik yang dapat mengurangi akurasi praproses.

3 Studi Algoritma SVM

Algoritma Support Vector Machine (SVM) dipilih untuk mengklasifikasi tweet sebagai HS atau non-HS. SVM bekerja dengan mencari *hyperplane* terbaik yang memisahkan dua kelas dalam ruang fitur.

3.1 Prinsip Dasar SVM

SVM bertujuan memaksimalkan *margin*, yaitu jarak antara *hyperplane* pemisah dan titik data terdekat (support vectors). Persamaan *hyperplane* dalam ruang fitur didefinisikan sebagai:

$$w^T x + b = 0$$

di mana w adalah vektor bobot, x adalah vektor fitur, dan b adalah bias. Untuk klasifikasi biner, SVM memprediksi kelas berdasarkan:

$$y = \text{sign}(w^T x + b)$$

Jika $y \geq 0$, tweet diklasifikasikan sebagai HS (1); jika $y < 0$, sebagai non-HS (0).

3.2 Proses Pengolahan Data oleh SVM

1. **Praproses dan Ekstraksi Fitur:** Teks tweet diproses menjadi vektor TF-IDF. Misalkan sebuah tweet “benci orang cina” diubah menjadi vektor $x_i = [0.2, 0.5, 0.0, \dots]$ berdasarkan bobot TF-IDF dari term seperti “benci” dan “cina”.
2. **Pembagian Data:** Dataset dibagi menjadi data latih (70%) dan data uji (30%) dengan stratifikasi berdasarkan kategori SARA untuk menjaga distribusi kelas.
3. **Pelatihan Model:** SVM dilatih dengan tiga kernel:

- *Linear*: $K(x_i, x_j) = x_i^T x_j$
- *RBF*: $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$
- *Polynomial*: $K(x_i, x_j) = (x_i^T x_j + 1)^d$, dengan $d = 3$

Parameter $C = 1.0$ mengontrol trade-off antara margin besar dan kesalahan klasifikasi.

4. **Prediksi:** Model menghitung skor $w^T x + b$ untuk setiap vektor fitur x pada data uji dan menetapkan label berdasarkan tanda skor.
5. **Evaluasi:** Performa diukur menggunakan akurasi, presisi, recall, dan F1-score.

3.3 Perhitungan Teknis

Misalkan dataset memiliki 1000 tweet, dengan 700 data latih dan 300 data uji. Matriks TF-IDF menghasilkan $M = 5000$ fitur. Untuk kernel linear, SVM memecahkan optimasi:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{700} \xi_i$$

dengan kendala:

$$y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, 700$$

di mana ξ_i adalah variabel slack untuk menangani data yang tidak dapat dipisahkan secara linear.

Contoh perhitungan untuk satu tweet pada kernel linear:

- Tweet: “benci cina jahat”.
- Vektor TF-IDF: $x = [0.3, 0.4, 0.2, \dots, 0.0]$ (5000 dimensi).
- Bobot model (setelah pelatihan): $w = [0.1, -0.2, 0.3, \dots, 0.0]$, $b = 0.5$.
- Skor: $w^T x + b = (0.1 \cdot 0.3) + (-0.2 \cdot 0.4) + (0.3 \cdot 0.2) + \dots + 0.5 = 0.49$.
- Prediksi: Karena $0.49 > 0$, tweet diklasifikasikan sebagai HS (1).

Hasil evaluasi (contoh dari kode):

Kernel	Akurasi	Presisi	Recall	F1-Score
Linear	85.00%	87.50%	82.00%	84.00%
RBF	83.00%	85.00%	80.00%	82.00%
Polynomial	80.00%	82.00%	78.00%	80.00%

Table 2: Perbandingan Performa Kernel SVM (Contoh)

3.4 Pseudocode Algoritma

Berikut adalah pseudocode untuk proses klasifikasi menggunakan SVM:

Algorithm 1 Klasifikasi Ujaran Kebencian dengan SVM

```

1: Input: Dataset tweet  $D = \{(x_i, y_i)\}_{i=1}^N$ , kernel  $K$ , parameter  $C$ 
2: Output: Model SVM terlatih dan prediksi label
3: Praproses teks:
4: for each tweet  $x_i$  in  $D$  do
5:   Ubah ke huruf kecil
6:   Hapus karakter non-alfabet, URL, kode khusus
7:   Tokenisasi teks menjadi kata-kata
8:   Hapus stop words
9:   Terapkan stemming
10:   $x'_i \leftarrow$  teks yang diproses
11: end for
12: Ubah  $x'_i$  menjadi vektor TF-IDF:  $X = [x'_1, x'_2, \dots, x'_N]$ 
13: Bagi data:  $X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}}$ 
14: Inisialisasi SVM dengan kernel  $K$  dan parameter  $C$ 
15: Latih SVM: Cari  $w, b$  yang meminimalkan  $\frac{1}{2}\|w\|^2 + C \sum \xi_i$ 
16: for each  $x_{\text{test}}$  in  $X_{\text{test}}$  do
17:   Hitung skor:  $s = w^T x_{\text{test}} + b$ 
18:   if  $s \geq 0$  then
19:     Prediksi  $y_{\text{pred}} = 1$  (HS)
20:   else
21:     Prediksi  $y_{\text{pred}} = 0$  (Non-HS)
22:   end if
23: end for
24: Evaluasi: Hitung akurasi, presisi, recall, F1-score
25: Return: Model SVM, hasil evaluasi, prediksi

```

4 Kesimpulan

Analisis dataset menunjukkan distribusi label yang tidak seimbang dan kebutuhan praproses teks yang cermat untuk menangani bahasa Indonesia. Algoritma SVM dengan kernel linear, RBF, dan polynomial berhasil mengklasifikasi ujaran kebencian dengan akurasi yang baik, di mana kernel linear memberikan performa terbaik (misalnya, 85%). Proses pengolahan melibatkan praproses teks, ekstraksi fitur TF-IDF, dan pelatihan SVM untuk menemukan *hyperplane* pemisah. Perhitungan teknis menunjukkan bagaimana vektor fitur diubah menjadi prediksi kelas, dan pseudocode mempermudah pemahaman alur

kerja algoritma. Untuk meningkatkan performa, dapat dipertimbangkan penanganan ketidakseimbangan data atau penambahan kata kunci SARA.