



PROJECT REPORT
SENTIMENT ANALYSIS ON MOVIE REVIEW USING
RANDOM FOREST AND LOGISTIC REGRESSION
ALGORITHM

LADY VIONA SUGIANTO
18.K1.0023

Faculty of Computer Science
Soegijapranata Catholic University
2022

HALAMAN PENGESAHAN



Judul Tugas Akhir: : Sentiment Analysis on Movie Review Using Random Forest and Logistic Regression Algorithm

Diajukan oleh : Lady Viona Sugianto

NIM : 18.K1.0023

Tanggal disetujui : 23 Mei 2022

Telah setuju oleh

Pembimbing : Hironimus Leong S.Kom., M.Kom.

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

Penguji 2 : Hironimus Leong S.Kom., M.Kom.

Penguji 3 : Rosita Herawati S.T., M.I.T.

Penguji 4 : Y.b. Dwi Setianto S.T., M.Cs.

Penguji 5 : R. Setiawan Aji Nugroho S.T., MCompIT., Ph.D

Penguji 6 : Yulianto Tejo Putranto S.T., M.T.

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.K1.0023

DECLARATION OF AUTHORSHIP

I, the undersigned:

Name : LADY VIONA SUGIANTO

ID : 18.K1.0023

declare that this work, titled "SENTIMENT ANALYSIS ON MOVIE REVIEW USING RANDOM FOREST AND LOGISTIC REGRESSION ALGORITHM", 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, June, 03, 2022



LADY VIONA SUGIANTO

18.K1.0023

STATEMENT PAGE OF SCIENTIFIC PUBLICATIONS FOR ACADEMIC INTEREST

I, the undersigned:

Name : LADY VIONA SUGIANTO

Study Program : Informatics Engineering

Faculty : Computer Science

Type of Work : Thesis

Approve to grant Soegijapranata Catholic University Semarang Non-exclusive Royalty-Free Rights for the scientific works entitled “SENTIMENT ANALYSIS ON MOVIE REVIEW USING RANDOM FOREST AND LOGISTIC REGRESSION ALGORITHM” along with existing tools (if needed). With these rights, database form, maintain, and publish this final project as long as it includes my name as the writer/creator and as the copyright owner.

This statement was made truthfully.

Semarang, June, 03, 2022



LADY VIONA SUGIANTO

18.K1.0023

ACKNOWLEDGMENT

First of all, I would like to thank God for His blessings and grace, so that I can complete a thesis entitled "Sentiment Analysis on Movie Review Using Random Forest and Logistic Regression Algorithm".

I am also very grateful to Hironimus Leong S. Kom., M. Kom as a supervisor at Soegijapranata University who has assisted in writing this thesis.

I would like to thank my family and friends who always support and motivate me so that this thesis can run well and smoothly.

Finally, I realize that the writing of this thesis is still far from perfect. Therefore, I ask for his criticism and suggestions for the sake of building the perfection of this thesis. Hopefully writing this thesis can be useful and provide more knowledge for readers.

Semarang, June, 03, 2022



LADY VIONA SUGIANTO

18.K1.0023

ABSTRACT

In the analysis of the movie classification of a movie, users tend to see whether a movie is good or not from the movie's rating. Most users are not interested in seeing other user reviews before, even though the rating is always tied to the reviews given by users to find out whether the film is good or not. Due to the large number of reviews left by users on a movie website, users have difficulty knowing the classification of positive reviews and negative reviews, and users also do not know the quality of a movie.

Therefore, in this project, sentiment analysis was made for a film review. By using 2 algorithms, namely Random Forest Classifier and Logistic Regression to determine the sentiment analysis of each film. Both algorithms determine positive and negative sentiments. The data structure used is a database, which will train data and test data for comparison purposes. The dataset is also converted to CSV for easier analysis. To evaluate the model, a comparison of algorithms based on accuracy was carried out.

From implementing the Random Forest algorithm with a data range of 650-2000, the resulting analysis sentiment is 9 positive and 1 negative. While the Logistic Regression algorithm produces 8 positive and 2 negative sentiment analyses. And from the comparison of accuracy, Random Forest is better and more suitable for this research because it has an average accuracy of 72,9261% while the average accuracy of Logistic Regression is 63,516%.

Keyword: sentiment analysis, movie review, logistic regression, random forest

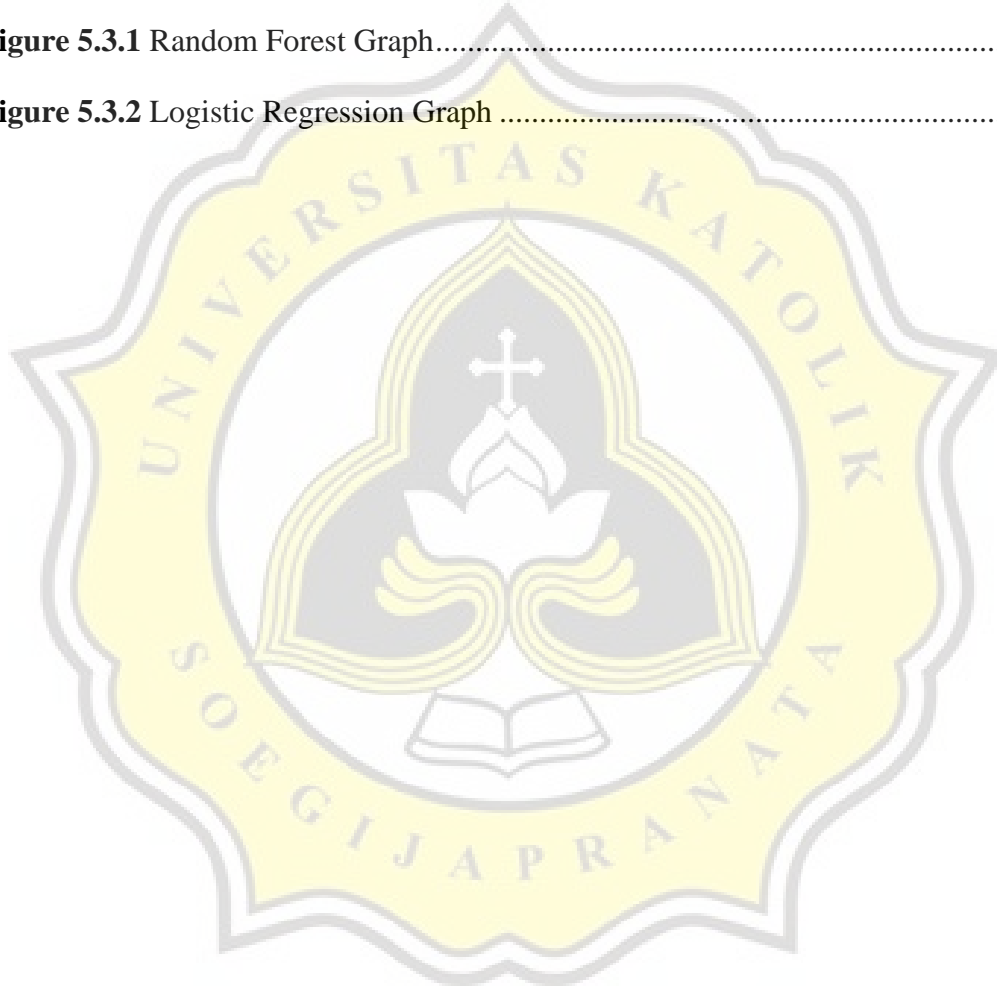
TABLE OF CONTENTS

COVER.....	i
APPROVAL AND RATIFICATION PAGE (Heading plain)	ii
DECLARATION OF AUTHORSHIP	iii
ACKNOWLEDGMENT	iv
ABSTRACT.....	vi
TABLE OF CONTENTS	vii
LIST OF FIGURE	ix
LIST OF TABLE	x
CHAPTER 1 INTRODUCTION.....	12
1.1. <i>Background</i>	12
1.2. <i>Problem Formulation</i>	13
1.3. <i>Scope</i>	13
1.4. <i>Objective</i>	13
CHAPTER 2 LITERATURE STUDY	14
CHAPTER 3 RESEARCH METHODOLOGY.....	19
3.1 <i>Data Collection</i>	19
3.2 <i>Algorithms</i>	19
3.3 <i>Design</i>	20
3.4 <i>Coding</i>	20
3.5 <i>Analysis</i>	21
CHAPTER 4 ANALYSIS AND DESIGN.....	22
4.1. <i>Random Forest</i>	22
4.1.1. <i>Getting Data</i>	22
4.1.2. <i>Data Preprocessing</i>	24

4.1.3.	<i>TF-IDF (Term Frequency-Inverse Document Frequency)</i>	25
4.1.4.	<i>Split Data</i>	31
4.1.5.	<i>Implementation of Random Forest Algorithm</i>	34
4.2.	Logistic Regression	58
4.2.1.	<i>Getting Data</i>	58
4.2.2.	<i>TF-IDF (Term Frequency-Inverse Document Frequency)</i>	59
4.2.3.	<i>Text Processing</i>	61
4.2.4.	<i>Feature Extraction</i>	66
4.2.5.	<i>Implementation of Logistic Regression Algorithm</i>	69
CHAPTER 5 IMPLEMENTATION AND RESULTS		78
5.1.	<i>Implementation</i>	78
5.1.1	<i>Random Forest</i>	78
5.1.2	<i>Logistic Regression</i>	91
5.2.	<i>Results</i>	100
5.3.	<i>Analysis</i>	102
CHAPTER 6 CONCLUSION		108
REFERENCES		109
APPENDIX		110

LIST OF FIGURE

Figure 4.1.1 Random Forest Workflow	22
Figure 4.1.5.1 Representation of Random Forest	34
Figure 4.2.1 Logistic Regression Design Scheme	58
Figure 4.2.5.1 Curve of Sigmoid Function	70
Figure 5.3.1 Random Forest Graph.....	106
Figure 5.3.2 Logistic Regression Graph	107



LIST OF TABLE

Table 4.1.1.1 Random Forest Dataset	23
Table 4.1.2.1 Data Preprocessing.....	24
Table 4.1.3.1 TempWord	26
Table 4.1.3.2 CountWord.....	27
Table 4.1.3.3 Calculation of TF-IDF.....	27
Table 4.1.3.5 TF-IDF Results.....	29
Table 4.1.4.1 Training Data	31
Table 4.1.4.2 Training Label.....	32
Table 4.1.4.3 Testing Data	33
Table 4.1.4.4 Testing Label.....	33
Table 4.1.5.1 Bootstrap Indices.....	34
Table 4.1.5.2 OOB Indices.....	35
Table 4.1.5.3 Bootstrap Data.....	35
Table 4.1.5.4 Bootstrap Label	36
Table 4.1.5.5 OOB Data.....	37
Table 4.1.5.6 OOB Label	37
Table 4.1.5.7 Contents of Child in Feature IDX = 1 in Loop 1	39
Table 4.1.5.8 Contents of Child in Feature IDX = 1 in Loop 2	40
Table 4.1.5.9 Contentes of Child in Feature IDX = 0	40
Table 4.1.5.10 Terminal node 1	48
Table 4.1.5.11 Testing Data	49
Table 4.1.5.12 OOB Data.....	51
Table 4.1.5.13 Split Point from Node 1	51

Table 4.1.5.14 OOB Score from Node 1	51
Table 4.1.5.15 Terminal node 2	52
Table 4.1.5.16 Split Point from Node 2	54
Table 4.1.5.17 OOB Score from Node 2.....	54
Table 4.1.5.18 Terminal Node 3	55
Table 4.1.5.19 Split Point from Node 3	56
Table 4.1.5.20 OOB Score from Node 3.....	56
Table 4.1.5.21 Tree_Is.....	56
Table 4.2.1.1 Logistic Regression Dataset.....	58
Table 4.2.2.1 TF-IDF Logistic Results	60
Table 4.2.3.1 Positive Data	61
Table 4.2.3.2 Negative Data.....	62
Table 4.2.4.1 Word Dict.....	67
Table 5.2.1 Results of Random Forest.....	100
Table 5.2.2 Results of Logistic Regression.....	101
Table 5.3.1 Analysis of Random Forest.....	102
Table 5.3.2 Analysis of Logistic Regression.....	105