



**PROJECT REPORT  
ANALYZING AND PREDICTING SENTIMENT RESULT  
FOR AMAZON PRODUCTS USING TEXTBLOB AND  
SUPPORT VECTOR MACHINE**

**TAN, DITYA YOSAPUTRA SISWOKO  
18.K1.0018**

**Faculty of Computer Science  
Soegijapranata Catholic University  
2022**

## HALAMAN PENGESAHAN



Judul Tugas Akhir: : Analyzing and predicting sentiment result for Amazon products using TextBlob and Support Vector Machine

Diajukan oleh : Tan, Ditya Yosaputra Siswoko

NIM : 18.K1.0018

Tanggal disetujui : 15 Juli 2022

Telah setuju oleh

Pembimbing : Yulianto Tejo Putranto S.T., M.T.

Penguji 1 : Yonathan Purbo Santosa S.Kom., M.Sc

Penguji 2 : Yulianto Tejo Putranto S.T., M.T.

Penguji 3 : Hironimus Leong S.Kom., M.Kom.

Penguji 4 : R. Setiawan Aji Nugroho S.T., MCompIT., Ph.D

Penguji 5 : Rosita Herawati S.T., M.I.T.

Penguji 6 : Y.b. Dwi Setianto S.T., M.Cs.

Ketua Program Studi : Rosita Herawati S.T., M.I.T.

Dekan : Dr. Bernardinus Harnadi S.T., M.T.

Halaman ini merupakan halaman yang sah dan dapat diverifikasi melalui alamat di bawah ini.

[sintak.unika.ac.id/skripsi/verifikasi/?id=18.K1.0018](http://sintak.unika.ac.id/skripsi/verifikasi/?id=18.K1.0018)

# HALAMAN PERNYATAAN PUBLIKASI KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS

Yang bertanda tangan dibawah ini:

Nama : Tan, Ditya Yosaputra Siswoko

Program Studi : Teknik Informatika

Fakultas : Ilmu Komputer

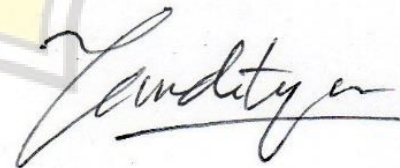
Jenis Karya : Skripsi

Menyetujui untuk memberikan kepada Universitas Katolik Soegijapranata Semarang Hak Bebas Royalti Noneklusif atas karya ilmiah yang berjudul “Analyzing and predicting sentiment result for Amazon products using TextBlob and Support Vector Machine” beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneklusif ini Universitas Katolik Soegijapranata berhak menyimpan, mengalihkan media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan tugas akhir ini selama tetap mencantumkan nama saya sebagai penulis / pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Semarang, 23 Juli 2022

Yang menyatakan



Tan, Ditya Yosaputra Siswoko

## DECLARATION OF AUTHORSHIP

I, the undersigned:

Name : Tan, Ditya Yosaputra Siswoko

ID : 18.K1.0018

declare that this work, titled "ANALYZING AND PREDICTING SENTIMENT RESULT FOR AMAZON PRODUCTS USING TEXTBLOB AND SUPPORT VECTOR MACHINE", and the work presented in it is my own. I confirm that:

1. This work was done wholly or mainly while in candidature for a research degree at Soegijapranata Catholic University
2. Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
3. Where I have consulted the published work of others, this is always clearly attributed.
4. Where I have quoted from the work of others, the source is always given.
5. Except for such quotations, this work is entirely my own work.
6. I have acknowledged all main sources of help.
7. Where the work is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Semarang, July 23, 2022



Tan, Ditya Yosaputra Siswoko

18.K1.0018

## ACKNOWLEDGMENT

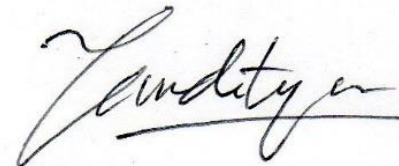
Praise to God the Almighty for His gift and mercy that I have finished tesis titled “ANALYZING AND PREDICTING SENTIMENT RESULT FOR AMAZON PRODUCTS USING TEXTBLOB AND SUPPORT VECTOR MACHINE” as the final requirement in accomplishing undergraduate at Informatics Engineering Departement Computer Science Faculty on Soegijapranata Catholic University Semarang.

For all the challenges and the difficulties I faced during on accomplishing this tesis, with any supports and guidance from any sides, so I should express my gratitude to :

1. My supervisor, Mr. Yulianto Tejo Putranto for supports, guides and directs me to finishing my project patiently.
2. All of lecturer from Informatics Engineering Departement who have teach, educate, shape, and motivate me into I am now to give me a way to the future on my study.
3. All of my family and relatives who supports me on undergoing study period with both materially and morally so that I can study well without any problems for giving me the best on my future.

Though of my limitations and mistake on this creation and research. Hopefully the constructive critics and suggestions could perfecting this tesis and encourage us to able to perfect and better research in the future.

Semarang, July 23, 2022



Tan, Ditya Yosaputra Siswoko

## ABSTRACT

Nowadays, numerous of reviews from buyers spreaded widely on online shopping platform, especially Amazon which it give an important information for knowing how good the condition and service of this item from various available categories. The bulk of reviews from buyers complicate to indentify all of sentiments of many comments, so that program was made to identify and analyze all of sentiments of many comments using sentiment analysis method.

This program identifies sentiment from each of reviews based of average valuation from automatic polarity of reviews using TextBlob library from Python programming language and the overall of each comments which changed into certain numbers such as -1, -0,5, 0, 0,5, 1 for 1, 2, 3, 4, 5 overalls. If there are a comment which have helpful votes, so there are additional changed overall values to count the average of this sentiment valuation. The average valuation of sentiment will determine whether sentiment are positive or negative. Then the sentiment result of reviews are processed into Support Vector Machine (SVM) algorithm which later measuring accuracy level and predict other reviews from another prediction dataset to prove exactness of it.

From the project result that performed, the Polynomial Kernel from SVM algorithm is the highest accuracy of 92,85% on prediction result in Without Stemming experiment and for highest examination of sentiment identification prediction accuracy is Radial Basis Function (RBF) from SVM algorithm with 92,64% which 667 of 720 from another dataset are correctly predicted in Without Undersampling experiment.

*Keyword: Sentiment Analysis, TextBlob library, Support Vector Machine, Python programming language*

## TABLE OF CONTENTS

<b>COVER</b> .....	<b>i</b>
<b>HALAMAN PENGESAHAN</b> .....	<b>ii</b>
<b>HALAMAN PERNYATAAN PUBLIKASI KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS</b> .....	<b>iii</b>
<b>DECLARATION OF AUTHORSHIP</b> .....	<b>iv</b>
<b>ACKNOWLEDGMENT</b> .....	<b>v</b>
<b>TABLE OF CONTENTS</b> .....	<b>vii</b>
<b>LIST OF FIGURE</b> .....	<b>x</b>
<b>LIST OF TABLE</b> .....	<b>1</b>
<b>CHAPTER 1 INTRODUCTION</b> .....	<b>3</b>
1.1. Background .....	3
1.2. Problem Formulation.....	4
1.3. Scope.....	4
1.4. Objective .....	4
<b>CHAPTER 2 LITERATURE STUDY</b> .....	<b>5</b>
<b>CHAPTER 3 RESEARCH METHODOLOGY</b> .....	<b>9</b>
3.1. Literature Study.....	9
3.2. Datasets and Code .....	9
3.3. Implementation and Analysis .....	9
3.3.1. Data Crawling and Collecting .....	9
3.3.2. Preprocessing Data.....	10
3.3.3. TextBlob .....	12
3.3.4. Sentiment Mean.....	12
3.3.5. Undersampling Data.....	15



3.3.6.	TF-IDF .....	16
3.3.7.	SVM Algorithm .....	17
3.3.8.	Prediction Dataset.....	18
<b>CHAPTER 4 ANALYSIS AND DESIGN .....</b>		<b>19</b>
4.1.	Analysis.....	19
4.2.	Design .....	23
<b>CHAPTER 5 IMPLEMENTATION AND RESULTS.....</b>		<b>25</b>
5.1.	Implementation .....	25
5.1.1.	Dataset preprocessing.....	25
5.1.2.	Sentiment determiner .....	28
5.1.3.	Balancing the dataset.....	29
5.1.4.	Train test split .....	30
5.1.5.	Term Frequency–Inverse Document Frequency (TF-IDF) .....	30
5.1.6.	Support Vector Machine (SVM).....	30
5.1.7.	Dataset Prediction.....	31
5.2.	Results.....	32
5.2.1.	Train and Test Accuracy on Main Dataset .....	32
5.2.2.	Train and Test Accuracy on Main Dataset without Stemming.....	33
5.2.3.	Train and Test Accuracy on Main Dataset without Lemmatization .....	34
5.2.4.	Train and Test Accuracy on Main Dataset without Helpful Votes Calculation.....	35
5.2.5.	Train and Test Accuracy on Main Dataset without Undersampling .....	36
5.2.6.	Accuracy on Prediction Dataset .....	37
5.2.7.	Accuracy on Prediction Dataset without Stemming .....	38
5.2.8.	Accuracy on Prediction Dataset without Lemmatization .....	38
5.2.9.	Accuracy on Prediction Dataset without Helpful Votes Calculation .....	39



5.2.10. Accuracy on Prediction Dataset without Undersampling.....	39
5.2.11. Conclusion from whole Accuracy Results.....	40
<b>CHAPTER 6 CONCLUSION .....</b>	<b>42</b>
<b>REFERENCES.....</b>	<b>43</b>
<b>APPENDIX.....</b>	<b>a</b>



## LIST OF FIGURE

<b>Figure 3.1</b> The result of all sentiment calculation with overalls in diagram.....	14
<b>Figure 3.1</b> Example from a dataset that contains A and B group .....	15
<b>Figure 3.2</b> Different result after 1 <sup>st</sup> sampling (left) and 2 <sup>nd</sup> sampling (right).....	16
<b>Figure 4.1</b> Flowchart of analyzing and predicting sentiment result.....	23



## LIST OF TABLE

<b>Table 4.1.</b> Main Dataset (After remove asin, name, title, verified column) .....	19
<b>Table 4.2.</b> Prediction dataset .....	21
<b>Table 5.1.</b> Accuracy evaluation on main dataset.....	32
<b>Table 5.2.</b> Accuracy evaluation on main dataset (5-fold) .....	32
<b>Table 5.3.</b> Accuracy evaluation on main dataset (10-fold) .....	32
<b>Table 5.4.</b> Accuracy evaluation on main dataset without stemming.....	33
<b>Table 5.5.</b> Accuracy evaluation on main dataset without stemming (5-fold) .....	33
<b>Table 5.6.</b> Accuracy evaluation on main dataset without stemming (10-fold) .....	33
<b>Table 5.7.</b> Accuracy evaluation on main dataset without lemmatization.....	34
<b>Table 5.8.</b> Accuracy evaluation on main dataset without lemmatization (5-fold) .....	34
<b>Table 5.9.</b> Accuracy evaluation on main dataset without lemmatization (10-fold) .....	34
<b>Table 5.10.</b> Accuracy evaluation on main dataset without helpful votes calculation .....	35
<b>Table 5.11.</b> Accuracy evaluation on main dataset without helpful votes calculation (5-fold) .....	35
<b>Table 5.12.</b> Accuracy evaluation on main dataset without helpful votes calculation (10- fold).....	35
<b>Table 5.13.</b> Accuracy evaluation on main dataset without undersampling.....	36
<b>Table 5.14.</b> Accuracy evaluation on main dataset without undersampling (5-fold) .....	36
<b>Table 5.15.</b> Accuracy evaluation on main dataset without undersampling (10-fold) .....	36
<b>Table 5.16.</b> Predicted accuracy on 720 reviews .....	37
<b>Table 5.17.</b> Predicted accuracy on 720 reviews in percentage.....	37
<b>Table 5.18.</b> Predicted accuracy on 720 reviews without stemming .....	38
<b>Table 5.19.</b> Predicted accuracy on 720 reviews without stemming in percentage.....	38

<b>Table 5.20.</b> Predicted accuracy on 720 reviews without lemmatization .....	38
<b>Table 5.21.</b> Predicted accuracy on 720 reviews without lemmatization in percentage ...	38
<b>Table 5.22.</b> Predicted accuracy on 720 reviews without helpful votes calculation .....	39
<b>Table 5.23.</b> Predicted accuracy on 720 reviews without helpful votes calculation in percentage .....	39
<b>Table 5.24.</b> Predicted accuracy on 720 reviews without undersampling .....	39
<b>Table 5.25.</b> Predicted accuracy on 720 reviews without undersampling in percentage...	40
<b>Table 5.26.</b> Highest accuracy Comparison between 5 experiments.....	40
<b>Table 5.27.</b> Tendency of 4 other experiments against normal experiments.....	41

