

# PEMANFAATAN *TEXT SUMMARIZATION* DENGAN *SUPPORT VECTOR MACHINE* DAN *K-NEAREST NEIGHBOR* PADA ANALISIS SENTIMEN UNTUK MEMPERMUDAH PENGGUNA MEMBACA *REVIEW GAME STEAM*

Hilarius Bryan, Rolly Intan, Hans Juwiantho

Program Studi Informatika, Fakultas Teknologi Industri, Universitas Kristen Petra

Jln. Siwalankerto 121- 131 Surabaya 60236

Telp. (031)-2983455, Fax. (031)-8417658

c14170045@john.petra.ac.id rintan@petra.ac.id, Hans.juwiantho@petra.ac.id

## ABSTRAK

Bertambah banyaknya perkembangan *game* sejalan dengan pertumbuhan para penikmatnya. Biasanya para penikmat ini yang sering disebut pemain memiliki platform khusus untuk melihat perkembangan *game* terbaru. Salah satu yang sering menjadi incaran adalah Steam, dimana platform tersebut menyediakan informasi lengkap seperti *review*, harga, tanggal rilis, dan sebagainya untuk pengguna yang ingin membeli *game*. Biasanya sebelum membeli sebuah *game* pengguna akan melihat *review* terlebih dahulu. Banyaknya *review* di dalam Steam menyebabkan kesulitan bagi pengguna untuk mencari informasi. Terdapat penelitian sebelumnya yang menyelesaikan problem tersebut dengan melakukan *text summarization* terhadap *review* menjadi beberapa aspek dan analisis sentimen akan tetapi terdapat kekurangan yaitu pengguna harus mengerti beberapa istilah dalam *game*. Dari permasalahan tersebut dilakukan *text summarization* untuk merangkum informasi dan sentimen analisis untuk menilai suatu nilai dari *game* tersebut.

Agar mendapat rangkuman informasi dengan baik perlu melalui beberapa proses pengolahan data. Untuk proses pengumpulan data *review game* didapatkan melalui API Steam yang telah tersedia. Setelah terkumpul akan dilakukan *preprocessing* untuk mengatasi data yang bervariasi dan tidak konsisten yang bisa mempengaruhi proses *training*. *Preprocessing* meliputi *Tokenization*, *Stopwords Removal*, dan *Stemming*. Proses *text summarization* pada tiap *review game* untuk *feature* to vector menggunakan TF-IDF dan *Sentiment Score* untuk mendapatkan kalimat utama sebelum dilakukan proses *training* menggunakan SVM. Proses klasifikasi menggunakan metode KNN dimana melakukan perbandingan dari tiap data *review game* apakah data tersebut lebih mendekati positif atau negatif, sehingga membantu pengguna ketika melihat informasi *game* menjadi lebih singkat dan mudah.

Pengukuran keberhasilan metode ini dalam menjawab permasalahan dengan melakukan pengujian data dengan *Confusion Matrix* dan survey kepada pengguna Steam. Penggunaan *text summarization* yang dilakukan terhadap tiap *review game* kurang berperan dalam meningkatkan hasil analisis sentimen, karena metode yang kurang sesuai serta data *review game* yang berbentuk abstrak. Akurasi analisis sentimen masih lebih baik ketika tidak dilakukan *text summarization* pada data. Sejumlah 50 orang diminta untuk memberikan pernyataan seputar hasil dari analisis sentimen dan *text summarization*. Hasil yang diperoleh 40 dari 50 user mengatakan aplikasi membantu membaca *review game* dan 10 lainnya tidak.

**Kata Kunci:** Analisis Sentimen, *Text summarization*, KNN, SVM, API Steam

## ABSTRACT

*Today the development of the game is increasing and in line with the growth of the players. Usually, these players who are often called players have a special platform to see the latest game developments. One that is often targeted is Steam, where the platform provides complete information such as reviews, prices, release dates, and so on for users who want to buy games. Usually before buying a game the user will see a review first. The number of reviews on Steam makes it difficult for users to find information. From these problems, text summarization was carried out to summarize information and sentiment analysis to assess the value of the game.*

*In order to get a good summary of the information, it is necessary to go through several data processing processes. The game review data collection process is obtained through the available Steam API. Once collected, preprocessing will be carried out to overcome the varied and inconsistent data that can affect the training process. Preprocessing includes Tokenization, Stopwords Removal, and Stemming. The text summarization process for feature to vector uses TF-IDF and Sentiment Score to get the main sentence before the training process using SVM is carried out. The classification process uses the KNN method which compares each game review data whether the data is closer to positive or negative, thus helping users when viewing game information becomes shorter and easier. Measurement of the success of this method in answering problems by testing data with the Confusion Matrix and surveying Steam users. The use of text summarization for each game review has little role in improving the results of sentiment analysis, because the method is not suitable and the game review data is in the form of an abstract. The accuracy of sentiment analysis is still better when text summarization is not carried out on the data. A total of 50 people were asked to provide statements regarding the results of sentiment analysis and text summarization. The results obtained by 40 out of 50 users said the application helped read game reviews and 10 others did not.*

**Keywords:** *Sentiment analysis, Text summarization, KNN, SVM, Steam API.*

## 1. INTRODUCTION

*Game* merupakan salah satu hiburan yang berkembang sangat pesat hingga sekarang. *Game* dimainkan melalui berbagai perangkat seperti *Personal Computer (PC)*, *Console* dan *Mobile Device*. Beberapa tahun terakhir jumlah pengguna Steam yang memberikan *review* telah meningkat secara signifikan. Hal ini disebabkan beberapa faktor yaitu jumlah *video game* dan platform yang memfasilitasi proses *review*. Data *review* game yang muncul oleh pengguna Steam membuat ekstraksi informasi menjadi lebih sulit jika harus diselesaikan secara tepat waktu dan efisien [10]. Steam merupakan salah satu distributor game digital untuk PC dan *Console*. Steam memiliki jumlah pengguna yang banyak, terbukti dari data pengguna yang aktif setiap bulan pada tahun 2019 mencapai 95 juta [14]. Tercatat 30.000 game sudah dirilis oleh Steam, tidak termasuk *software* dan *Downloadable Content (DLC)* [2]. Dengan pertumbuhan game yang pesat dan jumlah dan isi tiap *review* yang banyak membawa beberapa kesulitan kepada pengguna Steam dalam mendapatkan informasi dari tiap *review* [8].

Steam memiliki fitur *review game* dimana *user* dapat menilai kemenarikan sebuah produk *game*. *Review* tersebut memberikan informasi bermanfaat mengenai sebuah *game*, khususnya bagi pengguna yang ingin *download* (membeli). Permasalahan dari *review* tersebut, terkadang jumlah *review* sebuah *game* sangat banyak seperti Portal 2 memiliki jumlah *review* positif sebanyak 210,582 dan negatif 2646 [13]. Berdasarkan pengalaman peneliti, terdapat beberapa game yang memiliki informasi menjerumuskan dan membuat pengalaman bermain tidak sesuai dengan *review*. Dari hasil wawancara kepada 10 pengguna Steam, *review* memiliki peran yang cukup esensial dalam menentukan pembelian suatu *game*. Contoh dari salah satu survei pengguna Steam, terdapat game berjudul *Cyberpunk* yang memiliki *rating* dan *review* cukup tinggi. Setelah dimainkan, pengalaman *game* tersebut kontradiksi dengan apa yang diberikan pada *review*. Terjadi banyak *bug* dan *glitch* dalam *game* tersebut sehingga memberikan pengalaman bermain yang buruk. Pemain juga mengalami kerugian uang karena membeli *game* dengan harga tidak sesuai ekspektasi. Selain itu, tidak semua *review* mengandung informasi yang tersusun dan terpilah dengan jelas mengenai positif dan negatif dari sebuah produk. Hal ini menyulitkan *user* untuk menilai sebuah produk *game*. Oleh sebab itu, diperlukan sebuah aplikasi yang dapat membantu *user* untuk merangkum dan menilai sentimen *review* dari suatu game secara tepat. Salah satu tantangan dalam pengelompokan ini adalah bagaimana mengidentifikasi *false positif* dan *false negatif* dari sebuah kalimat.

Analisis Sentimen adalah sebuah teknik mengekstrak informasi berupa pandangan/sentimen seseorang terhadap suatu kejadian. Dengan analisis sentimen *user* terbantu untuk menilai suatu pendapat lebih akurat. Untuk menghasilkan analisis yang akurat perlu dilakukan beberapa tahapan. Pertama kali tahapan *preprocessing* yang bertujuan menghilangkan spesial karakter seperti koma, bintang, dan lainnya, menghilangkan angka numerik, menghilangkan *stopwords*, dan yang terakhir membetulkan salah pengejaan. Setelah melalui *preprocessing compound sentence* akan berubah menjadi *simple sentence*. Dari *simple sentence* tersebut akan dilakukan *summarization* menggunakan *Support Vector Machine*. Tujuan dilakukan *summarization* adalah menyederhanakan *review* agar kalimat yang ambigu menjadi lebih

tegas. Hal ini dinilai dapat meningkatkan akurasi pada proses analisa sentimen. Nantinya *simple sentence* akan menjadi inputan untuk analisis sentimen yang menggunakan *K-Nearest Neighbor*. Pada proses *K-Nearest Neighbor* dilakukan penilaian kalimat tergolong dalam positif, negatif, *false positif* dan *false negatif*.

Penelitian sebelumnya melakukan pendekatan yang berbeda pada proses analisis sentimen *review game* di platform Steam. Beberapa penelitian melakukan analisis *review game* pada Steam menggunakan *Naïve Bayes* dan *Decision Tree Classifier* dengan akurasi yang 75% [15]. Penelitian *Summarizing Game Reviews: First Contact* melakukan eksperimen dengan *review game* Steam yang bertujuan untuk menghasilkan *summarization* dari sentimen analisis dan identifikasi aspek. Proses pada penelitian melakukan *summarization* berdasarkan aspek *game* seperti *graphics*, *gameplay*, *audio*, *community*, *performance*, dan *story* [8].

Dari semua penelitian yang pernah dilakukan sebelumnya, belum ada yang spesifik meneliti text summarization berdasarkan analisis sentimen yang dapat menghasilkan pemecahan kalimat menjadi positif, negatif, *false positif* dan *false negatif*. Tujuannya untuk membantu user Steam dalam memberikan informasi sejauh mana sebuah game dapat memberikan ketertarikan pada pengguna. Proses yang dilakukan adalah pemecahan compound sentence menjadi *simple sentence*, kemudian dilakukan text summarization terhadap tiap *review game*, dan menjadi inputan untuk analisis sentimen. Dilakukan juga klasifikasi untuk sentimen agar memberikan hasil analisis *review* tersebut. Proses text summarization menggunakan metode *Support Vector Machine*. Sementara untuk Proses klasifikasi menjadi negatif, positif, *false negatif* dan *false positif* menggunakan metode *K-Nearest Neighbor*. Agar data yang didapatkan memiliki hasil yang akurat sesuai dengan persepsi pengguna perlu pengecekan data dan pemaknaan. Penilaian positif dan negatif sebuah kalimat dilakukan oleh pengguna *game* dan pakar bahasa. Dengan adanya penelitian ini diharapkan dapat mempermudah user dalam membaca *review* sebuah game di Steam.

## 2. LANDASAN TEORI

### 2.1 Natural Language Processing (NLP)

*Natural Language Processing (NLP)* adalah bidang ilmu dari *artificial intelligence* yang berhubungan dalam interaksi antara komputer dan manusia menggunakan *natural language*. Tujuan dari *Natural Language Processing* adalah untuk membaca, dekripsi, dan memahami bahasa manusia untuk membantu proses automasi.

Menurut Indurkha [7], *Natural Language Processing* terdapat beberapa teknik yang digunakan dalam mengenali bahasa manusia. Teknik yang digunakan terdapat 3 yaitu *syntax*, *semantics*, *pragmatics*. Penjelasan ilmu bahasa tersebut dapat diuraikan sebagai berikut:

- *Syntax*

*Syntax* mengacu pada pengaturan sistematis kata-kata dalam sebuah kalimat. Dalam analisis *syntax* NLP, bahasa alami digunakan untuk mengikuti aturan tata bahasa. *Syntax* adalah susunan kata-kata yang linier dalam sebuah kalimat.

- *Semantics*

*Semantics* adalah tafsir makna kalimat dalam bahasa. Interpretasi makna kalimat diperoleh dengan mengenali struktur gramatikal yang menyusun kalimat tersebut. Dengan cara ini, dengan mengidentifikasi tata bahasa kalimat, Anda bisa mendapatkan makna dari kalimat tersebut.

- *Pragmatics*

*Pragmatics* menjelaskan bagaimana pernyataan yang ada berhubungan dengan situasi yang ada. Dengan mempertimbangkan berbagai aspek seperti konteks dan tujuan kalimat. Kemudian, makna suara atau teks tersebut sesuai dengan konteks yang ingin disampaikan.

## 2.2 Text Mining

Penambangan teks adalah proses penggalian pola atau informasi dari banyak string atau pengaturan tidak terstruktur. Jika digabungkan menjadi sebuah kalimat, susunan ini memiliki arti. Namun, informasi ini tidak dapat ditangkap secara otomatis oleh komputer, oleh karena itu, penambangan teks adalah proses penggalian informasi atau makna dari teks, yang memungkinkan komputer untuk mengetahui makna teks tersebut [3].

## 2.3 Text Summarization

Peringkasan teks otomatis (*automated text summarization* atau ATS) adalah pembuatan ringkasan dari sebuah teks secara otomatis dengan memanfaatkan aplikasi yang dijalankan pada komputer. Sebuah sistem peringkasan diberi masukan berupa teks, kemudian melakukan peringkasan, dan menghasilkan keluaran berupa teks yang lebih singkat dari teks asli. Hasil peringkasan mengandung poin-poin penting atau informasi utama dari teks sumber. Terdapat dua pendekatan pada peringkasan teks, yaitu ekstraksi (*shallower approaches*) dan abstraksi (*deeper approaches*). Pada teknik ekstraksi, sistem menyalin unit-unit teks yang dianggap paling penting atau paling informatif dari teks sumber menjadi ringkasan. Unit-unit teks yang disalin dapat berupa klausa utama, kalimat utama, atau paragraf utama. Sedangkan teknik abstraksi melibatkan parafrase dari teks sumber. Teknik abstraksi mengambil intisari dari teks sumber, kemudian membuat ringkasan dengan menciptakan kalimat-kalimat baru yang merepresentasikan intisari teks sumber dalam bentuk berbeda dengan kalimat-kalimat pada teks sumber. Kalimat adalah satuan bahasa terkecil, dalam wujud lisan atau tulisan, yang mengungkapkan pikiran yang utuh. Berdasarkan jumlah sumbernya, sebuah ringkasan dapat dihasilkan dari satu sumber (*single-document*) atau dari banyak sumber (*multi-document*). Peringkasan *single-document* masukannya berupa sebuah teks dan keluarannya berupa sebuah teks baru yang lebih singkat [11].

## 2.4 Sentimen Analisis

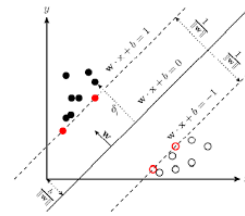
Analisis sentimen atau penggalian opini adalah pendeteksian sikap dari suatu objek atau orang. Analisis sentimen dapat digunakan untuk mendapatkan persentase sentimen positif dan negatif untuk individu, perusahaan, institusi, produk, atau situasi tertentu. Nilai analisis sentimen dapat dibedakan menjadi 3 yaitu emosi positif, emosi negatif dan emosi netral atau pendalaman, sehingga dapat diketahui siapa atau kelompok mana yang menjadi sumber emosi positif atau emosi negatif).

Analisis sentimen bertujuan untuk mengevaluasi emosi, sikap, pendapat, dan evaluasi pembicara atau penulis tentang produk atau tokoh masyarakat. Oleh karena itu, perlu dilakukan beberapa penelitian, khususnya di bidang review produk, sebelum mulai menentukan elemen produk yang dibahas sebelum proses penambangan opini [6].

## 2.5 Support Vector Machine

Support Vector Machine (SVM) merupakan salah satu metode dalam *supervised learning* yang biasanya digunakan untuk klasifikasi dan regresi. SVM melakukan pembelajaran dari korelasi berbagai data pada training set untuk memprediksi kelas baru. SVM digunakan untuk mencari *hyperplane* terbaik dengan

memaksimalkan jarak antar kelas [1]. *Hyperplane* adalah sebuah fungsi yang dapat digunakan untuk pemisah antar kelas.



**Gambar 1.** Proses *support vector machine* [5]

Dalam 2-D fungsi yang digunakan untuk klasifikasi antar kelas disebut sebagai *line whereas*, fungsi yang digunakan untuk klasifikasi antar kelas dalam 3-D disebut *plane similarly*, sedangkan fungsi yang digunakan untuk klasifikasi di dalam ruang kelas dimensi yang lebih tinggi disebut *hyperplane*. *Hyperplane* diperoleh dengan memaksimalkan margin pada masing-masing class. Terdapat sekumpulan titik atau data yang disebut juga dengan support vector yang dipisahkan oleh garis yang disebut dengan *hyperplane*. Support vector tersebut merupakan data yang letaknya paling dekat dengan *hyperplane* yang ada. *Hyperplane* yang memiliki nilai optimal mempunyai margin terbesar ke support vector.

## 2.6 K-Nearest Neighbor

*K Nearest Neighbors* adalah salah satu algoritma yang dilakukan untuk melakukan klasifikasi berdasarkan contoh dasar yang tidak membangun. Algoritma tersebut bergantung pada label kategori yang melekat pada dokumen pelatihan meniru dengan dokumen tes. Dari dokumen tes tersebut sistem menentukan nilai *k* tetangga terdekat dengan dokumen lainnya. [9]. Tujuan dari algoritma ini adalah mengklasifikasikan obyek baru berdasarkan atribut dan training sample. Classifier tidak menggunakan model apapun untuk dicocokkan dan hanya berdasarkan pada memori. Diberikan titik *query*, akan ditemukan sejumlah *k* obyek atau (titik training) yang paling dekat dengan titik *query*. Klasifikasi menggunakan *voting* terbanyak diantara klasifikasi dari *k* obyek. Algoritma *k-nearest neighbor*(K-NN) menggunakan klasifikasi ketetanggaan sebagai nilai prediksi dari *query instance* yang baru [12]. Persamaan *K-Nearest Neighbor* dengan rumus berikut.

$$pd_i = \sqrt{\sum_{i=1}^p (X_{2i} - X_{1i})^2}$$

**Rumus 1.** *K-Nearest Neighbor* [12]

Keterangan:

- $X_1$  : Sampel data
- $X_2$  : Data uji/ *testing*
- $i$  : Variabel data
- $d$  : Jarak
- $p$  : Dimensi Data

## 2.7 K-Fold Cross Validation

*K-fold cross validation* adalah teknik yang dapat digunakan untuk melakukan pengujian terhadap suatu data. *K-fold cross validation*

merupakan salah satu metode yang digunakan untuk mengetahui rata-rata keberhasilan dari suatu sistem dengan cara melakukan perulangan dengan mengacak atribut masukan sehingga sistem tersebut teruji untuk beberapa atribut input yang acak. K-fold cross validation diawali dengan membagi data sejumlah n-fold yang diinginkan. Dalam proses cross validation data akan dibagi dalam n buah partisi dengan ukuran yang sama  $N_1, N_2, N_3$  dan selanjutnya proses testing dan training dilakukan sebanyak n kali. Dalam iterasi ke-I partisi akan menjadi data testing dan sisanya akan menjadi data training [9].

### 2.8 Confusion Matrix

Confusion matrix adalah suatu metode yang digunakan untuk melakukan perhitungan akurasi pada *text mining*. Confusion matrix digambarkan dengan tabel yang menyatakan jumlah data uji benar diklasifikasikan dan data uji yang salah diklasifikasikan [4]. Berikut ada table *confusion matrix*.

Tabel 1. Confusion Matrix

Sumber: Han, J., & Kamber, M. (2006).

Confusion Matrix		Prediksi	
		Positif “+”	Negatif “-”
Aktual	Positif “+”	TP (True Positive)	FN (False Negative)
	Negatif “-”	FP (False Positive)	TN (True Negative)

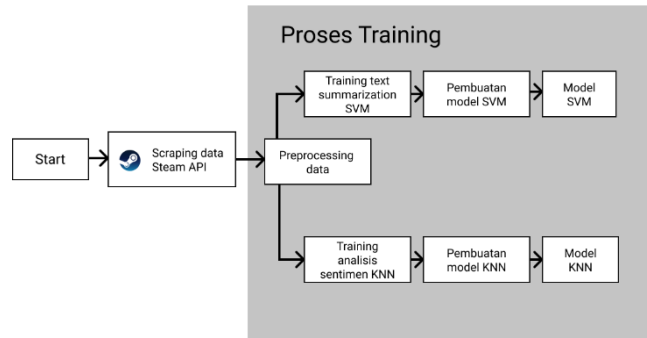
## 3. METODE

### 3.1 Data Training

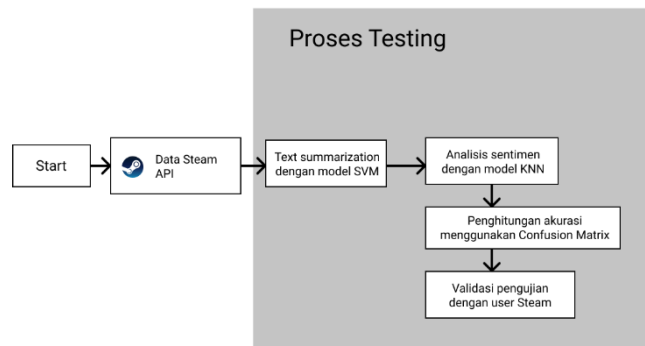
Proses *text summarization* pada penelitian bukan melakukan ringkasan secara keseluruhan pada *review game* tetapi merangkum tiap *review game* untuk menentukan dalam suatu game apakah lebih banyak review positif, negatif, false positif, dan false negatif.

Proses *scraping data* mengambil dari Steam API untuk mendapatkan data set yang berhubungan dengan game setelah itu dilakukan labeling dibantu oleh pakar Sheryl Keren Muliadie S.S, B.A. Proses data *preprocessing* adalah proses mengolah data yang telah dikumpulkan dari Steam, sebelum dilakukan proses training dan klasifikasi. Proses data *preprocessing* perlu dilakukan untuk mengatasi data yang memiliki variasi terlalu banyak dan tidak konsisten. Data *preprocessing* yang dilakukan pada skripsi ini meliputi *Tokenization, StopwordsRemoval, Stemming, TF-IDF*.

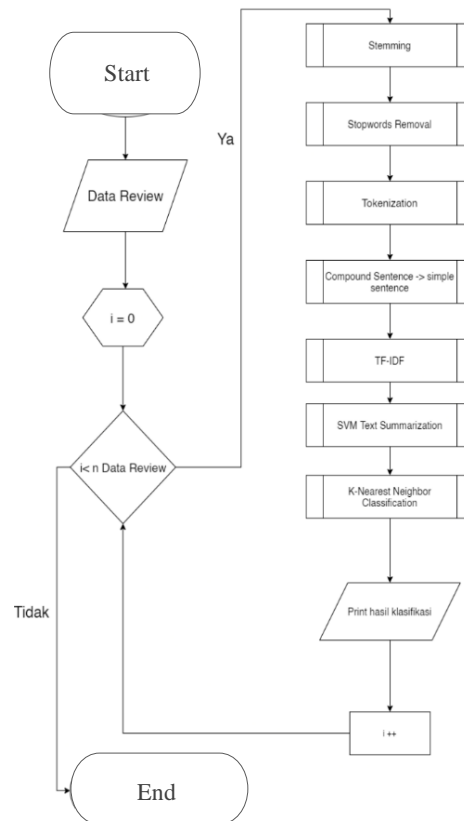
Setelah proses data *preprocessing* dilakukan proses pemecahan *compound sentence* menjadi *simple sentence* kemudian dari pemecahan tersebut dilakukan *text summarization* untuk menjadi inputan analisis sentimen. Setelah proses tersebut dilakuakn proses *training*, proses *training* adalah proses *classifier* mempelajari input data yang diberikan. Sebelum input data digunakan dalam proses *training*, data *training* akan diolah terlebih dahulu dengan data *preprocessing*. Proses training pada penelitian ini adalah *semi supervised* data yang digunakan pada proses training akan diberikan label/class terlebih dahulu oleh pakar, dan pengguna *Steam* dan oleh computer. Proses secara keseluruhan dapat dilihat pada gambar 1 dan gambar 2. Klasifikasi dibagi menjadi 4 kategori sentimen positif, negatif, false positif, dan false negatif. Gambaran umum dari proses training dapat dilihat pada Gambar 3.



Gambar 2. Block diagram proses training



Gambar 3. Block diagram proses testing



Gambar 4. Flowchart proses training

### 3.2 Tokenizing

Proses *Tokenization* adalah merupakan proses yang memotong kalimat pada data anotasi menjadi kata atau karakter, kemudian dari hasil tersebut akan dilakukan pengambilan kata positif dan negatif. Data *review* akan diolah, kemudian akan dilakukan pemecahan data *review* dengan menggunakan NLTK. Dalam tahap tokenization juga dilakukan penghilangan karakter – karakter seperti dan tanda baca yang lainnya. Kemudian dilakukan pemecahan kalimat *review* berdasarkan kata sambung.

### 3.3 Stopwords Removal

Proses *stopwords removal* menghapus kata – kata umum yang tidak memiliki makna. Proses ini dilakukan untuk mengurangi variasi kata tanpa mengubah maknanya. Proses *stopwords removal* pada aplikasi ini menggunakan *library* NLTK dalam *library* tersebut sudah memiliki daftar makna kata yang tidak memiliki arti apabila dihapus.

### 3.4 Stemming

Proses *stemming* adalah menghapus imbuhan pada kata dalam suatu *review* sehingga setiap kata tersisa hanya kata pokok. Proses ini dilakukan untuk mengurangi variasi kata tanpa mengubah maknanya. Proses *stemming* pada aplikasi ini menggunakan *library* NLTK. *Library*. Dalam NLTK terdapat fungsi *stemming* yang dapat membantu *preprocessing* data untuk menghapus imbuhan kata tanpa menghilangkan maknanya.

### 3.5 TF-IDF (Term Frequency – Inverse Document Frequency)

Metode TF-IDF merupakan metode pembobotan yang merupakan integrasi antar *term frequency* (tf), dan *inverse document* (idf). Apabila sebuah kata semakin sering muncul pada sebuah dokumen maka kata tersebut semakin penting. Pada proses ini menggunakan *library* Scikit yang disediakan oleh python. *Library* ini akan melakukan perhitungan pada suatu *review* untuk menentukan kata yang penting untuk membantu dalam proses *text summarization*.

### 3.6 Sentiment Score

Proses penghitungan sentiment score merupakan metode untuk melakukan penghitungan *word to vector*. Proses penghitungan dilakukan dengan menghitung tiap kata yang mengandung emosi menjadi vector, apabila terdapat kata positif di dalam sebuah kalimat maka score akan menjadi +1. Bila terdapat kata negative dalam sebuah kalimat maka score akan menjadi -1.

### 3.7 Support Vector Machine Text Summarization

Metode yang digunakan untuk *text summarization* adalah *Support Vector Machine*, metode ini bertujuan untuk mendapatkan bidang pemisah / *hyperplane* yang terbaik dengan memiliki prinsip *Empirical Risk Minimization* (ERM) yaitu meminimalkan *error* pada data pelatihan. Inputan dalam proses SVM yaitu fitur-fitur dari kalimat yang diambil yaitu TF-IDF, *positive score*, *negative score*, dari fitur tersebut akan dilakukan pemisahan dan mendapatkan kalimat utama. Dalam proses *text summarization* menggunakan metode *Support Vector Machine* menggunakan 3 *feature* yaitu TF-IDF, positif *sentiment score*, dan negatif *sentiment score*. Proses *text summarization* menggunakan 3 *feature* tersebut untuk menghitung nilai vector dari suatu kalimat dari *feature* tersebut akan diklasifikasi menggunakan *hyperplane* yang akan mendapatkan hasil kalimat *summarization*. Inputan dari proses ini adalah kalimat dari data *training* kemudian dilakukan proses *learning* yang akan mendapatkan hasil kalimat utama.

### 3.8 K-nearest Neighbor Classification

Metode yang digunakan untuk klasifikasi adalah *K-Nearest Neighbor*. Proses training pada *K-Nearest Neighbor* menggunakan *library* Scikit. Inputan berupa vector dari data uji yang akan dihitung jaraknya dari data training, setelah melakukan penghitungan akan mendapatkan hasil jarak terdekat dari data training tersebut untuk mendapatkan hasil klasifikasi. Pada *library* tersebut akan membantu proses *learning* dari data yang sudah diringkas dari metode sebelumnya dan data yang sudah dilabeli terlebih dahulu.

## 4. PENGUJIAN

Pengujian dilakukan untuk mendapatkan metode mana yang perlu dilakukan dari beberapa metode yang telah diusulkan, sehingga menghasilkan akurasi klasifikasi tertinggi. Jumlah perbandingan data juga dihitung sebagai bahan pertimbangan.

Pengujian pada skripsi ini terbagi menjadi dua yaitu, pengujian dengan *k-Cross Validation* dan pengujian data tes diluar *dataset training*. Pada pengujian *confusion matrix*. Pengujian 2 kali dilakukan dengan menggunakan analisis sentimen saja dan analisis sentimen, *text summarization* dari dataset yang telah dipisah menggunakan metode *K-Nearest Neighbor* dan *Support Vector Machine*. Pengujian *text summarization* dilakukan dengan *semantic scoring* dan *grammar scoring* dari pakar dengan skala 1-5. Pengujian dilakukan dengan menggunakan data tes yang telah diambil diklasifikasikan dengan model yang telah di dapat dari proses *training*.

Pengujian terakhir dilakukan kepada 50 user steam untuk validasi apakah aplikasi berhasil membantu pengguna dalam membaca *review game* dan membantu keputusan analisis sentimen

Tabel 2. Hasil *scoring text summarization*  
Sumber: Hasil survei

Score	Semantic Scoring	Grammar Scoring
1	0	0
2	0	0
3	2	2
4	7	7
5	1	1
Rata – rata	7.8	7.8

Tabel 3. Hasil pengujian KNN dan SVM  
Sumber: Hasil testing aplikasi

Data	Akurasi	Precision	Recall	Fscore
500 positif & 500 negatif	0.668	0.694	0.919	0.791
600 positif & 400 negatif	0.56	0.562	0.9	0.692
400 positif & 600 negatif	0.636	0.604	0.78	0.681

**Tabel 4. Hasil 3-Fold Cross Validation**  
**Sumber: Hasil testing aplikasi**

Metode	Akurasi	Precision	Recall	Fscore
KNN	0,637	0,636	0,835	0,729
KNN & SVM	0,599	0,645	0,725	0,69

**Tabel 5. Hasil 5-Fold Cross Validation**  
**Sumber: Hasil testing aplikasi**

Iterasi	Akurasi	Precision	Recall	Fscore
KNN	0,623	0,655	0,814	0,715
KNN & SVM	0,602	0,664	0,729	0,688

**Tabel 6. Hasil survei pengujian aplikasi**  
**Sumber: Survei**

Pertanyaan	Ya	Tidak
Apakah hasil analisis sentimen dari website sesuai pendapat user	45	5
Apakah aplikasi analisis sentimen dan text summarization membantu user ketika membaca review game	40	10

## 5. KESIMPULAN

Berdasarkan hasil survey 40 pengguna Steam merasa hasil *text summarization* lebih sederhana dan menentukan isi dari *review* tersebut sehingga membantu dalam membaca *review game*. Pada pengujian menggunakan *k-cross validation* 3 ketika menggunakan analisis sentimen mendapatkan hasil rata-rata 61 % kemudian ketika menggunakan analisis sentimen dan *text summarization* mendapatkan hasil rata-rata 59%. Pada pengujian menggunakan *k-cross validation* 5 ketika menggunakan analisis sentimen mendapatkan hasil rata-rata 62 % kemudian ketika menggunakan analisis sentimen dan *text summarization* mendapatkan hasil rata-rata 60%.

Kesimpulan dari hasil percobaan penggunaan Text Summarization tiap *review game* kurang berperan untuk meningkatkan hasil dari analisis sentiment, malah menjadi penghambat. Akurasi analisis sentiment yang didapat setelah summarization malah menurun ketimbang data yang diproses langsung, Text Summarization menggunakan metode SVM kurang sesuai dengan dataset yang mengandung banyak kata tidak baku dan kalimat ambigu, sehingga hasil kurang optimal. Akurasi sentiment masih lebih baik ketika dataset tidak melalui proses summarization, karena konteks kalimat belum berubah.

Saran untuk pengembangan lebih lanjut adalah menambah dataset *review game* untuk meningkatkan hasil dan akurasi, Menggunakan metode yang lebih sesuai untuk *text summarization* seperti

textrdrank dan LSA. Summarization dibuat bentuk poin tidak hanya memotong kalimat.

## 6. REFERENCES

- [1] Bam, S. B., & Shahi, T. B. (2014). Named Entity Recognition for Nepali Text Using Support Vector Machines. *Intelligent Information Management*, 21-25.
- [2] Bolding, Jonathan. *PC Gamer* 14 January 2019. <https://www.pcgamer.com/steam-now-has-30000-games/>. Bolding, J. (2019, January 14). *PC Gamer*. Retrieved from [pcgamer.com](https://www.pcgamer.com/steam-now-has-30000-games/): <https://www.pcgamer.com/steam-now-has-30000-games/>
- [3] Christianto, M., Andjarwirawan, J., & Tjondrowiguno, A. (1-5). Aplikasi Analisa Sentimen Pada Komentar Berbahasa Indonesia Dalam Objek Video di Website Youtube. *Jurnal Infra*, 2020.
- [4] Han, J., & Kamber, M. (2006). *Data Mining: Concepts and Techniques*. Illinois: University of Illinois at Urbana-Champaign.
- [5] Haritama W, A. A. (2017). Penerapan Model Mesin Belajar Support Vector Machines Pada Automatic Scoring untuk Jawaban Singkat. 15.
- [6] Haryanto, M. R. (2019). Analisis Sentimen Tokoh Politik di Media Sosial Twitter menggunakan Metode Naive Bayes dan Simple Additive Weighting.
- [7] Indurkha, N., & Damerou, F. J. (2010). *Handbook of Natural Language Processing*. 2nd edition, (pp. 3-7).
- [8] Kosmopoulos, A., Liapis, A., Giannakopoulos, G., & Pittaras, N. (2020). Summarizing Game Reviews: First Contact. *EETN Conference on Artificial Intelligence 2020*. Athens: SETN 2020
- [9] Nugroho, M. A., & Santoso, H. A. (2016). Klasifikasi Dokumen Komentar Pada Situs Youtube Menggunakan Algoritma K-Nearest Neighbor. 4-6.
- [10] Panagiotopoulos, G., Giannakopoulos, G., & Liapis, A. (2019). A Study on Video Game Review Summarization. *Proceedings of the Multiling 2019 Workshop* (pp. 35-41). Varna: INCOMA Ltd.
- [11] Rani, R., & Tandon, S. (2018). Chat Summarization and Sentiment Analysis Techniques in Data Mining. *2018 4th International Conference on Computing Sciences (ICCS)*.
- [12] Sahara, S. (2016). Penerapan Metode K-Nearest Neighbors Untuk Analisis Sentiment Review Game pada Android. *Jurnal Evolusi Volume 4*, 38-41.
- [13] Steam. *Steam*. 2021. <https://steamcommunity.com/dev> Steam. (2021). *Steam*. Retrieved from Steam Web API Documentation: <https://steamcommunity.com/dev>
- [14] Strickland, D. (2019, February 8). *TweakTown*. Retrieved from [tweaktown.com](https://www.tweaktown.com/news/70495/steam-grew-to-nearly-95-million-monthly-active-users-in-2019/index.html): <https://www.tweaktown.com/news/70495/steam-grew-to-nearly-95-million-monthly-active-users-in-2019/index.html>
- [15] Zuo, Z. (2018). Sentiment Analysis of Steam Review Datasets using Naive Bayes and Decision Tree Classifier.