

**WORD EMBEDDING PADA SENTIMEN ANALISIS UNTUK
MEREPRERENTASIKAN KATA MENJADI VEKTOR**

SKRIPSI



disusun oleh

**Dimas Midyan Alam
17.11.0982**

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2021**

**WORD EMBEDDING PADA SENTIMEN ANALISIS UNTUK
MEREPRERENTASIKAN KATA MENJADI VEKTOR**

SKRIPSI

Untuk memenuhi sebagian persyaratan
mencapai gelar Sarjana
pada Program Studi Informatika



disusun oleh

Dimas Midyan Alam

17.11.0982

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2021

PERSETUJUAN

SKRIPSI

WORD EMBEDDING PADA SENTIMEN ANALISIS UNTUK MEREPRESENTASIKAN KATA MENJADI VEKTOR

yang dipersiapkan dan disusun oleh

Dimas Midyan Alam

17.11.0982

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 29 Januari 2021

Dosen Pembimbing,

Mardhiya Hayaty, S.T., M.Kom.,

NIK : 190302108

PENGESAHAN**SKRIPSI****WORD EMBEDDING PADA SENTIMEN ANALISIS UNTUK
MEREPRESENTASIKAN KATA MENJADI VEKTOR**

yang dipersiapkan dan disusun oleh

Dimas Midyan Alam

17.11.0982

telah dipertahankan di depan Dewan Penguji

pada tanggal 17 Februari 2021

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Dina Maulina, M.Kom.
NIK : 190302250

Nuraini, M.Kom.
NIK : 190302066

Mardhiya Hayaty, S.T., M.Kom.
NIK. 190302108

Skripsi ini telah diterima sebagai salah satu persyaratan

untuk memperoleh gelar Sarjana Komputer

Tanggal 17 Februari 2021

DEKAN FAKULTAS ILMU KOMPUTER

Krisnawati, S.Si, M.T.
NIK. 190302038

PERNYATAAN

Saya yang bertanda tangan di bawah ini menyatakan bahwa skripsi ini merupakan karya saya sendiri (ASLI) dan isi pada skripsi ini tidak terdapat karya yang pernah diajukan oleh orang lain untuk memperoleh gelar akademis di suatu institusi pendidikan tinggi manapun, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis dan/atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar Pustaka.

Segala sesuatu yang terkait dengan masalah dan karya yang telah dibuat adalah menjadi tanggung jawab saya pribadi.

Yogyakarta, 5 Maret 2021



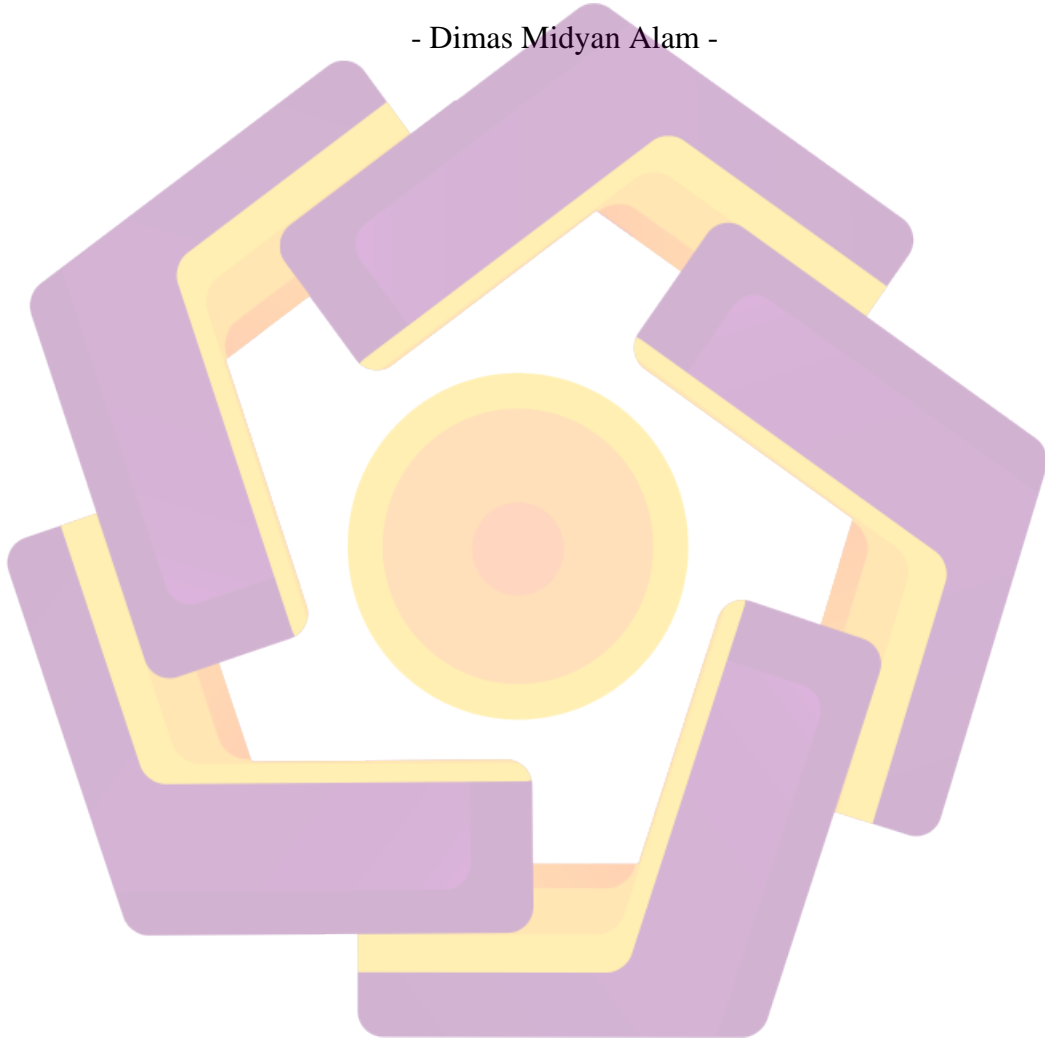
Dimas Midyan Alam

NIM 17.11.0982

MOTTO

“Setiap manusia memiliki jalan ceritanya masing-masing, jangan pernah samakan
jalan ceritamu dengan jalan cerita orang lain”

- Dimas Midyan Alam -



PERSEMBAHAN

Puji syukur saya panjatkan kepada ALLAH SWT atas berkah dan karunia yang telah diberikan, serta junjungan Nabi Besar kita Nabi Muhammad SAW, sehingga skripsi ini terselesaikan dengan baik dan lancar. Dengan ini saya persembahkan skripsi ini kepada semua pihak yang terlibat secara langsung maupun tidak langsung, yaitu kepada :

1. Kedua orang tua saya (Jeng Yanti dan Abdul Hamid), kedua kakak saya (Dewi dan Aby), bude saya (Bude Nuk), dan orang-orang disekitar saya yang tidak dapat disebutkan satu persatu. Selalu mendoakan saya, selalu mensupport saya dalam mengerjakan skripsi ini, dan selalu memberikan motivasi untuk maju terus pantang mundur.
2. Dosen pembimbing saya Ibu Mardhiya Hayaty, S.T., M.Kom., yang telah membimbing saya dari awal sampai akhir pembuatan skripsi.
3. Dosen - dosen Universitas Amikom Yogyakarta yang telah memberikan banyak ilmu pengetahuan dari semester awal hingga akhir.
4. Teman - teman khususnya Kelas 17-IF-01 yang telah menemani dan selalu memberikan banyak cerita setiap harinya dalam perkuliahan selama 6 semester.
5. Teman teman Uwuhnisty (Tedy, Seto, Anif, Sileh, Alwi, Bijis, Daniel, Manul, Agus, Farhan, Dito, Uwuh, Deri, dan Andri) yang sxesalu memberikan support kepada saya,

KATA PENGANTAR

Penulis panjatkan puji dan syukur kepada ALLAH SWT atas berkat dan rahmat-Nya, sehingga penulis dapat menyelesaikan skripsi dengan baik yang berjudul **“WORD EMBEDDING PADA SENTIMEN ANALISIS UNTUK MEREPRESENTASIKAN KATA MENJADI VEKTOR”** disusun sebagai salah satu syarat utama untuk menyelesaikan program sarjana pada Universitas AMIKOM Yogyakarta. Penyelesaian skripsi ini juga tidak lepas dari bantuan berbagai pihak, karena itu pada kesempatan ini penulis ingin menyampaikan rasa hormat dan terima kasih kepada :

1. Prof. Dr. M. Suyanto, MM. selaku Rektor Universitas AMIKOM Yogyakarta.
2. Ibu Krisnawati, S.Si, M.T. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
3. Bapak Sudarmawan, M.T. selaku Ketua Program Studi Informatika Universitas AMIKOM Yogyakarta.
4. Ibu Mardhiya Hayaty, S.T., M.Kom. selaku dosen pembimbing yang selalu meluangkan tenaga, waktu, dan pikirannya dalam membimbing saya.
5. Ibu Dina Maulina, M.Kom. dan Ibu Nuraini, M.Kom. selaku dosen penguji. Terimakasih atas segala saran yang diberikan selama pengujian untuk memperbaiki penelitian menjadi lebih baik lagi.

Penulis menyadari bahwa skripsi ini masih banyak kekurangan. Maka, penulis menerima segala kritik dan saran yang membangun dari semua pihak. Semoga skripsi ini bisa bermanfaat baik bagi penulis serta pembaca. Atas saran dan kritik, penulis ucapkan terima kasih.

Yogyakarta, 17 Februari 2021

Dimas Midyan Alam

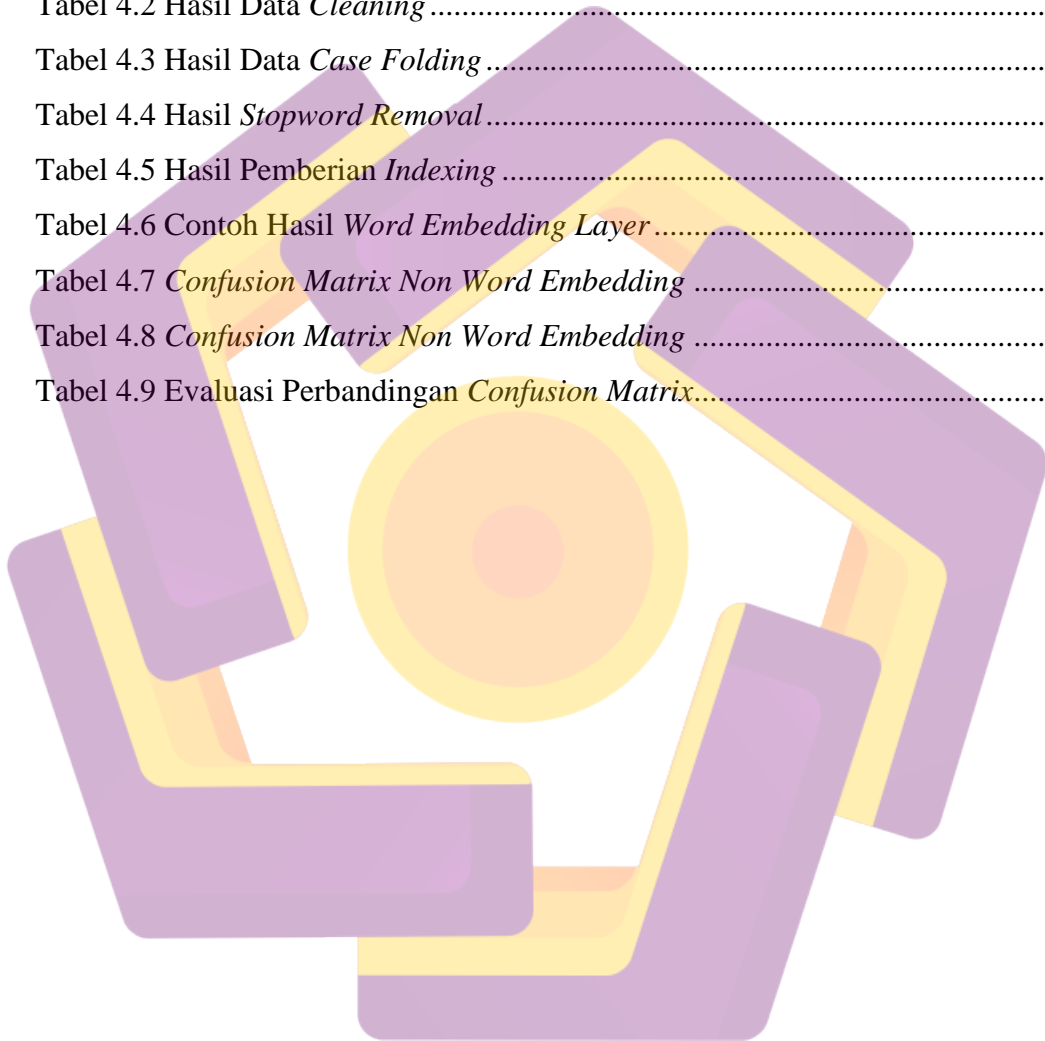
DAFTAR ISI

JUDUL.....	i
LEMBAR PERSETUJUAN	ii
LEMBAR PENGESAHAN	iii
PERNYATAAN	iv
MOTTO	v
PERSEMBAHAN.....	vi
KATA PENGANTAR	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	x
DAFTAR GAMBAR.....	xi
INTISARI	xii
ABSTRACT.....	xiii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah.....	3
1.3 Batasan Masalah	3
1.4 Tujuan Penelitian	3
1.5 Manfaat Penelitian	3
1.6 Metode Penelitian	4
1.7 Sistematika Penulisan	5
BAB II LANDASAN TEORI.....	7
2.1 Tinjauan Pustaka.....	7
2.2 Dasar Teori.....	10
2.2.1 <i>Natural Language Processing</i>	10
2.2.2 <i>Review</i>	10
2.2.3 <i>Data Mining</i>	11
2.2.4 <i>Text Mining</i>	12
2.2.5 <i>Deep Learning</i>	13
2.2.6 Analisis Sentimen.....	13
2.2.7 <i>Preprocessing Data</i>	14
2.2.8 <i>Word Embedding</i>	15
2.2.9 <i>Word2Vec CBOW (Continuous Bag-of-Word)</i>	15

2.2.10	<i>Long Short-Term Memory (LSTM)</i>	18
2.2.11	<i>Batch Size dan Epoch</i>	19
2.2.12	<i>Confusion Matrix</i>	20
BAB III METODOLOGI PENELITIAN		23
3.1	Tahapan Penelitian	23
3.2	Alat Penelitian	23
3.2.1	Perangkat Keras (<i>Hardware</i>)	24
3.2.2	Perangkat Lunak (<i>Software</i>)	24
3.3	Pengumpulan Data	24
3.4	<i>Annotation</i>	24
3.5	<i>Preprocessing Data</i>	25
3.6	Data Training dan Data Testing	26
3.7	Ekstraksi Fitur	26
3.8	Implementasi Algoritma Klasifikasi	27
3.9	Evaluasi	28
BAB IV IMPLEMENENTASI DAN PEMBAHASAN		29
4.1	Pengumpulan Data	29
4.2	Dataset	30
4.3	Preprocessing Data	31
4.3.1	<i>Data Cleaning</i>	32
4.3.2	<i>Case Folding</i>	33
4.3.3	<i>Stopword</i>	34
4.3.4	<i>Tokenizing</i>	35
4.4	Ekstraksi Fitur	37
4.5	Implementasi Algoritma LSTM	41
4.6	Hasil Training Data	44
4.7	Pengukuran Algoritma	46
BAB V PENUTUP		52
5.1	Kesimpulan	52
5.2	Saran	52
DAFTAR PUSTAKA		54

DAFTAR TABEL

Tabel 2.1 Perbandingan Penelitian.....	8
Tabel 2.2 <i>Confusion Matrix</i>	21
Tabel 4.1 Contoh Dataset Ulasan Pada <i>Marketplace</i>	31
Tabel 4.2 Hasil Data <i>Cleaning</i>	32
Tabel 4.3 Hasil Data <i>Case Folding</i>	33
Tabel 4.4 Hasil <i>Stopword Removal</i>	34
Tabel 4.5 Hasil Pemberian <i>Indexing</i>	37
Tabel 4.6 Contoh Hasil <i>Word Embedding Layer</i>	41
Tabel 4.7 <i>Confusion Matrix Non Word Embedding</i>	48
Tabel 4.8 <i>Confusion Matrix Non Word Embedding</i>	50
Tabel 4.9 Evaluasi Perbandingan <i>Confusion Matrix</i>	51



DAFTAR GAMBAR

Gambar 2.1 Arsitektur CBOW	17
Gambar 2.2 Pengulang RNN yang berisi satu layer	18
Gambar 2.3 Pengulang dalam LSTM berisi empat layer	19
Gambar 3.1 Diagram Alur Tahapan Penelitian	23
Gambar 4.1 Data Pada Website Kaggle	29
Gambar 4.2 Dataset	30
Gambar 4.3 Script Data Cleaning	32
Gambar 4.4 Script Case Folding	33
Gambar 4.5 Script Stopword Removal	34
Gambar 4.6 Script Tokenizing	35
Gambar 4.7 Script Membuat Data Latih dan Data Uji	36
Gambar 4.8 Hasil Vector Pad Sequences	38
Gambar 4.9 Script Generate Word2Vec	39
Gambar 4.10 Script Menampilkan Pendekatan Antar Kata	40
Gambar 4.11 Script Train Long Short Term Memory (LSTM)	42
Gambar 4.12 Total Parameter	43
Gambar 4.13 Architecture Model Long Sort Term Memory (LSTM)	44
Gambar 4.14 Script Fit Network Long Sort Term Memory (LSTM)	44
Gambar 4.15 Hasil Training Dan Validasi Non Embedding	45
Gambar 4.16 Hasil Training Dan Validasi Word Embedding	45
Gambar 4.17 Script Evaluasi Model	46
Gambar 4.18 Script Prediksi Model	47
Gambar 4.19 Script Confusion Matrix	47
Gambar 4.20 Hasil Confusion Matrix Non Word Embedding	48
Gambar 4.21 Hasil Confusion Matrix Word Embedding	49

INTISARI

Analisis sentimen adalah sebuah proses yang memahami, mengekstraksi, dan mengolah data teks secara otomatis untuk menemukan jenis sentimen pada teks tersebut, merupakan salah satu solusi mengatasi masalah untuk mengelompokan opini atau review menjadi opini positif atau negatif secara otomatis Analisis sentimen berguna untuk memudahkan pengguna pada proses memahami sentimen sehingga dapat melakukan penentuan keputusan pada suatu objek.

Pada penelitian ini menerapkan Word Embedding dengan Metode Word2Vec dan menggunakan algoritma klasifikasi Long Short Term Memory (LSTM). Metode Word2Vec digunakan untuk merepresentasikan kata kata dalam bentuk matematis. Word2Vec merupakan sebuah algoritma untuk mempelajari posisi kedekatan semantic antar kata dari sebuah teks masukan. Penelitian ini bertujuan untuk mengetahui berapa hasil akurasi *Word Embedding* dengan metode *Word2Vec* dan tanpa menggunakan *Word Embedding*.

Berdasarkan penelitian yang dilakukan dengan pembagian data *training* 80% dan data *testing* 20% dari total jumlah dataset 10022, diperoleh nilai akurasi analisis sentimen yang menggunakan word embedding adalah 0,92 atau 92%, sedangkan nilai akurasi tanpa menggunakan word embedding adalah 0,88 atau 88%.

Kata-kunci: Sentimen Analisis, Word2Vec, Word Embedding, LSTM, Cosine Similarity

ABSTRACT

Sentiment analysis is a process that understands, extracts, and processes text data automatically to find the type of sentiment in the text, which is a solution to the problem of grouping opinions or reviews into positive or negative opinions automatically. Sentiment analysis is useful for making it easier for users to process. Understand sentiment so that you can make decisions on an object.

This research applies Word Embedding with the Word2Vec method and uses the Long Short Term Memory (LSTM) classification algorithm. The Word2Vec method is used to represent words in mathematical form. Word2Vec is an algorithm for learning semantic closeness between words from an input text. This study aims to see the accuracy results of Word Embedding using the Word2Vec method and without using Word Embedding.

Based on research conducted by sharing 80% training data and testing data 20% of the total number of 10022 datasets, the value of sentiment analysis using word embedding was 0.92 or 92%, while the value without using word embedding was 0.88 or 88%.

Keywords: Sentimen Analisis, Word2Vec, Word Embedding, LSTM, Cosine Similarity