

Implementasi *Machine Learning* pada Aplikasi Pendeteksi Konten Pornografi

Asep Abdul Sofyan¹, Siti Maisaroh², Iwan Sumpena³

^{1,2,3}Teknik Informatika, STMIK Bina Sarana Global, Tangerang

Email: ¹asepabdulsofyan@stmikglobal.ac.id, ²maekayla27@gmail.com, ³sumpenaiwan@gmail.com

Salah satu perkembangan teknologi adalah keberadaan Internet. Internet sudah menjadi bagian tak terpisahkan dari kehidupan manusia modern saat ini. Tidak bisa dipungkiri, selain hal positif yang bisa didapat dari internet ada pula dampak negatif yang ditimbulkan yaitu sumber informasi yang terkandung di internet terkadang mengandung unsur negative, salah satunya pornografi. Kecerdasan yang sedang berkembang bisa dimanfaatkan untuk hal ini, yaitu suatu mesin atau sistem yang dapat mengklasifikasi apakah suatu situs mengandung konten pornografi atau tidak. Banyak metode yang dapat dipakai dalam masalah ini, salah satunya adalah dengan mengklasifikasi text yang ada dalam situs tersebut dan dianalisa atau diklasifikasikan. Cara mengambil data nya sendiri bisa memakai teknik scraping, yaitu teknik yang dapat mengambil konten text format html dari suatu situs. Teknik ini bisa digunakan dalam mendeteksi apakah suatu situs mengandung konten pornografi atau tidak. Untuk proses klasifikasinya sendiri terdapat banyak algoritma yang bisa digunakan salah satunya adalah algoritma *Naive Bayes*. Berdasarkan latar belakang diatas, penulis mengajukan penelitian berjudul “Implementasi *Machine Learning* pada aplikasi pendeteksi konten pornografi”. Diharapkan dengan adanya penelitian ini mampu memperoleh sebuah aplikasi yang dapat mendeteksi konten pornografi pada suatu situs tanpa perlu membuka halaman situs tersebut. Dengan begitu situs yang terdeteksi mengandung konten pornografi bisa ditindak lebih lanjut supaya tidak bisa di akses.

Kata kunci: *Machine Learning*, Pornografi, *Naive Bayes*, Klasifikasi, Situs.

Abstract- *One technological development is the existence of the Internet. The internet has become an inseparable part of modern human life today. It is undeniable, besides the positive things that can be obtained from the internet, there are also negative impacts caused by the source of information contained on the internet sometimes containing negative elements, one of which is pornography. Developing intelligence can be utilized for this, namely a machine or system that can classify whether a site contains pornographic content or not. Many methods can be used in this problem, one of which is to classify the text that is on the site and analyzed or classified. How to retrieve the data yourself can use scraping techniques, which are techniques that can retrieve HTML format text content from a site. This technique can be used to detect whether a site contains pornographic content or not. For the classification process itself, many algorithms can be used, one of which is the Naive Bayes algorithm. Based on the above background, the authors propose a study entitled*

"Implementation of Machine Learning in pornographic content detection applications". It is expected that with this research being able to obtain an application that can detect pornographic content on a site without the need to open the site's pages. That way the site that is detected to contain pornographic content can be followed up so that it cannot be accessed.

Keywords: *Machine Learning, Pornography, Naive Bayes, Classification, Site.*

I. PENDAHULUAN

Kemajuan teknologi yang berkembang pesat telah banyak membantu manusia menyelesaikan banyak pekerjaan dengan lebih cepat dan efisien. Salah satu perkembangan teknologi adalah keberadaan Internet. Internet sudah menjadi bagian tak terpisahkan dari kehidupan manusia modern saat ini. Tentu hal ini merupakan dampak positif dari kemajuan teknologi. Namun tidak bisa dipungkiri, ada pula dampak negatif yang ditimbulkan yaitu sumber informasi yang terkandung di internet terkadang mengandung unsur negative, salah satunya pornografi.

Di negara kita pornografi merupakan sesuatu yang ilegal karena sifatnya yang merusak bagi setiap orang yang mengaksesnya, terutama kalangan remaja dan anak-anak. Karena itu di Indonesia sudah dibentuk Undang-Undang yang mengatur tentang pornografi yaitu UU No. 44 Tahun 2008 Tentang Pornografi (UU Pornografi). Sementara itu juga melalui Kementerian Komunikasi dan Informatika pemerintah sudah berusaha mencegah penyebaran konten negatif khususnya konten pornografi di Indonesia melalui internet, salah satunya dengan membuat Layanan DNS Nawala.

Perkembangan Artificial intelligence (AI)^[2] saat ini mengubah peradaban manusia dimana sekarang banyak pekerjaan yang awalnya dikerjakan oleh manusia kini sudah di ambil alih oleh sebuah sistem atau machine. Semakin canggih sebuah teknologi akan semakin mempermudah suatu pekerjaan. Kecerdasan Buatan sangat bermanfaat untuk mengembangkan metode dan sistem untuk menyelesaikan suatu masalah yang bisa diselesaikan oleh manusia. Misalnya pencarian tempat, bidang bisnis^[3], rumah tangga dan dapat meningkatkan kinerja sistem informasi yang berbasis komputer^[1].

II. METODE PENELITIAN

Adapun tahapan yang dilakukan penulis dalam pelaksanaan penelitian ini adalah sebagai berikut:

1. Studi Literatur Pada tahap ini peneliti mengumpulkan dan mempelajari bahan referensi seperti dokumen dan literatur yang berkaitan dengan topik penelitian. Literatur yang dijadikan bahan referensi dapat berupa buku, skripsi, jurnal, artikel, dan beberapa sumber lainnya yang diperoleh dari internet.
2. Analisis Masalah Pada tahap ini penulis melakukan analisis terhadap literatur yang sebelumnya telah dikumpulkan untuk memperoleh pemahaman mengenai metode yang dilakukan pada penelitian ini dan masalah yang ingin diselesaikan dalam penelitian ini.
3. Implementasi pada tahap ini dilakukan proses implementasi algoritma Naive Bayes dalam aplikasi komputer menggunakan bahasa pemrograman Asp.net core.
4. Pengujian Pada tahap ini dilakukan pengujian data yang telah ada untuk memastikan bahwa implementasi algoritma Naive Bayes dalam mengidentifikasi konten pornografi memberikan hasil yang terbaik.
5. Dokumentasi dan Penyusunan Laporan Pada tahap ini dilakukan dokumentasi hasil dan penyusunan laporan hasil dari analisis dan implementasi dari penelitian yang dilakukan dalam bentuk jurnal.

Berikut adalah gambaran dan detail -detail mengenai Klasifikasi Text:

1. Text Preprocessing

Pengklasifikasian artikel berita secara otomatis bisa dikategorikan sebagai text mining. Proses text mining^[4] dibagi menjadi 3 tahap utama, yaitu proses awal terhadap teks (*text preprocessing*), transformasi teks ke dalam bentuk antara (*text transformation/feature generation*), dan penemuan pola (*pattern discovery*)^[5]. Masukan awal dari proses ini adalah suatu data teks dan keluarannya berupa pola sebagai hasil interpretasi.

2. Text Transformation/Feature Generation

Tahapan ini bertujuan untuk mempersiapkan teks menjadi data yang akan mengalami pengolahan pada tahapan berikutnya. Tindakan yang dilakukan meliputi tindakan kompleks dan tindakan sederhana. Contoh tindakan yang bersifat kompleks pada tahap ini adalah *part-of-speech (pos) tagging*, membangkitkan parse tree^[6]. Contoh tindakan yang bersifat sederhana adalah proses parsing sederhana terhadap teks, yaitu memecah suatu kalimat menjadi sekumpulan kata. Selain itu pada tahapan ini biasanya juga dilakukan *case folding*, yaitu pengubahan karakter huruf menjadi huruf kecil^[7]

3. Stemming Bahasa Indonesia

Dalam bahasa Indonesia, afiks/imbuhan terdiri dari sufiks (akhiran), infiks (sisipan) dan prefiks (awalan). Pada penelitian ini proses *stemming* yang dibangun^[8] hanya menangani kata yang mengalami penambahan prefiks dan sufiks. Hal ini dilakukan karena proses penambahan infiks dalam bahasa Indonesia jarang

terjadi sehingga tidak ada pengaruh yang signifikan terhadap akurasi sistem. Selain itu, penanganan kata yang mengandung infiks relative sulit dan membebani waktu komputasi sistem^[9].

Terdapat 5 aturan tahap dalam proses stemming pada bahasa Indonesia sebagai berikut:

1. Penanganan terhadap partikel infleksional. yaitu : lah, kah dan tah. Contoh : duduklah, apakah.
2. Penanganan terhadap kata ganti infleksional, yaitu : ku, mu, nya. Contoh : sepedamu, mobilnya.
3. Penanganan terhadap prefiks derivasional pertama, yaitu : meng dan semua variasinya, peng dan semua variasinya, di, ter, dan ke. Contoh : membakar, pegukur, kekasih
4. Penanganan terhadap prefiks derivasional kedua, yaitu : ber dan semua variasinya serta per dan semua variasinya. Contoh : berlari, belajar, perjelas.
5. Penanganan terhadap sufiks derivasional, yaitu: kan, an, i. Contoh : makanan, gantikan, tandai.

Karena struktur morfologi dalam bahasa Indonesia yang rumit, maka kelima tahap aturan diatas tidak cukup untuk menangani proses stemming bahasa Indonesia. Kesulitan dalam membedakan suatu kata yang mengandung imbuhan baik prefiks maupun sufiks dengan suatu kata dasar yang salah satu suku katanya merupakan bagian dari imbuhan, terutama dengan kata dasar yang mempunyai sukukata lebih besar dari dua^[10].

4. Pattern Discovery

Tahap penemuan pola atau pattern discovery adalah tahap terpenting dari seluruh proses text mining. Tahap ini berusaha menemukan pola atau pengetahuan dari keseluruhan teks. Terdapat dua teknik pembelajaran pada tahap pattern discovery ini, yaitu *unsupervised* dan *supervised learning*^[11]. *Supervised learning* melakukan klasifikasi suatu data baru berdasarkan data latih. *Unsupervised learning* data latih dikelompokkan berdasarkan ukuran kemiripan pada suatu kelas^[12].

5. Naïve Bayes Classifier

Naïve bayes classifier termasuk ke dalam algoritma pembelajaran bayes. Algoritma pembelajaran bayes menghitung probabilitas eksplisit untuk menggambarkan hipotesa yang dicari. Suatu data pada *naïve bayes classifier* direpresentasikan dengan konjungsi dari nilai-nilai atribut dan sebuah fungsi target $f(x)$ yang dapat memiliki nilai apapun dari himpunan set domain V ^[4]. Sistem dilatih menggunakan data latih lengkap berupa pasangan nilai-nilai atribut dan nilai target kemudian sistem akan diberikan sebuah data baru dalam bentuk $\langle a_1, a_2, a_3, \dots, a_n \rangle$ dan sistem diberi tugas untuk menebak nilai fungsi target dari data tersebut^[13].

III. HASIL DAN PEMBAHASAN

Penelitian tentang deteksi konten negatif khususnya konten pornografi sudah pernah dilakukan dan dikembangkan sebelumnya. Salah satunya oleh Rendra Mahardika dari Universitas Sumatra Utara dengan judul "Identifikasi Konten Pornografi Berbahasa Indonesia Menggunakan Algoritma Support Vector Machine (SVM)". Pada penelitian tersebut Rendra Mahardika membuat aplikasi menggunakan python 3.7.0 dan algoritma Support Vector Machine dalam pembuatannya.

Pada penelitian yang dilakukan Rendra Mahardika batasan ruang lingkup permasalahan yang diteliti adalah :

1. Data latih dan data uji merupakan teks dari isi laman web berbahasa Indonesia.
2. Identifikasi ditujukan terhadap laman web yang dinilai sebagai konten pornografi.
3. Konten pornografi yang akan diidentifikasi adalah konten yang berbentuk cerita pornografi.
4. Identifikasi tidak dapat dilakukan pada kesalahan penulisan kata yang bercampur dengan karakter angka maupun tanda baca.

Adapun kesimpulan dalam penelitian yang dilakukan Rendra Mahardika dari sisi saya sebagai peneliti selanjutnya:

1. Aplikasi yang dibuat menggunakan bahasa pemrograman Python.
2. Aplikasi yang dibuat menggunakan algoritma Support Vector Machine.
3. Proses pengambilan data latih dan data uji dilakukan diluar sistem yang dibuat dengan mengemas data ke dalam dataset terlebih dahulu.
4. Perancangan sistem utama yang dibuat terdiri dari halaman halaman pengujian dan halaman hasil pengujian

Alur proses implementasi .dalam Aplikasi Pendeteksi Konten Pornografi terdiri dari beberapa proses yaitu :

1. Proses Pengambilan data

Dalam proses mendapatkan data dari situs yang menjadi target pendeteksian konten pornografi, penulis menggunakan teknik scraping dalam program yang dibuat. Penulis memakai library HtmlAgilityPack 1.11.17 supaya bisa mengambil data text dari url/situs yang dituju.

2. Proses Pembersihan Data

Dalam proses pembersihan data yang telah diambil dari situs yang dituju, data akan melewati beberapa proses sampai data itu bersih dan dapat disimpan di database. Diantaranya :

a. Proses Tokenization

Dalam proses ini akan dihilangkan semua karakter kecuali huruf dalam data, dan hanya akan menyisakan kumpulan kata yang dipisahkan oleh karakter spasi saja. Penulis menggunakan fungsi

Regex yang sudah tersedia dalam framework .net core.

b. Proses Transformasi Text ke List.

Dalam proses ini sekumpulan kata akan diubah bentuknya kedalam bentuk list string menggunakan teknik split dengan menggunakan karakter spasi sebagai pemisah, lalu dilakukan proses looping untuk memasukan setiap kata ke dalam list string.

c. Proses Firlter Data Berdasarkan Panjangnya.

Dalam proses ini data yang sudah menjadi list akan diproses looping kembali dan dilakukan pemisahan kata berdasarkan panjangnya.

d. Proses Transformasi ke Hurup kecil.

Dalam proses ini list string akan dilakukan proses looping kembali supaya bisa menggunakan teknik replace setiap huruf besar dan akan diubah ke huruf kecil semua.

e. Proses Filter Berdasarkan Dictionary Data

Dalam proses ini list string akan filter untuk setiap kata yang sama dengan kata yang ada di data dictionary akan dieliminasi. Proses ini hanya menggunakan teknik if else dalam proses for each.

f. Proses Stemming Data

Dalam proses ini setiap kata disusun berdasarkan jenis kata dan jumlahnya ,lalu diberi label per situs apakah positif atau negative.

3. Proses Klasifikasi

Dalam proses kklasifikasi apakah suatu situs terdeteksi mengandung konten pornografi atau tidak, penulis menggunakan algoritma naive bayes dalam proses ini. Rumus atau persamaan yang digunakan adalah sebagai berikut :

$$P(w_k | v_j) = \frac{n_k + 1}{n + |kata|}$$

Dengan : w_k = Kelas /label .

v_j = Data testing.

n = adalah jumlah total kata yang terdapat di dalam data tekstual yang memiliki nilai fungsi target yang sesuai.

n_k = adalah jumlah kemunculan kata pada suatu kelas/label pada semua data tekstual yang memiliki nilai fungsi target yang sesuai.

IV. KESIMPULAN DAN SARAN

A.Kesimpulan

Berdasarkan pembahasan dan hasil penelitian yang telah dibahas pada bab sebelumnya maka dapat disimpulkan kesimpulan –kesimpulan sebagai berikut :

1. Pencegahan penyebaran konten pornografi dapat dicegah dengan cara mendeteksi isi konten situs tanpa membuka situs tersebut, yaitu dengan metode scraping. Hal ini sudah dibuktikan dengan implementasi pembuatan aplikasi pendeteksi konten pornografi yang telah dibuat. Aplikasi ini telah berhasil mengambil data pada konten suatu situs dan data tersebut juga telah berhasil di analisis untuk klasifikasi apakah mengandung konten pornografi atau tidak berdasarkan data model yang dimasukkan terlebih dahulu.
2. Berdasarkan tujuan yang diuraikan pada bab 1, Penulis telah berhasil membuat aplikasi pendeteksi konten pornografi dan telah berhasil mendeteksi apakah situs yang menjadi target mengandung konten pornografi atau tidak.

B.Saran

Adapun saran yang ingin Penulis sampaikan dalam penelitian adalah sebagai berikut :

1. Bagi peneliti selanjutnya diharapkan untuk bisa mengeksplorasi lebih banyak algoritma guna mendapatkan hasil yang lebih baik lagi untuk judul serupa.
2. Bagi pihak kampus diharapkan lebih giat lagi untuk merangsang peneliti selanjutnya supaya akan lebih banyak lagi peneliti yang mengambil judul penelitian yang bertemakan tentang Kecerdasan Buatan.

DAFTAR PUSTAKA

- [1] A. Rozi, Zaenal, dan Community, SmitDev.” *Bootstrap Design Framework*”. Jakarta: Elex Media Komputindo, 2015.
- [2] Dr. Suyanto. “*Data Mining untuk Klasifikasi dan Klusterisasi Data*”. Bandung: Informatika Bandung, 2019.
- [3] Fauzi, Rizki Ahmad. “*Sistem Informasi Akuntansi (Berbasis Akuntansi)*”. Yogyakarta: Deepublish, 2017.
- [4] Teguh Wahyono. “*Fundamental of Python for Machine Learning*”. Yogyakarta: Penerbit Gava Media, 2018.
- [5] Suyanto. “*Machine Learning Tingkat Dasar dan Lanjutan*”. Bandung: Informatika Bandung, 2018.
- [6] Marisa, Fitri. “*Web Programming (Client Side and Server Side)*”. Yogyakarta: Deepublish, 2017.
- [7] Nofriansyah, Dicky. “*Konsep data mining vs sistem pendukung keputusan*”. Yogyakarta: Deepublish, 2014.
- [8] Pratama, I Putu Agus Eka. “*Sistem Informasi dan Implementasinya*”, Bandung: Informatika Bandung, 2014.
- [9] Rusdiana., dan Irfan. “*Sistem Informasi Manajemen*”. Bandung: Pustaka Setia, 2014.
- [10] S, Rosa. A., dan M. Shalahuddin. “*Rekayasa Perangkat Lunak*”. Bandung: Informatika Bandung, 2015.
- [11] Luz, S. 2006. *Machine Learning of Text Categorization*. Trinity College, Department of Computer Science.
- [12] Mitchell, T. M. 1997. *Machine Learning*, Singapore, McGraw –Hill.
- [13] Dumais, S., Platt, J., Heckerman, D. & Sahami, M. 2002. *Inductive Learning Algorithms and Representations for Text Categorization*.