



PROJECT REPORT

HATE SPEECH PREDICTION USING K-MEANS ALGORITHM

LIM, ALEXANDRE NOVANDRYAN PRATAMA

17.K1.0010

**Faculty of Computer Science
Soegijapranata Catholic University
2021**



HALAMAN PENGESAHAN

Judul Tugas Akhir: : Hate Speech Prediction Using K-Means Algorithm
Diajukan oleh : Lim Alexandre Np
NIM : 17.K1.0010
Tanggal disetujui : 09 Juli 2021
Telah setuju oleh
Pembimbing : Hironimus Leong S.Kom., M.Kom.
Penguji 1 : Hironimus Leong S.Kom., M.Kom.
Penguji 2 : R. Setiawan Aji Nugroho S.T., MCompIT., Ph.D
Penguji 3 : Y.b. Dwi Setianto
Penguji 4 : Rosita Herawati S.T., M.I.T.
Penguji 5 : Yonathan Purbo Santosa S.Kom., M.Sc
Penguji 6 : Yulianto Tejo Putranto S.T., M.T.
Ketua Program Studi : Rosita Herawati S.T., M.I.T.
Dekan : R. Setiawan Aji Nugroho S.T., MCompIT., Ph.D

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

sintak.unika.ac.id/skripsi/verifikasi/?id=17.K1.0010

HALAMAN PERNYATAAN PUBLIKASI KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS

Yang bertanda tangan dibawah ini:

Nama : LIM, ALEXANDRE NOVANDRYAN PRATAMA

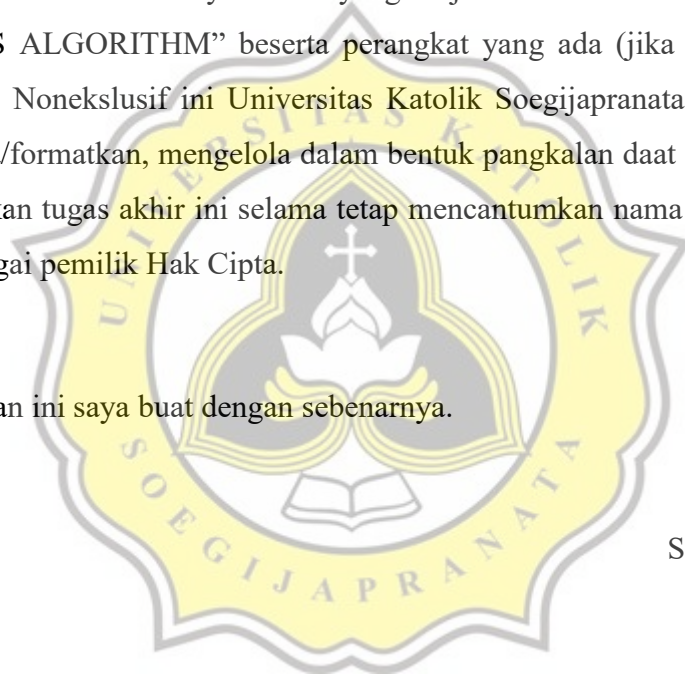
Program Studi : Teknik Informatika

Fakultas : Ilmu Komputer

Jenis Karya : Penelitian

Menyetujui untuk memberikan kepada Universitas Katolik Soegijapranata Semarang Hak Bebas Royalti Noneksklusif atas karya ilmiah yang berjudul “HATE SPEECH PREDICTION USING K-MEANS ALGORITHM” beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneksklusif ini Universitas Katolik Soegijapranata berhak menyimpan, mengalihkan media/formatkan, mengelola dalam bentuk pangkalan daat (*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, 9 Juli 2021

Yang menyatakan

LIM, ALEXANDRE NOVANDRYAN PRATAMA

DECLARATION OF AUTHORSHIP

I, the undersigned:

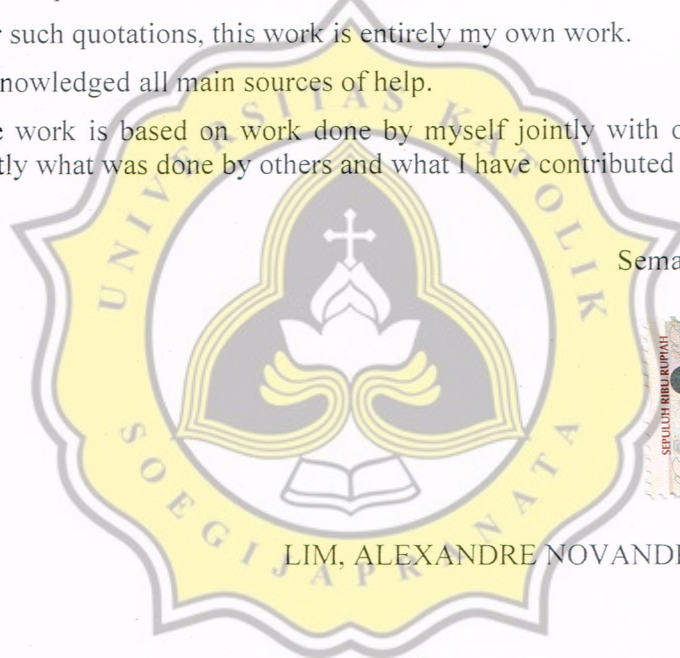
Name : LIM, ALEXANDRE NOVANDRYAN PRATAMA

ID : 17.K1.0010

declare that this work, titled "HATE SPEECH PREDICTION USING K-MEANS 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, July 9th, 2021



LIM, ALEXANDRE NOVANDRYAN PRATAMA

17.K1.0010

ACKNOWLEDGMENT

I have received a myriad of support, advice, and assistance throughout this document writing. I would like to thank my supervisor Hironimus Leong S.Kom., M.Kom for guiding and lead me to make this topic. I would also like to thank my friends the supports with advice to finish this document.

I would like to thank my family for giving me pep talk, patience, ceaseless love, support, and advices throughout my study in Soegijapranata Catholic University. You all gave me great escape to rest my mind from my thesis. And finally to the God Almighty Himself that bless, and allow me to finish this thesis.



Semarang, 9 Juli 2021

A handwritten signature in black ink, appearing to read 'Alexandre Novandryan Pratama'.

LIM, ALEXANDRE NOVANDRYAN PRATAMA

ABSTRACT

In the social media nowadays, there are lots of posts that shares a story about the user on something. Either it's to share a moment, or an opinion. So does freedom speech goes, not a few who abuse the freedom speech to take down others. Mostly those who abuse the freedom speech use hate speech to make the interlocutor feels uncomfortable.

This research is discusses about the usage of data mining algorithm and twitter data to predict hate speech in a post or a tweet. The K-Means that being used in this research is to define the Hate Speech in the dataset. The k is set to 2 to differentiate the first cluster is Hate Speech, and the second is Not Hate Speech.

The final results offered is in the form of a percentage of accuracy and comparison of the amount of data. From various comparison of data, the highest accuracy that being achieved is 80% followed by 66,7%, etc. However, by the results shown that different methods may varies different results in accuracy. But in overall the most stable results is by using N-Gram Tri-gram.

Keyword: K-Means, Hate Speech, Clustering, TF-IDF



TABLE OF CONTENTS

COVER.....	i
APPROVAL AND RATIFICATION PAGE.....	ii
DECLARATION OF AUTHORSHIP.....	iii
ACKNOWLEDGMENT.....	v
ABSTRACT.....	vi
TABLE OF CONTENTS.....	vii
LIST OF FIGURE.....	ix
LIST OF TABLE.....	x
CHAPTER 1 INTRODUCTION.....	1
1.1 Background.....	1
1.2 Problem Formulation.....	1
1.3 Scope.....	1
1.4 Objective.....	2
CHAPTER 2 LITERATURE STUDY.....	3
CHAPTER 3 RESEARCH METHODOLOGY.....	7
3.1 Journal.....	7
3.2 Data.....	7
3.3 Program Design.....	7
3.4 Coding.....	9
3.5 Implementation and Analysis.....	9
3.6 Conclusion and Report Writing.....	9
CHAPTER 4 ANALYSIS AND DESIGN.....	10
4.1 Analysis.....	10
4.1.1 Data Collecting.....	10
4.1.2 Processing Data.....	12
4.2 Design.....	20

CHAPTER 5 IMPLEMENTATION AND TESTING.....	21
5.1 Implementation.....	21
5.2 Testing.....	27
CHAPTER 6 CONCLUSION.....	31
REFERENCES.....	32
APPENDIX.....	b



LIST OF FIGURE

Figure 1.1. Process Flowchart.....	20
---	----



LIST OF TABLE

Table 1.1 : Data Table.....	11
Table 1.2 : Stopwords Table.....	12
Table 1.3 : TF-IDF Weight.....	17
Table 1.4 : Centroids.....	17
Table 1.5 : First Iteration.....	18
Table 1.6 : Second Iteration.....	19
Table 1.7 : Example Table.....	26
Table 1.8 : List Table.....	27
Table 1.9 : Accuracy Table.....	28
Table 1.10 : Effective Accuracy Table.....	29
Table 1.11 : N-Gram Results.....	30
Table 1.12 : Precision Recall Table.....	30

