

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/320371563>

# Text Mining Untuk Klasifikasi Kategori Cerita Pendek Menggunakan Naïve Bayes (NB)

Article · August 2017

CITATIONS

0

READS

132

1 author:



**Oman Somantri**

Politeknik Harapan Bersama, Tegal, Indonesia

15 PUBLICATIONS 2 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Information System [View project](#)



Power System [View project](#)

# Text Mining Untuk Klasifikasi Kategori Cerita Pendek Menggunakan *Naïve Bayes* (NB)

Oman Somantri<sup>#1</sup>

<sup>#</sup>Program Studi Teknik Informatika, Politeknik Harapan Bersama Tegal  
Jln. Mataram No.09 Pesurungan Lor Kota Tegal, Indonesia

<sup>1</sup>oman.somantri@poltektegal.ac.id

**Abstract**— *Determination of the category of a short story requires a slightly long process, in other way we must read a whole or at least a half of the contents of the short story to know the entire contents from the beginning to the end. These constraints require a solution to overcome by using Naïve Bayes algorithm (NB) to serve as the solution of the existing problems. Naïve Bayes, used as a model, resulted with accuracy of 78.59%. Evaluation was conducted by comparing the level of accuracy produced with other models of Support Vector Machine (SVM). The result of the research show that level of accuracy NB greater than Support Vector Machine (SVM) with accuracy level 64,36%. Based on the results of research conducted can be concluded that Naïve Bayes has a higher level of accuracy than the Support Vector Machine (SVM) for the short story category classification.*

**Keywords**— *Naïve Bayes, Support Vector Machine, short story, Model*

**Abstrak**— Penentuan kategori sebuah cerita pendek memerlukan sebuah proses yang sedikit lama, dimana kita harus membaca secara keseluruhan atau minimal setengah dari isi dari cerpen tersebut karena untuk dapat mengetahui seluruh isi konten dari cerpen tersebut adalah dengan cara membaca isi cerpen mulai dari awal sampai akhir bacaan cerpen. Kendala tersebut memerlukan sebuah solusi untuk mengatasinya, maka pada penelitian ini diusulkan sebuah model dengan menggunakan algoritme *Naïve bayes* (NB) untuk dijadikan sebagai solusi dari permasalahan yang ada. *Naïve Bayes* digunakan sebagai model dan menghasilkan tingkat akurasi sebesar 78,59%. Evaluasi dilakukan dengan membandingkan tingkat akurasi yang dihasilkan dengan model lain yaitu *Support Vector Machine* (SVM). Dari Hasil penelitian memperlihatkan bahwa tingkat akurasi NB lebih besar dibandingkan dengan *Support Vector Machine* (SVM) dengan tingkat akurasi 64,36%. Berdasarkan hasil penelitian yang dilakukan dapat disimpulkan bahwa *Naïve Bayes* mempunyai tingkat akurasi lebih tinggi dibandingkan dengan *Support Vector Machine* (SVM) untuk klasifikasi kategori cerpen.

**Kata Kunci**— *Naïve Bayes, Support Vector Machine, Cerpen, Model*

## I. PENDAHULUAN

Sebagai salah satu bagian dari kebudayaan Indonesia, cerpen merupakan karya sastra yang paling banyak diminati oleh banyak orang. Sebuah cerpen akan dapat diminati orang apabila isi dari cerpen tersebut menarik, dan dapat membawa orang yang membacanya hanyut ke dalam isi dari cerita tersebut. Berbagai macam latar belakang pembaca cerpen saat

ini, mulai dari remaja, anak-anak, dewasa, maupun para orang tua. Perbedaan latar belakang inilah tentunya menjadikan sebuah cerpen memiliki segmentasi yang berbeda sesuai dengan karakteristik pembacanya yang menyesuaikan dengan usia dan latar belakang dari pembaca sehingga cerpen memiliki banyak kategori sesuai dengan isi dari cerpen tersebut seperti kategori cerpen anak, dongeng, fiksi, pendidikan, dewasa, romantis, dan lainnya. Cerpen adalah cerita fiktif yang belum pasti kebenarannya serta ceritanya relatif pendek dan cerpen bukanlah suatu analisis argumentatif [1].

Untuk dapat menentukan sebuah cerpen masuk kedalam kategori cerpen tertentu bukanlah hal yang mudah, sudah tentu orang harus membaca keseluruhan atau minimal sebagian isi dari cerpen tersebut kemudian barulah dapat diketahui cerpen tersebut masuk kedalam kategori apa. Hal ini yang menjadikan kesulitan dalam menentukan sebuah cerpen masuk kedalam kategori tertentu, sedangkan terkadang banyak orang yang tidak bisa membaca terlebih dahulu isi dari cerpen tersebut. Permasalahan kadang terjadi banyak para orang tua yang ingin memberikan sebuah cerita cerpen kepada anaknya akan tetapi karena belum diketahui cerpen tersebut masuk kedalam kategori apa, terkadang isi cerpen tidak sesuai dengan umur usia anak, ini merupakan salah satu contoh kasus yang sering terjadi. Berdasarkan permasalahan tersebut maka perlu sebuah solusi yang dapat mengatasinya sehingga dapat dijadikan sebagai pendukung keputusan dalam menentukan kategori sebuah cerpen.

Dalam bidang komputerisasi yang termasuk kedalam *machine learning*, *Naïve Bayes* dan *Support Vector Machine* (SVM) merupakan metode yang digunakan untuk klasifikasi teks dalam *text mining*. Sebagai salah satu metode komputasi yang efisien dan mempunyai *performance predictive* yang baik, *naïve bayes* merupakan salah satu metode klasifikasi teks yang populer [2]. *Naïve Bayes* merupakan algoritme yang sering digunakan dalam pengkategorian teks, dimana konsep dasarnya adalah menggabungkan probabilitas kata-kata dan kategori sebuah dokumen [3] [8].

Penelitian terkait dengan klasifikasi teks dengan menggunakan *naïve bayes* sudah dilakukan oleh peneliti sebelumnya, seperti yang dilakukan oleh Nurul. S.A, (2016) melakukan penelitian untuk membandingkan *Naïve Bayes* dan *Support Vector Machine* (SVM) untuk klasifikasi emosi pada teks bahasa Indonesia [4]. Hamzah. A, (2012) melakukan penelitian klasifikasi teks dengan *naïve bayes classifier* (NBC) untuk pengelompokan teks berita dan abstrak akademis [5].

Selanjutnya Winarsih, N. A. S., & Supriyanto, C. (2016) meneliti untuk mengevaluasi metode klasifikasi deteksi emosi pada teks Indonesia [6]. Sedikit berbeda dengan yang dilakukan oleh Jamal, N., dkk. (2012) meneliti klasifikasi puisi dengan menggunakan *Support Vector Machine* (SVM) [7].

Dari semua penelitian yang telah dilakukan berbeda dengan penelitian yang sudah ada sebelumnya, perbedaan pada penelitian ini adalah pada proses preprosesing data dan metode yang digunakan untuk klasifikasi kategori cerpen. Berdasarkan dari kelebihan yang dimiliki maka pada penelitian ini mengusulkan *Naïve Bayes* sebagai metode yang diusulkan dan diterapkan untuk pengklasifikasian jenis kategori cerita pendek sehingga didapatkan sebuah model yang tepat untuk menghasilkan tingkat akurasi yang terbaik untuk klasifikasi kategori cerita pendek.

## II. NAÏVE BAYES

### A. Naïve Bayes (NB)

*Naive Bayes* merupakan salah satu metode yang banyak digunakan berdasarkan beberapa sifatnya yang sederhana, metode ini mengklasifikasikan data berdasarkan probabilitas  $P$  atribut  $x$  dari setiap kelas  $y$  data [8]. *Naive Bayes* adalah metode yang digunakan dalam statistika untuk menghitung peluang dari suatu hipotesis, *naive bayes* menghitung peluang suatu kelas berdasarkan pada atribut yang dimiliki dan menentukan kelas yang memiliki probabilitas paling tinggi. *Naive Bayes* mengklasifikasikan kelas berdasarkan pada probabilitas sederhana dengan mengasumsikan bahwa setiap atribut dalam data tersebut bersifat saling terpisah. Pada model probabilitas setiap kelas  $k$  dan jumlah atribut  $a$  yang dapat dituliskan seperti persamaan dibawah ini:

$$P = (y_1 | x_1, x_2, \dots, x_a) \quad (1)$$

Perhitungan *naive bayes* yaitu probabilitas dari kemunculan dokumen  $X_a$  pada kategori kelas  $Y_k$   $P(x_a/y_k)$ , dikali dengan probabilitas kategori kelas  $P(y_k)$ . Dari hasil kali tersebut kemudian dilakukan pembagian terhadap probabilitas kemunculan dokumen  $P(x_a)$ . Sehingga didapatkan rumus perhitungan *Naive Bayes* dituliskan pada persamaan:

$$P(y_k | x_a) = \frac{P(y_k)P(x_a|y_k)}{P(x_a)} \quad (2)$$

Kemudian dilakukan proses pemilihan kelas yang optimal maka dipilih nilai peluang terbesar dari setiap probabilitas kelas yang ada. Maka didapatkan rumus untuk memilih nilai terbesar seperti pada persamaan berikut:

$$y(x_i) = \arg \max P(y) \prod_{i=1}^a P(x_i | y) \quad (3)$$

Pembobotan suatu atribut kelas dapat meningkatkan pengaruh prediksi. Dengan memperhitungkan bobot atribut terhadap kelas maka yang menjadi dasar ketepatan klasifikasi bukan hanya probabilitas melainkan juga pada bobot setiap atribut kelas.

## III. METODOLOGI PENELITIAN

### A. Dataset Penelitian

*Dataset* yang digunakan dalam penelitian ini diambil dari [www.cerpenmu.com](http://www.cerpenmu.com). Data *online* ini adalah berupa teks yang berbentuk cerita pendek yang sudah ditentukan kategorinya, yaitu kategori cerpen anak dan kategori cerpen dongeng. *Dataset* adalah data yang dibuat antara tahun 2015 sampai dengan 2016 dengan jumlah data sebanyak 121 cerpen.

### B. Preprocessing Data

Sebelum *dataset* dimasukan kedalam model yang diusulkan, terlebih dahulu dilakukan preprosesing data. Pada tahapan ini dilakukan beberapa hal, diantaranya adalah *tokenized*, *transform cases*, *filter tokens*, *filter stopword* dan *Stem* [9].

1) *Tokenized*: merupakan proses untuk memisah-misahkan kata. Hasil dari pemisahan tersebut dinamakan token.

2) *Transform cases*: Merupakan proses untuk merubah bentuk kata-kata, pada proses ini karakter dijadikan menjadi huruf kecil atau *lower case* semua.

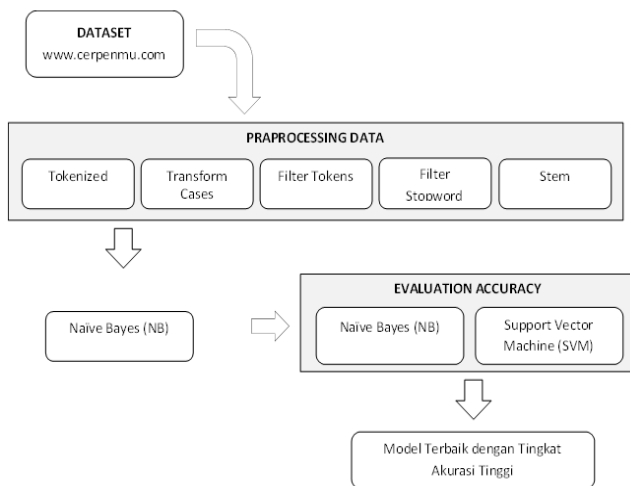
3) *Filter tokens*: Proses pengambilan kata-kata yang penting dari token yang sudah dihasilkan berdasarkan jumlah karakter. Pada proses ini parameter yang digunakan adalah *min chars* = 3, dan *max chars* = 25.

4) *filter stopword*: Proses menghilangkan kata-kata yang sering muncul namun tidak memiliki pengaruh apapun dalam ekstraksi klasifikasi teks. Pada proses ini kata yang termasuk adalah seperti penunjuk waktu, kata tanya dan kata sambung.

5) *Stem*: Merupakan proses pengubahan bentuk kata menjadi kata dasar. Metode ini merupakan proses pengubahan bentuk kata menjadi kata dasar yang menyesuaikan dengan strukturb yang digunakan dalam proses *stemming*.

### C. Model Yang diusulkan

Pada penelitian ini model yang diusulkan adalah *Naive Bayes* sebagai algoritme pembelajaran. Untuk mendapatkan tingkat akurasi yang sesuai, proses validasi dilakukan dengan menggunakan *K-Fold Cross Validation* dengan harapan hasil validasi eksperimen dapat menghasilkan hasil yang terbaik. Model yang diusulkan kemudian di evaluasi dengan cara membandingkan hasil tingkat akurasi yang didapatkan dengan model yang lain yaitu *Support Vector Machine* (SVM).



Gambar. 1 Model Yang Diusulkan

#### D. Validasi

Evaluasi model pada tahapan ini menggunakan evaluasi *matrixs Confusion* [10], seperti pada tabel 1 dibawah:

TABEL I  
CONFUSION MATRIXS

	Hasil Prediksi	
	Positive	Negative
Positive	True Positive	False Positive
Negative	False Negative	True Negative

Untuk menghitung tingkat akurasi digunakan persamaan sebagai berikut:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (4)$$

Keterangan:

- True Positive (TP)
- False Positive (FP)
- False Negative (FN)
- True Negative (TN)

#### IV. HASIL DAN PEMBAHASAN

Penelitian menggunakan *tools* Rapid Miner 5.3 untuk analisis data, dan komputer dengan spesifikasi CPU Intel Core i5 2,67 GHz, memori RAM 4 GB, sistem operasi Windows 7 profesional SP1 32-bit.

##### A. Hasil Ekperimen Naïve Bayes

Pada Eksperimen terhadap model yang digunakan, eksperimen dilakukan dengan menggunakan cerpen yang sudah ditetapkan sesuai dengan kategori cerpen yaitu cerpen anak dan cerpen dongeng dengan jumlah dataset sebanyak 121 cerpen.

Berikut ini adalah contoh konten isi cerpen yang digunakan sebagai salah satu dataset penelitian, sebagai contoh cerpen yang termasuk dalam kategori cerpen dongeng.

Hai namaku Sabila Salwa Putri Wahyuhadi cukup dipanggil Salwa

"Salwa bangun! cepat bangun terus shalat subuh" bentak bunda. Aku bangun dan shalat subuh sehabis shalat subuh aku tidur. Dan paginya aku siap siap mandi, makan lalu berangkat sekolah. Ini hari pertamaku masuk sekolah kelas 3. aku masuk dan perkenalan setelah itu pelajaran lalu istirahat.

Pas istirahat, aku kenalan sama teman-teman yang lain yang paling aku sukai adalah aca, gadis kecil memakai kerudung dengan pipi yang menggelembung.

Sudah 2 minggu aku sekolah disana rasanya sangat menyenangkan apalagi aca, dia sudah menjadi sahabatku sejak 6 hari yang lalu.

"Heh ca kita main bareng yuk!" ajakku pada aca. Tapi, dia tidak menjawab pertanyaanku dan murung dia sangat sedih. Aku juga tidak tahu mengapa dia sedih. Dan aca pun langsung pergi. Aku mengikutinya dan ternyata dia mengkhianatiku. Dia bersahabat dengan silvia ayu musuh terbesarku.

Aku sedih dan langsung pulang ke rumah dan akhirnya aku memutuskan untuk tidak bersahabat lagi dengan aca. Walaupun sangat berat.

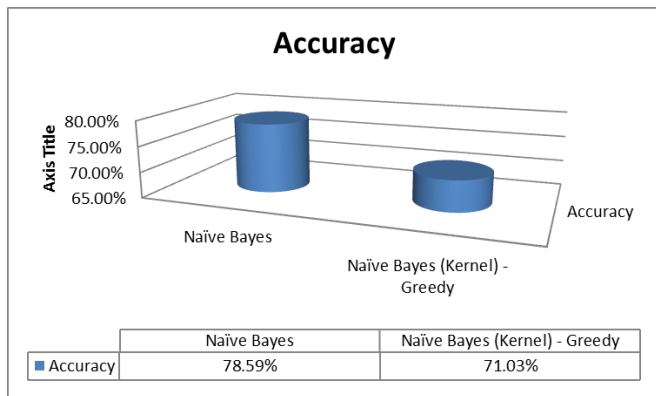
Hasil ekperimen terhadap dataset yang sudah didapatkan dengan menggunakan model *Naive Bayes* dapat diperlihatkan seperti pada tabel 2 berikut:

TABEL II  
HASIL EKPERIMEN MODEL NAIVE BAYES

No	Hasil Ekperimen	
	Model	Accuracy
1	Naïve Bayes	78.59%
2	Naïve Bayes (Kernel) - Greedy	71.03%

Pada tabel 2 diatas diperlihatkan bahwa tingkat akurasi dari hasil ekperimen menunjukkan bahwa *naïve bayes* menghasilkan tingkat akurasi sebesar 78,59%. Terjadinya kesalahan dalam klasifikasi kategori cerpen sehingga mengakibatkan tingkat akurasi yang dihasilkan menjadi kecil hal ini disebabkan oleh model yang diusulkan masih belum sesuai dengan yang diinginkan, hal ini terjadi akibat dari berbagai aspek seperti perbedaaan pada model yang digunakan, proses praprosesing data, *setting* parameter model yang digunakan, dan berbagai aspek lainnya yang dianggap mempengaruhi pada tingkat akurasi klisifikasi yang dihasilkan.

Gambar 2 memperlihatkan grafik hasil eksperimen yang dihasilkan dengan menggunakan model algoritme *Naive Bayes*.



Gambar. 2 Hasil eksperimen Naïve Bayes

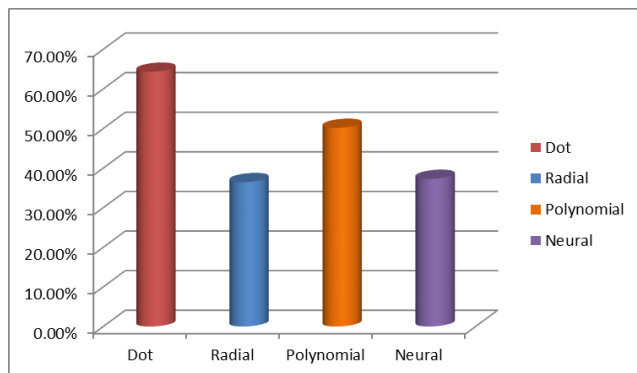
### B. Hasil Eksperimen Support Vector Machine

Pada model *Support Vector Machine* (SVM), eksperimen dilakukan sehingga menghasilkan hasil dapat diperlihatkan seperti pada tabel 3 berikut:

TABEL III  
HASIL EKPERIMEN MODEL SUPPORT VECTOR MACHINE (SVM)

No	Hasil Ekperimen SVM	
	Type Kernel	Accuracy
1	Dot	64.36%
2	Radial	36.47%
3	Polynomial	50.19%
4	Neural	37.31%

Apabila kita melihat tabel 3, terlihat bahwa tingkat akurasi yang paling tinggi SVM adalah dengan menggunakan *kernel type dot* yaitu sebesar 64,36%. Apabila dibuatkan pada grafik maka akan tampak seperti pada gambar 3.

Gambar. 3 Hasil eksperimen *Support Vector Machine* (SVM)

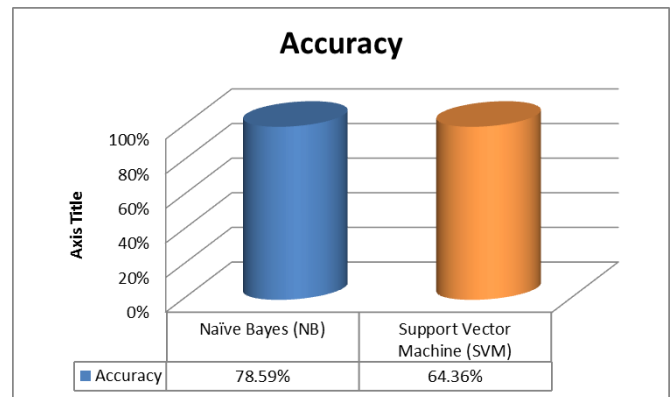
### C. Evaluasi Model

Setelah didapatkannya hasil eksperimen yang dilakukan kedalam model *Naïve Bayes* dan *Support Vector Machine*, maka didapatkan hasilnya seperti pada tabel 4 dibawah:

TABEL IV  
HASIL EKPERIMEN NAÏVE BAYES & SVM

No	Hasil Ekperimen Optimasi	
	Model	Accuracy
1	Naïve Bayes (NB)	78.59%
3	Support Vector Machine (SVM)	64.36%

Berdasarkan hasil yang diperoleh, terlihat bahwa *Naïve bayes* menghasilkan tingkat akurasi sebesar 78,59%, sedangkan *Support Vector Machine* (SVM) menghasilkan tingkat akurasi sebesar 64,34%.



Gambar. 4 Hasil eksperimen Naïve Bayes dan SVM

## V. KESIMPULAN

Hasil eksperimen yang dilakukan terhadap model yang diusulkan yaitu *Naïve Bayes*, menghasilkan tingkat akurasi sebesar 78,59% lebih besar dibandingkan dengan model lain yaitu *Support Vector Machine* (SVM). Dari hasil penelitian dapat disimpulkan bahwa *Naïve Bayes* memiliki tingkat akurasi yang lebih baik dibandingkan dengan SVM untuk pengklasifikasian kategori cerita pendek. Meskipun demikian bahwa hasil yang didapatkan belumlah sempurna sehingga perlu adanya sebuah optimasi dalam model tersebut. Maka saran untuk penelitian selanjutnya adalah melakukan eksperimen dengan model-model lain dan melakukan *setting* parameter yang lebih tepat sehingga didapatkan tingkat akurasi yang lebih baik.

## DAFTAR REFERENSI

- [1] Sumardjo, Jacob dan Saini K.M. 1988. Apresiasi Kesusastraan. Jakarta: PT.Gramedia.
- [2] Chen, J., Huang, H., Tian, S., & Qu, Y. (2009). Feature selection for text classification with Naïve Bayes. *Expert Systems with Applications*, 36, 5432–5435.
- [3] Zhang, & Gao, F. (2011). An Improvement to NB for Text Classification. *Procedia Engineering*, 15, 2160–2164.
- [4] Nurul, S. A. (2016). Perbandingan Metode Naïve Bayes Dan Support Vector Machine Untuk Klasifikasi Emosi Pada Teks Bahasa Indonesia. *Skripsi, Fakultas Ilmu Komputer*.
- [5] Hamzah, A. (2012). Klasifikasi teks dengan naïve bayes classifier (nbc) untuk pengelompokan teks berita dan abstract akademis. In *Prosiding Seminar Nasional*.
- [6] Winarsih, N. A. S., & Supriyanto, C. (2016). Evaluation of classification methods for Indonesian text emotion detection. In *Technology of Information and Communication (ISemantic), International Seminar on Application for* (pp. 130-133). IEEE.

- [7] Jamal, N., Mohd, M., & Noah, S. A. (2012). Poetry classification using support vector machines. *Journal of Computer Science*, 8(9), 1441.
- [8] McCallum, A. & Nigam, K. (1998). A comparison of event models for naïve Bayes text classification. In AAAI-98 workshop on learning for text categorization.
- [9] Aggarwal, C. C., & Zhai, C. (Eds.). (2012). *Mining text data*. Springer Science & Business Media.
- [10] Davis, J., & Goadrich, M. (2006, June). The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning* (pp. 233-240). ACM.

**Oman Somantri**, lahir pada tahun 1985 Sumedang, menerima gelar Sarjana Komputer (S.Kom) dari STMIK Sumedang jurusan Teknik Informatika pada tahun 2011, dan gelar Magister Komputer (M.Kom) dari Universitas Dian Nuswantoro (UDINUS) jurusan Teknik Informatika pada tahun 2015. Saat ini mengajar sebagai dosen di Politeknik Harapan Bersama Tegal. Minat penelitian adalah *Intelligent System, Machine Learning, Data Mining* dan *Text Mining*.