

**REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI
INDONESIA MENGGUNAKAN METODE WORD2VEC**

SKRIPSI



Disusun oleh

Andrian Tri Muryanto

17.11.0977

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

2020

REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI

INDONESIA MENGGUNAKAN METODE WORD2VEC

SKRIPSI

Untuk memenuhi sebagian persyaratan

Mencapai gelar Sarjana

pada Program Studi Informatika



Disusun oleh

Andrian Tri Muryanto

17.11.0977

PROGRAM SARJANA

PROGRAM STUDI INFORMATIKA

FAKULTAS ILMU KOMPUTER

UNIVERSITAS AMIKOM YOGYAKARTA

YOGYAKARTA

2020

PERSETUJUAN

SKRIPSI

REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI INDONESIA MENGGUNAKAN METODE WORD2VEC

yang dipersiapkan dan disusun oleh

Andrian Tri Muryanto

17.11.0977

telah disetujui oleh Dosen Pembimbing Skripsi

pada tanggal 26 Juni 2020

Dosen Pembimbing,

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

NIK : 190302108

PENGESAHAN

SKRIPSI

REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI INDONESIA MENGGUNAKAN METODE WORD2VEC

yang dipersiapkan dan disusun oleh

Andrian Tri Muryanto

17.11.0977

telah dipertahankan di depan Dewan Penguji

pada tanggal 18 November 2020

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Ainul Yaqin, M. Kom.

NIK : 190302098

M. Rudyanto Arief, S.T, M.T.

NIK : 190302098

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

NIK. 190302108

Skripsi ini telah diterima sebagai salah satu persyaratan

untuk memperoleh gelar Sarjana Komputer

Tanggal

DEKAN FAKULTAS ILMU KOMPUTER

Krisnawati, S.Si, M.T.

NIK. 190302038

PERNYATAAN

Saya yang bertandatangan 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 tanggungjawab saya pribadi.

Yogyakarta, 21 November 2020



Andrian Tri Muryanto

NIM 17.11.0977

MOTTO

“What will people say, and I don’t care about that”

- deddy corbuzier -

“Biarkan orang lain menganggapmu rendah, karena mereka hanya memandang dirimu dari covernya saja”

- Andrian Tri Muryanto -



PERSEMBAHAN

Puji syukur saya panjatkan kepada Allah SWT atas berkah dan karunia-NYA 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 dan kakak kakak saya yang selalu mendoakan dan selalu mensupport saya dalam mengerjakan skripsi ini, dan selalu memberikan motivasi untuk maju terus.
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 ilmu dari semester awal sampai akhir selama kuliah
4. Teman teman khususnya Kelas 17-IF-01 yang telah menemani dan selalu memberikan semangat juang dalam mengerjakan skripsi ini
5. Teman teman SRAWUNG CLUB (Damar Komir, Aldo bill bae, Dhimskut, Lutfi edan, Fariz paryono, Rizky Roy, Bappocx POY, Dalipon, Padhe Dedidot, Alif Part-min, Anjasopo anjay, Vina mbak pina pathok, IFA ipak ipol, RIA Riaw, Puspa jaelani, Ema Lia) yang sesalu memberikan support kepada saya.
6. Teman teman Uwuhnisty (Tedi beer, Setoin, Bangdim, Anip, Hendra Sileh, Alwi, Bijis, Daniel, A.K.A Manul, Farhan Koton, Deri OC) yang sesalu memberikan support kepada saya,

KATA PENGANTAR

Puji syukur kehadiran Allah SWT atas segala rahmat dan hidayah-NYA yang telah diberikan kepada peneliti sehingga dapat menyelesaikan penelitian ini. Dan tidak lupa kita panjatkan shalawat seta salam kepada junjungan nabi besar kita Nabi Muhammad SAW, yang telah menjadi suri tauladan yang baik bagi umatnya dan untuk berbuat kebijakan.

Skripsi yang berjudul **“REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI INDONESIA MENGGUNAKAN METODE WORD2VEC”** ini disusun sebagai salah satu syarat kelulusan bagi setiap mahasiswa Universitas AMIKOM Yogyakarta. Selain itu juga merupakan suatu bukti bahwa mahasiswa telah menyelesaikan kuliah jenjang program Strata-1 dan untuk memperoleh gelar Sarjana Komputer.

Penyelesaian skripsi ini juga tidak lepas dari bantuan berbagai pihak, karena itu pada kesempatan ini penulis menyampaikan rasa hormat dan terimakasih 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 bijaksana memberikan bimbingan, nasehat serta waktunya selama penulisan skripsi ini.
5. Bapak Ainul Yaqin, M. Kom dan Bapak M. Rudyanto Arief, S.T, M.T selaku dosen penguji, terimakasih atas saran yang diberikan selama pengujian untuk memperbaiki penelitian ini menjadi lebih baik lagi.

Peneliti menyadari bahwa pembuatan skripsi ini masih banyak kekurangan dan kelemahannya. Oleh karena itu peneliti berharap kepada semua pihak agar dapat menyampaikan kritik dan saran yang membangun untuk menambah kesempurnaan skripsi ini. Namun peneliti tetap berharap skripsi ini akan bermanfaat bafi semua pihak yang membacanya. Apa bila terdapat kesalahan semoga Allah SWT melimpahkan magfirah-NYA. *Aamiin yaa Kholi.*

Yogyakarta, 21 November 2020

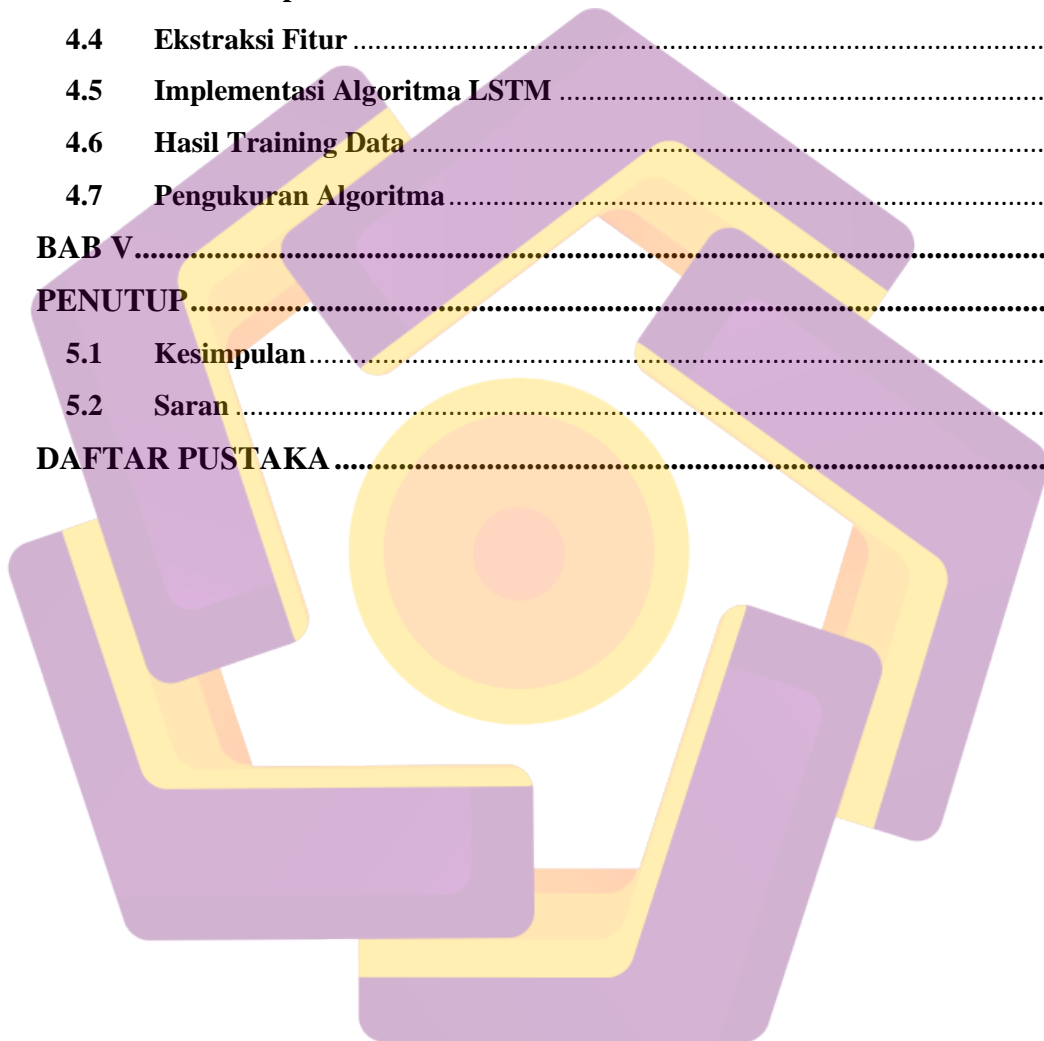
Andrian Tri Muryanto

DAFTAR ISI

REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI INDONESIA MENGGUNAKAN METODE WORD2VEC	1
REPRESENTASI KATA MENJADI VECTOR PADA ULASAN HOTEL DI INDONESIA MENGGUNAKAN METODE WORD2VEC	1
PERSETUJUAN.....	i
PERNYATAAN.....	iii
MOTTO	iv
PERSEMBAHAN.....	v
KATA PENGANTAR.....	vi
DAFTAR ISI.....	viii
DAFTAR TABEL	xi
DAFTAR GAMBAR.....	xii
INTISARI	xiii
ABSTRACT	xiv
BAB I.....	1
PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah	3
1.3 Batasan Masalah	3
1.4 Maksud dan Tujuan Penelitian	4
1.5 Metode Penelitian.....	4
1.5.1 Pengumpulan Data.....	4
1.5.2 Annotation	5
1.5.3 Preprocessing Data	6
1.5.4 Ekstraksi Fitur	6
1.5.5 Implementasi Algoritma Klasifikasi.....	6
1.5.6 Evaluasi.....	7
1.6 Sistematika Penulisan.....	7
BAB II	9
LANDASAN TEORI.....	9

2.1	Tinjauan Pustaka	9
2.2	Dasar Teori	12
2.2.1	<i>Natural Language Processing</i>	12
2.2.2	<i>Data Mining</i>	12
2.2.3	<i>Text Mining</i>	13
2.2.4	<i>Deep Learning</i>	15
2.2.5	Analisis Sentimen	15
2.2.6	Review	16
2.2.7	<i>Word Embedding</i>	17
2.2.8	<i>Word2Vec CBOW (continous bag-of-word)</i>	17
2.2.9	<i>Scraping</i>	19
2.2.10	<i>Preprocessing Data</i>	20
2.2.11	<i>Long Short Term Memory</i>	21
2.2.12	Batch Size dan Epoch	24
2.2.13	<i>Cofusion Matrix</i>	25
2.2.14	<i>Indexing</i>	26
2.2.15	<i>Optimasi Adam</i>	27
BAB III	29
METODOLOGI PENELITIAN	29
3.1	Tahapan Penelitian	29
3.2	Alat Penelitian	29
3.2.1	Perangkat Keras (Hardware)	29
3.2.2	Perangkat Lunak (Software).....	30
3.3	Pengumpulan Data.....	30
3.4	Annotation	30
3.5	Preprocessing Data	31
3.6	Data Training dan Data Testing	32
3.7	Ekstraksi Fitur	33
3.8	Implementasi Algoritma Klasifikasi.....	33
3.9	Evaluasi	34
BAB IV	35
IMPLEMENTASI DAN PEMBAHASAN	35
4.1	Pengumpulan Data.....	35

4.2	Dataset	36
4.3	Preprocessing Data	38
4.3.1	Data Cleaning	38
4.3.2	Case Folding	39
4.3.3	Tokenizing	39
4.3.4	Stopword Removal	40
4.4	Ekstraksi Fitur	43
4.5	Implementasi Algoritma LSTM	47
4.6	Hasil Training Data	50
4.7	Pengukuran Algoritma	52
BAB V	57
PENUTUP	57
5.1	Kesimpulan	57
5.2	Saran	58
DAFTAR PUSTAKA	59



DAFTAR TABEL

Tabel 2. 1 Perbandingan Penelitian.....	10
Tabel 2. 2 Confusion Matrix	25
Tabel 4. 1 Contoh dataset ulasan hotel di Indonesia	37
Tabel 4. 2 Sebelum preprocessing dan sesudah preprocessing	41
Tabel 4. 3 Hasil Pemberian indexing.....	43
Tabel 4. 4 Hasil Pad Sequences.....	44
Tabel 4. 5 Contoh Hasil Word Embedding layer.....	46
Tabel 4. 6 Confusion Matrix non Word Embedding	54
Tabel 4. 7 Confusion Matrix Word Embedding.....	54
Tabel 4. 8 Evaluasi Perbandingan Confusion Matrix.....	55

DAFTAR GAMBAR

Gambar 2. 1 Arsitektur CBOW.....	18
Gambar 2. 2Pengulang RNN yang berisi satu layer	22
Gambar 2. 3 Pengulang dalam LSTM berisi empat Layer.....	22
Gambar 3. 1 Diagram Alur Tahapan Penelitian.....	29
Gambar 4. 1 Scraping data menggunakan WebHarvy	36
Gambar 4. 2 Dataset	37
Gambar 4. 3 Script Data Cleaning	39
Gambar 4. 4Script Case Folding	39
Gambar 4. 5 Script Tokenizing.....	40
Gambar 4. 6 Script Stopword Removal.....	41
Gambar 4. 7 Script Pembagian Dataset.....	42
Gambar 4. 8 Script Generate Word2Vec.....	45
Gambar 4. 9 Script menampilkan pendekatan antar kata	46
Gambar 4. 10 Script Train LSTM.....	48
Gambar 4. 11 Total parameter	49
Gambar 4. 12 Architecture Model LSTM.....	49
Gambar 4. 13 Script fit network LSTM.....	50
Gambar 4. 14 hasil training dan Validasi non embedding.....	51
Gambar 4. 15 Hasil training dan Validasi Embedding	51
Gambar 4. 16 Script evaluasi Model	52
Gambar 4. 17 Script Prediksi Model.....	53
Gambar 4. 18 Script Confusion Matrix.....	53

INTISARI

Hotel merupakan salah satu produk pariwisata yang sangat penting untuk dipertimbangkan baik dari segi fasilitas dan pelayanan. Saat ini sudah banyak website wisata yang menyediakan fasilitas internet untuk menuliskan opini dan pengalaman pribadinya secara online. Analisa sentimen atau opinion mining merupakan salah satu solusi mengatasi masalah untuk mengelompokkan opini atau review menjadi opini positif atau negatif secara otomatis.

Analisis sentimen adalah sebuah proses yang memahami, mengekstraksi, dan mengolah data teks secara otomatis untuk menemukan jenis sentimen pada teks tersebut. Analisis sentimen berguna untuk memudahkan pengguna pada proses memahami sentiment sehingga dapat melakukan penentuan keputusan pada suatu objek.

Penelitian ini menerapkan Metode Word2Vec untuk dilakukannya kinerja sentiment analisis. Metode Word2Vec dapat digunakan untuk merepresentasikan kata kata dalam bentuk matematis. Word2Vec merupakan sebuah algoritma untuk mempelajari posisi kedekatan semantic antar kata dari sebuah teks masukan. Word2vec memiliki dua arsitektur yaitu Skip-gram dan CBOW .

Berdasarkan hasil penelitian diperoleh nilai akurasi analisis sentimen untuk ulasan hotel di Indonesia yang menggunakan word embedding adalah 0,90 atau 90%, sedangkan nilai akurasi dari tanpa word embedding adalah 0,88 atau 88%. Namun saat menggunakan word embedding, kecepatan epoch membutuhkan waktu proses yang lebih lama. Jika tidak menggunakan word embedding, kecepatan pemrosesan akan lebih cepat.

Kata Kunci : Analisis Sentimen, *Word Embedding*, *Word2Vec*, *LSTM*

ABSTRACT

Hotels are one of the tourism products that are very important to consider both in terms of facilities and services. Currently, there are many tourist websites that provide internet facilities to write opinions and personal experiences online. Sentiment analysis or opinion mining is one solution to overcome problems to group opinions or reviews into positive or negative opinions automatically.

Sentiment analysis is a process that understands, extracts, and processes text data automatically to find the type of sentiment in the text. Sentiment analysis is useful for making it easier for users in the process of understanding sentiment so that they can make decisions on an object.

This study applies the Word2Vec method to perform sentiment analysis performance. The Word2Vec method can be used to represent words in mathematical form. Word2Vec is an algorithm for studying semantic proximity positions between words from an input text. Word2vec has two architectures namely Skip-gram and CBOW.

Based on the results of the study, the accuracy value of sentiment analysis for hotel reviews in Indonesia that uses word embedding is 0.90 or 90%, while the accuracy value for no word embedding is 0.88 or 88%.

However, when using word embedding, epoch speed requires a longer processing time. If you don't use word embedding, the processing speed will be faster.

Keywords: *Sentiment Analysis, Word Embedding, Word2Vec, LSTM*