

**WORD EMBEDDING-BERT UNTUK EKSTRAKSI FITUR
PADA ANALISIS SENTIMEN**

SKRIPSI



disusun oleh
Aditya Rahman
17.11.1639

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

**WORD EMBEDDING-BERT UNTUK EKSTRAKSI FITUR
PADA ANALISIS SENTIMEN**

SKRIPSI

Untuk memenuhi sebagian persyaratan
Mencapai gelar Sarjana
pada Program Studi Informatika



disusun oleh
Aditya Rahman
17.11.1639

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

PERSETUJUAN

SKRIPSI

**WORD EMBEDDING-BERT UNTUK EKSTRAKI FITUR
PADA ANALISIS SENTIMEN**

yang dipersiapkan dan disusun oleh

Aditiya Rahman

17.11.1639

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 5 Maret 2021

Dosen Pembimbing,



Mardhiva Havaty, S.T., M.Kom.,

NIK : 190302108

PENGESAHAN

SKRIPSI

**WORD EMBEDDING-BERT UNTUK EKSTRAKI FITUR
PADA ANALISIS SENTIMEN**

yang dipersiapkan dan disusun oleh

Aditiya Rahman

17.11.1639

telah dipertahankan di depan Dewan Penguji
pada tanggal 18 Februari 2021

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Uyock Anggoro Saputro, M.Kom.

NIK : 190302419

Arif Dwi Laksito, M.Kom.

NIK : 190302150

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

NIK. 190302108

Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer

Tanggal 18 Februari 2021

DEKAN FAKULTAS ILMU KOMPUTER

Krisnawati, S.Si, M.T.

NIK. 190302038

PERNYATAAN

Saya yang bertandatangan dibawah ini menyatakan bahwa, skripsi ini merupakan karya saya sendiri (ASLI), dan isi dalam 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 naskah dan karya yang telah dibuat adalah menjadi tanggungjawab saya pribadi.

Pangkalpinang, 5 Maret 2021



Aditiya Rahman

NIM. 17.11.1639

MOTTO

“Dreaming as you will live forever. Live as you will die today.”

- James Dean

“You have limited time. Don't waste it by living someone else's life.”

- Steve Jobs

“Pleasure in a work makes perfection at the results achieved.”

- Aristoteles

“The only way to do great work is to love what you do. If you haven't found it yet, keep looking. Don't settle.”

- Steve Jobs



PERSEMBAHAN

Alhamdulillah dengan kerja keras serta doa, skripsi ini dapat diselesaikan dengan baik. Segala puji dan syukur bagi Allah SWT yang tiada henti memberikan keberkahan. Dengan ini saya mempersembahkan skripsi ini kepada semua pihak yang terlibat secara langsung atau tidak langsung, yaitu untuk :

1. Kedua orang tua dan kakak – kakak saya yang selalu mendoakan, selalu menyemangati dan memberikan motivasi tiada henti kepada saya.
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 selama masa kuliah.
4. Teman-teman satu kloter penelitian yang memberikan dukungan, bantuan dan motivasi, serta menemani saya untuk menyelesaikan skripsi ini.
5. Teman-teman kelas 17-IF-11 yang telah menemani dan selalu memberikan semangat untuk menyelesaikan skripsi ini.

KATA PENGANTAR

Puji syukur penulis panjatkan kehadirat Allah SWT yang telah melimpahkan rahmat serta hidayah-Nya dan shawat serta salam juga tidak lupa penulis panjatkan kepada junjungan kita Nabi Muhammad SAW yang telah memberikan teladan mulia dalam menuntun umatnya sehingga penulis dapat menyelesaikan skripsi ini.

Skripsi yang berjudul **“WORD EMBEDDING BERT UNTUK EKSTRAKSI FITUR PADA ANALISIS SENTIMEN“** ini 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 bijaksana memberikan bimbingan, nasehat serta waktunya selama penulisan skripsi ini.
5. Bapak Uyock Anggoro Saputro, M.Kom. dan Bapak Arif Dwi Laksito, M.Kom. selaku dosen penguji. Terimakasih atas saran yang diberikan

selama pengujian untuk memperbaiki penelitian menjadi lebih baik lagi.

Penulis menyadari skripsi ini masih ada kekurangan. Maka, penulis menerima kritik dan saran yang membangun serta teguran dari semua pihak. Penulis menerima dengan lapang dada untuk kesempurnaan karya selanjutnya. Semoga skripsi yang sederhana ini bisa bermanfaat, khususnya bagi penulis dan pembaca yang budiman pada umumnya. Apabila terdapat kesalahan semoga Allah SWT melimpahkan magfirah-Nya. *Aamiin.*

Pangkalpinang, 5 Maret 2021



Aditiya Rahman

DAFTAR ISI

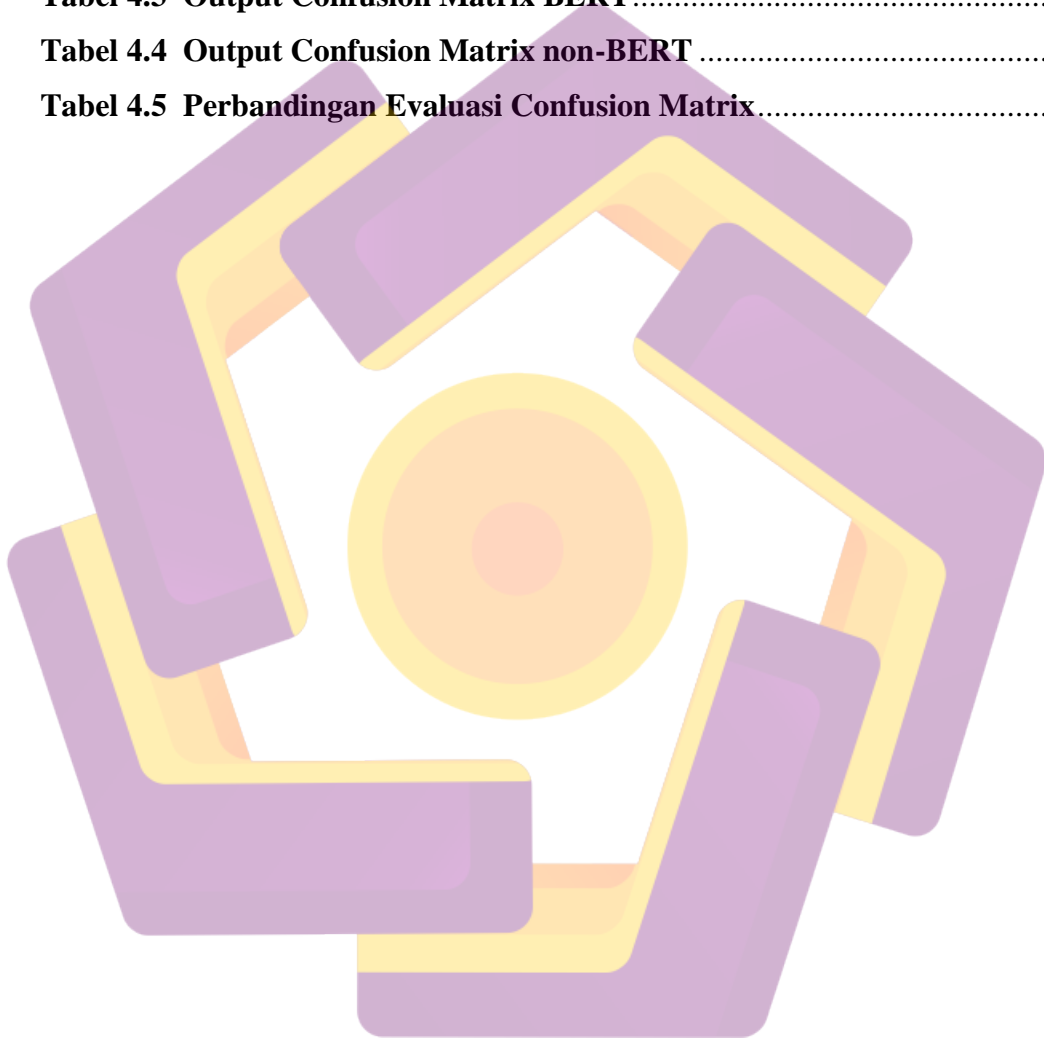
JUDUL.....	i
PERSETUJUAN.....	iii
PENGESAHAN.....	iv
PERNYATAAN.....	v
MOTTO.....	vi
PERSEMBAHAN.....	vii
KATA PENGANTAR.....	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR.....	xiv
INTISARI.....	xvi
ABSTRACT.....	xvii
BAB I.....	1
PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	2
1.3 Batasan Masalah.....	3
1.4 Maksud dan Tujuan Penelitian.....	3
1.5 Manfaat Penelitian.....	3
1.6 Metode Penelitian.....	3
1.6.1 Pengumpulan Data.....	3
1.6.2 Manual Labeling.....	4
1.6.3 Preprocessing Data.....	4
1.6.4 Ekstraksi Fitur.....	4

1.6.5	Evaluasi	5
1.7	Sistematika Penulisan	5
BAB II	7
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>Deep Learning</i>	10
2.2.3	<i>Data Mining</i>	10
2.2.4	<i>Text Mining</i>	11
2.2.5	Analisis Sentimen	13
2.2.6	<i>Review</i>	15
2.2.7	<i>Word Embedding</i>	15
2.2.8	<i>Vector Space Model</i>	19
2.2.9	<i>Indexing</i>	21
2.2.10	<i>Bidirectional Encoder Representations from Transformers (BERT)</i> 22	
2.2.11	<i>Linear Regression</i>	28
2.2.12	<i>Preprocessing Data</i>	29
2.2.12	Optimasi Adam	31
2.2.13	<i>Batch Size dan Epoch</i>	32
2.2.14	<i>Confusion Matrix</i>	33
BAB III	35
METODELOGI PENELITIAN	35
3.1	Tahap Penelitian	35
3.2	Alat Penelitian	35

3.2.1	Perangkat Keras (<i>Hardware</i>)	35
3.2.2	Perangkat Lunak (<i>Software</i>).....	36
3.3	Pengumpulan Data.....	36
3.4	Manual Labeling	37
3.5	<i>Preprocessing</i> Data	37
3.6	Data <i>Training</i> , Data <i>Valid</i> dan Data <i>Testing</i>	37
3.7	Ekstraksi Fitur	38
3.8	Implementasi Algoritma Klasifikasi.....	38
3.9	Evaluasi	39
BAB IV	40
IMPLEMENTASI DAN PEMBAHASAN	40
4.1	Pengumpulan Data.....	40
4.2	<i>Dataset</i>	41
4.3	Persiapan	42
4.4	Mempersiapkan Data <i>Pre-Train</i>	43
4.5	<i>Preprocessing</i> Data	44
4.5.1	<i>Case Folding</i>	44
4.5.2	<i>Tokenizing</i>	45
4.6	Permodelan.....	48
4.7	<i>Testing Model</i>	52
4.7	Pengukuran Akurasi <i>Testing</i>	55
BAB V	64
PENUTUP	64
5.1	Kesimpulan.....	64
5.2	Saran.....	64
DAFTAR PUSTAKA	65

DAFTAR TABEL

Tabel 2.1	Tabel Perbandingan.....	9
Tabel 2.2	Confusion matrix.....	33
Tabel 4.1	Contoh Ulasan hotel positif dan negatif	42
Tabel 4.2	Hasil Case Folding.....	45
Tabel 4.3	Output Confusion Matrix BERT.....	61
Tabel 4.4	Output Confusion Matrix non-BERT	62
Tabel 4.5	Perbandingan Evaluasi Confusion Matrix.....	63



DAFTAR GAMBAR

Gambar 2.1 Alur proses text mining	13
Gambar 2.2 Ilustrasi mapping sentimen word embedding	17
Gambar 2.3 Word embedding	18
Gambar 2.4 Representasi Model Ruang Vektor	19
Gambar 2.5 Representasi nilai kata dalam dokumen	20
Gambar 2.6 Visualisasi sederhana kinerja encoder dan decoder	23
Gambar 2.7 Ilustrasi proses fine-tuning pada BERT	25
Gambar 2.8 Representasi input pada BERT	25
Gambar 2.9 Visualisasi proses fine-tuning pada BERT	28
Gambar 3.1 Diagram Alir Tahapan Penelitian	35
Gambar 4.1 Scraping data menggunakan WebHarvy	40
Gambar 4.2 Dataset	41
Gambar 4.3 Perintah untuk memasang library transformers	42
Gambar 4.4 Perintah untuk cloning repository IndoNLU	43
Gambar 4.5 Kode untuk mengimpor library dan function	43
Gambar 4.6 Kode untuk mempersiapkan data pre-train	43
Gambar 4.7 Script Case Folding	44
Gambar 4.8 Kode untuk memeriksa kinerja class tokenizer	45
Gambar 4.9 Hasil keluaran proses tokenizing	46
Gambar 4.10 Deklarasi masing – masing file dataset	47
Gambar 4.11 Tahap tokenizing pada dataset	47
Gambar 4.12 Memberi konfigurasi tambahan pada dataset	48
Gambar 4.13 Optimizer Adam dan pengiriman model ke GPU	49
Gambar 4.14 Kode untuk menjalankan training	49
Gambar 4.15 Kode untuk evaluasi valid set	50
Gambar 4.16 Output yang ditampilkan ke layar pada saat training	50
Gambar 4.17 Visualisasi proses training dan evaluasi BERT	51
Gambar 4.18 Visualisasi proses training dan evaluasi non-BERT	52
Gambar 4.19 Kode untuk melakukan testing	52
Gambar 4.20 Fungsi forward_sequence_classification	53
Gambar 4.21 Isi variabel logits	54

Gambar 4.22 Output fungsi torch.topk()	54
Gambar 4.23 Menyalin label dari file pred.txt ke dalam array	55
Gambar 4.24 Perhitungan akurasi, confusion matrix, dan classification report	56
Gambar 4.25 Keluaran akurasi, classification report, dan grafik confusion matrix BERT	56
Gambar 4.26 Keluaran akurasi, classification report, dan grafik confusion matrix BERT dengan tahapan preprocessing lengkap	58
Gambar 4.27 Keluaran akurasi, classification report, dan grafik confusion matrix non-BERT.....	59
Gambar 4.28 Kode perbandingan label file asli dengan labeling BERT	60
Gambar 4.29 Grafik perbandingan labeling.....	60



INTISARI

Perkembangan teknologi informasi saat ini telah mempermudah kehidupan manusia dalam berbagai bidang, salah satunya di bidang pariwisata yaitu pemesanan hotel. Kita bisa mendapatkan banyak informasi mengenai hotel tersebut mulai dari fasilitas yang diberikan hingga ulasan oleh pengunjung yang sudah pernah menggunakan hotel tersebut. Analisis sentimen merupakan salah satu cabang dari *Natural Language Processing* (NLP) yang dapat membantu dalam menentukan kualitas layanan hotel yang ditawarkan dari ulasan yang telah diberikan pengguna. Penelitian ini menggunakan data ulasan hotel untuk melakukan analisis sentimen yang didapatkan dari situs Traveloka.

Penelitian ini memanfaatkan sebuah metode *deep learning* yaitu *Bidirectional Encoder Representation from Transformer* atau BERT sebagai metode *word embedding* untuk mempresentasikan kata menjadi vektor. Klasifikasi pada penelitian ini dilakukan dengan menambahkan layer *linear regression* pada layer paling atas di BERT.

Dari percobaan yang dilakukan dengan pembagian data *training* sebanyak 70%, data eval sebanyak 10%, dan data tes sebanyak 20% dari total 10.000 data, dapat dilihat bahwa metode BERT memiliki hasil akurasi sebesar 97%, sedangkan metode non-BERT menghasilkan akurasi sebesar 93%. Hasil yang diberikan oleh metode BERT lebih tinggi daripada hasil yang diberikan oleh metode non-BERT. Sehingga bisa disimpulkan bahwa ada pengaruh dari penggunaan metode BERT terhadap akurasi yang dihasilkan.

Kata-kunci: Analisis Sentimen, *Word Embedding*, *Deep Learning*, BERT.

ABSTRACT

The current development of information technology has made human life easier in various fields, one of which is tourism, namely hotel reservations. We can get a lot of information about the hotel from the facilities provided to reviews by visitors who have used the hotel. Sentiment analysis is a branch of Natural Language Processing (NLP) that can help determine the quality of hotel services offered from reviews that users have provided. This study uses hotel review data to analyze sentiments obtained from the Traveloka website.

This study utilizes a deep learning method, namely Bidirectional Encoder Representation from Transformer or BERT as a word embedding method to present words as vectors. The classification in this study was carried out by adding a linear regression layer to the top layer of BERT.

From experiments conducted with the distribution of training data as much as 70%, evaluation data as much as 10%, and test data as much as 20% of the total 10,000 data, it can be seen that the BERT method has an accuracy of 97%, while the non-BERT method produces an accuracy of 93 %. The results given by the BERT method are higher than the results given by the non-BERT methods. So it can be concluded that there is an effect of using the BERT method on the resulting accuracy.

Keywords: *Sentiment Analysis, Word Embedding, Deep Learning, BERT.*

