

**ANALISIS PENGARUH DATA NORMALISASI DAN NON
NORMALISASI PADA DETEKSI *CYBERBULLYING*
BERBAHASA INDONESIA DENGAN METODE
INDOBERT**

TUGAS AKHIR

Oleh:

MITA YANI NURMA PRATIWI

NIM. 1910817320003



PROGRAM STUDI TEKNOLOGI INFORMASI

FAKULTAS TEKNIK

UNIVERSITAS LAMBUNG MANGKURAT

BANJARMASIN

2023

**ANALISIS PENGARUH DATA NORMALISASI DAN NON
NORMALISASI PADA DETEKSI CYBERBULLYING BERBAHASA
INDONESIA DENGAN METODE INDOBERT**

TUGAS AKHIR

Diajukan Untuk Memenuhi Salah Satu Syarat
Sarjana Strata-1 Teknologi Informasi

Oleh:

MITA YANI NURMA PRATIWI

NIM. 1910817320003



PROGRAM STUDI TEKNOLOGI INFORMASI

FAKULTAS TEKNIK

UNIVERSITAS LAMBUNG MANGKURAT

BANJARMASIN

2023

LEMBAR PERNYATAAN

Yang bertanda tangan di bawah ini:

Nama : Mita Yani Nurma Pratiwi
NIM : 1910817320003
Fakultas : Teknik
Program Studi : Teknologi Informasi
Judul Tugas Akhir : Analisis Pengaruh Data Normalisasi dan Non
Normalisasi pada Deteksi *Cyberbullying*
Berbahasa Indonesia dengan Metode
IndoBERT
Pembimbing Utama : Dr. Ir. Yuslena Sari, S.Kom., M.Kom.

Dengan ini saya menyatakan bahwa dalam Tugas Akhir ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar akademik di suatu perguruan tinggi, dan sepanjang pengetahuan saya, juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar rujukan.

Banjarmasin, 28 Juli 2023



Mita Yani Nurma Pratiwi

LEMBAR PENGESAHAN
SKRIPSI PROGRAM STUDI S-1 TEKNOLOGI INFORMASI
Analisis Pengaruh Data Normalisasi dan Non Normalisasi pada
Deteksi Cyberbullying Berbahasa Indonesia dengan Metode
IndoBERT
oleh
Mita Yani Nurma Pratiwi (1910817320003)

Telah dipertahankan di depan Tim Pengaji pada 3 Juli 2023 dan dinyatakan

L U L U S

Komite Pengaji :

Ketua : Andreyan Rizky Baskara, S.Kom., M.Kom
NIP 199307032019031011

Anggota 1 : Eka Setya Wijaya, S.T., M.Kom
NIP 198205082008011010

Anggota 2 : Nurul Fathawah Mustamin, S.Pd., M.T
NIP 199110252019032018

Pembimbing : Dr. Ir. Yuslена Sari, S.Kom., M.Kom
Utama NIP 198411202015042002

Banjarbaru, 01 SEP 2023
diketahui dan disahkan oleh:



Koordinator Program Studi
S-1 Teknologi Informasi,

Dr. Ir. Yuslена Sari, S.Kom., M.Kom

NIP 198411202015042002

PERSETUJUAN TUGAS AKHIR

ANALISIS PENGARUH DATA NORMALISASI DAN NON NORMALISASI
PADA DETEKSI CYBERBULLYING BERBAHASA INDONESIA DENGAN
METODE INDOBERT

OLEH
MITA YANI NURMA PRATIWI
NIM. 19108171320003

Telah diperiksa dan terpenuhi semua persyaratan akademik, administrasi, dan
disetujui untuk dipertahankan di hadapan dewan pengaji

Banjarmasin, 27 Juni 2023

Pembimbing Utama,



Dr. Ir. Yuslena Sari S.Kom., M.Kom.
NIP. 198411202015042002

ABSTRAK

Kasus *cyberbullying* yang marak terjadi di media sosial merupakan salah satu dampak negatif dari kebebasan dan kemudahan yang ditawarkan oleh media sosial. Berdasarkan 2.777 hasil *polling* pada remaja Indonesia berusia 14-24 tahun, 45% dari mereka melaporkan pernah mengalami *cyberbullying*. Hal ini menunjukkan bahwa deteksi *cyberbullying* sangat diperlukan sebagai salah satu bentuk pencegahan guna mengurangi kasus *cyberbullying*. NLP (*Natural Language Processing*) merupakan salah satu cara yang dapat digunakan untuk membantu dalam mendeteksi *cyberbullying*. IndoBERT dan IndoBERTweet merupakan *pre-trained* model yang dapat digunakan untuk deteksi *cyberbullying*. Data yang digunakan untuk melatih model NLP umumnya melewati tahap normalisasi, yaitu tahap mengubah kalimat non baku ke dalam kalimat baku sesuai dengan kaidah bahasa. Pengguna media sosial di Indonesia sendiri jarang menggunakan bahasa baku dalam melakukan komunikasi. Mereka cenderung menggunakan bahasa non baku dengan berbagai istilah dan kosakata baru. Oleh karena itu, penelitian ini dilakukan untuk mengetahui pengaruh data yang melewati proses normalisasi dan tidak melewati proses normalisasi pada model IndoBERT dan IndoBERTweet dalam mendeteksi *cyberbullying*. Data yang digunakan berasal dari Twitter dengan jumlah data 13.446. Model IndoBERT dan IndoBERTweet dengan data yang tidak melewati proses normalisasi memiliki hasil performansi lebih baik dibandingkan dengan data yang melewati proses normalisasi. Model IndoBERT tersebut memiliki nilai *accuracy*, *f1-score*, *recall*, dan *precision* sebesar 0.8520, 0.8520, 0.8520, dan 0.8521. Sementara itu, model IndoBERTweet memiliki nilai *accuracy*, *f1-score*, *recall*, dan *precision* sebesar 0.8602, 0.8602, 0.8602, dan 0.8609. Secara keseluruhan, model dengan hasil performansi terbaik adalah model IndoBERTweet dengan data yang tidak dinormalisasi.

Kata kunci: *Cyberbullying*, IndoBERT, IndoBERTweet, NLP, *Text Classification*

ABSTRACT

Cyberbullying cases that have been prevalent on social media are one of the negative consequences of the freedom and convenience offered by these platforms. Based on a survey of 2,777 Indonesian teenagers aged 14-24 years, 45% of them reported experiencing cyberbullying. This underscores the importance of cyberbullying detection as a preventive measure to reduce such incidents. Natural Language Processing (NLP) is one of the methods that can be used to aid in cyberbullying detection. IndoBERT and IndoBERTweet are pre-trained models that can be employed for this purpose. Data used to train NLP models typically undergoes normalization, which involves converting non-standard sentences into standardized sentences according to language norms. However, social media users in Indonesia rarely use standardized language in their communications. They tend to utilize non-standard language with various new terms and vocabulary. Therefore, this study aimed to investigate the influence of data that undergoes normalization and data that does not undergo normalization on the performance of IndoBERT and IndoBERTweet models in detecting cyberbullying. The data used in this research was obtained from Twitter, consisting of 13,446 samples. The results show that IndoBERT and IndoBERTweet models trained on data that did not undergo normalization achieved better performance compared to those trained on normalized data. The IndoBERT model exhibited accuracy, f1-score, recall, and precision values of 0.8520, 0.8520, 0.8520, and 0.8521, respectively. On the other hand, the IndoBERTweet model demonstrated accuracy, f1-score, recall, and precision values of 0.8602, 0.8602, 0.8602, and 0.8609, respectively. Overall, the model with the best performance was the IndoBERTweet model with non-normalized data.

Keywords: Cyberbullying, NLP, IndoBERT, IndoBERTweet, Text Classification

LEMBAR PERSEMBAHAN

Penulis mempersembahkan Tugas Akhir ini kepada:

1. Ayah, Ibu, Alm. Nenek, Adik serta keluarga tercinta yang telah memberikan doa, motivasi, dan dukungan saat menghadapi suatu masalah dalam keberlangsungan penyelesaian Tugas Akhir ini.
2. Ibu Dr. Ir. Yuslena Sari, S.Kom., M.Kom. selaku Koordinator Program Studi Teknologi Informasi dan Dosen Pembimbing Utama yang selalu menyempatkan waktunya untuk memberikan bimbingan dan arahan kepada mahasiswa teknologi informasi dan penulis untuk segera menyelesaikan Tugas Akhir ini.
3. Bapak Andry Fajar Zulkarnain, S.ST., M.T selaku Dosen Pembimbing Akademik yang selalu menyempatkan waktu untuk memberikan bimbingan, arahan, dan dukungan kepada penulis selama mengikuti proses perkuliahan.
4. Bapak Andreyan Rizky Baskara, S.Kom., M.Kom, Bapak Eka Setya Wijaya, S.T., M.Kom., dan Ibu Nurul Fathanah Mustamin, S.Pd., M.T. selaku Dosen Pengaji yang telah memberikan banyak arahan, masukan, dan solusi untuk penelitian ini.
5. Seluruh Dosen beserta Staf Program Studi Teknologi Informasi yang turut membantu dan mengarahkan dalam penyelesaian Tugas Akhir ini.
6. Seluruh teman Program Studi Teknologi Informasi Angkatan 2019 (SAKTI) beserta adik dan kakak tingkat yang telah banyak membantu serta memberikan dukungan dari perkuliahan hingga penyelesaian Tugas Akhir ini.
7. Teman-teman yang telah membantu dalam melakukan pemrosesan data untuk penelitian ini.
8. Eugynia Jessica Virgynia Rahman, Nurul Hikmah, Vania Laili Rahmah, Rahmadani, Nur Azlina, Nina Hafizah, Resya Cahaya Maharani, dan Rif'at Maulana yang selalu memberikan semangat dan dukungan, menghibur, mendengarkan keluh kesah, dan membantu selama proses studi berjalan hingga proses penyelesaian studi dilakukan.

9. Teman seperjuangan saya sejak SMA hingga detik ini, Mariati, Nurul Firdaus, Talitha Az Zahra Rosadi, Siti Rosita Sari, Irma Maulida, Deftri Sekar Ningrum, dan Nur Jannah yang selalu memberikan semangat dan dukungan, menghibur, mendengarkan keluh kesah, dan membantu selama SMA hingga proses penyelesaian studi dilakukan.
10. Diri sendiri yang tidak menyerah untuk terus belajar hingga detik ini.

KATA PENGANTAR

Puji dan syukur penulis ucapkan kehadiran Allah Subhanahu wa ta'ala yang telah melimpahkan rahmat dan hidayah-Nya. Sehingga penulis dapat menyelesaikan Tugas Akhir dengan Judul "*Analisis Pengaruh Data Normalisasi dan Non Normalisasi pada Deteksi Cyberbullying Berbahasa Indonesia dengan Metode IndoBERT*". Penulis menyadari penyelesaian laporan ini tidak akan terwujud tanpa adanya bimbingan dan bantuan dari berbagai pihak. Oleh karena itu, dalam kesempatan ini penulis menyampaikan terima kasih yang sebesar-besarnya kepada:

1. Rektor Universitas Lambung Mangkurat, Bapak Prof. Dr. Ahmad, S.E., M.Si., yang memimpin dan memanajemen jalannya seluruh perkuliahan yang ada di Universitas Lambung Mangkurat.
2. Dekan Fakultas Teknik, Bapak Prof. Dr. Ir. Irphan Fitrian Radam, S.T., M.T., IPU, yang memberikan layanan terbaik dalam perkuliahan, terkhusus pada pelaksanaan Tugas Akhir di lingkungan Fakultas Teknik.
3. Koordinator Program Studi Teknologi Informasi serta Pembimbing Utama, Ibu Dr. Ir. Yuslena Sari, S.Kom. M.Kom. yang telah memberikan arahan, bimbingan, serta solusi dalam penyelesaian Tugas Akhir.
4. Dosen-dosen beserta staff di Program Studi Teknologi Informasi yang telah mengarahkan dan teman-teman yang membantu dalam proses penyelesaian Tugas Akhir.

Akhir kata, penulis menyampaikan terimakasih kepada semua pihak yang turut serta membantu dalam penyelesaian laporan Tugas Akhir ini. Penulis mengharap saran dan kritik yang membangun demi perbaikan dan penyempurnaan laporan ini. Semoga laporan ini dapat bermanfaat baik untuk pribadi, pembaca dan semua pihak yang membutuhkan.

Banjarmasin, Agustus 2023

Penulis

Mita Yani Nurma Pratiwi

DAFTAR ISI

LEMBAR PERNYATAAN	i
ABSTRAK	iv
ABSTRACT	v
LEMBAR PERSEMAHAN	vi
KATA PENGANTAR	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR	xv
DAFTAR LAMPIRAN.....	xvii
DAFTAR RIWAYAT HIDUP.....	xviii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	5
1.3 Batasan Masalah.....	5
1.4 Tujuan Penelitian.....	6
1.5 Manfaat Penelitian.....	6
BAB II TINJAUAN PUSTAKA	7
2.1 Landasan Teori	7
2.1.1 <i>Cyberbullying</i>	7
2.1.2 <i>Natural Language Processing (NLP)</i>	9
2.1.3 <i>Text Classification</i>	9
2.1.4 <i>Text Pre-processing</i>	10
2.1.5 <i>Transformers</i>	10
2.1.6 <i>Bidirectional Encoder Representation from Transformers (BERT)</i>	12

2.1.7	<i>Indonesian based on the BERT (IndoBERT)</i>	14
2.1.8	<i>Indonesian BERT Tweet (IndoBERTTweet)</i>	14
2.1.9	<i>Transfer Learning</i>	14
2.1.10	<i>Pre-trained Model (PTM)</i>	15
2.1.11	<i>Confusion Matrix</i>	15
2.2	Penelitian Terkait	17
2.3.1	Deteksi <i>Cyberbullying</i> pada Facebook Menggunakan Algoritma K- <i>Nearest Neighbor</i>	17
2.3.2	Deteksi <i>Cyberbullying</i> Pada Data <i>Tweet</i> Menggunakan Metode <i>Random Forest</i> dan Seleksi Fitur <i>Information Gain</i>	17
2.3.3	<i>CyberBERT: BERT for cyberbullying identification</i>	17
2.3.4	<i>Rapid Cyber-bullying detection method using Compact BERT Models</i> 18	
2.3.5	<i>Cyberbullying Detection using Pre-Trained BERT Model</i>	18
2.3.6	<i>Bidirectional Encoder Representations from Transformers for Cyberbullying Text Detection in Indonesian Social Media</i>	19
2.3.7	<i>IndoBERTTweet: A Pretrained Language Model for Indonesian Twitter with Effective Domain-Specific Vocabulary Initialization</i>	20
2.3	Kerangka Berpikir	25
2.3.1	<i>Indicators</i>	25
2.3.2	<i>Proposed Method</i>	25
2.3.3	<i>Objectives</i>	26
2.3.4	<i>Measurement</i>	26
BAB III	METODOLOGI PENELITIAN	27
3.1	Alat dan Bahan	27
3.1.1	Alat Penelitian.....	27
3.1.2	Bahan Penelitian	28

3.2 Alur Penelitian.....	28
3.2.1 Identifikasi Masalah.....	28
3.2.2 Studi Literatur	29
3.2.3 Pengumpulan Data	29
3.2.4 Implementasi Model	30
3.2.5 Pengujian.....	34
3.2.6 Analisis Hasil Uji.....	36
3.2.7 Implementasi Sistem.....	37
BAB IV HASIL DAN PEMBAHASAN	38
4.1 Pengumpulan Data	38
4.2 <i>Text Pre-processing</i>	38
4.3 Pelabelan Data.....	40
4.4 Pembagian Data.....	40
4.5 Eksperimen Model	41
4.5.1 Eksperimen Model IndoBERT dengan Data Normalisasi	41
4.5.2 Eksperimen Model IndoBERT dengan Data Non Normalisasi	44
4.5.3 Eksperimen Model IndoBERTweet dengan Data Normalisasi.....	47
4.5.4 Eksperimen Model IndoBERTweet dengan Data Non Normalisasi	51
4.6 Pengujian Model	54
4.7 Analisis Hasil Uji Model.....	57
4.8 Implementasi Model.....	59
BAB V KESIMPULAN DAN SARAN	60
5.1 Kesimpulan.....	60
5.2 Saran.....	61
DAFTAR PUSTAKA	62
LAMPIRAN.....	64

DAFTAR TABEL

Tabel 2. 1 Confusion Matrix	16
Tabel 2. 2 Ringkasan Penelitian Terkait	21
Tabel 3. 1 Alat Penelitian.....	27
Tabel 3. 2 Jumlah Data Yang Digunakan (Kaggle)	29
Tabel 3. 3 Jumlah Data untuk Pencarian dengan Kata Kunci.....	29
Tabel 3. 4 Contoh Data	30
Tabel 3. 5 Proses Case Folding.....	32
Tabel 3. 6 Proses Filtering	32
Tabel 3. 7 Proses Tokenizing.....	32
Tabel 3. 8 Proses Normalization.....	33
Tabel 3. 9 Contoh Confusion Matrix	36
Tabel 4. 1 Hasil Akhir Proses Text Pre-processing	39
Tabel 4. 2 Kriteria Pelabelan Data	40
Tabel 4. 3 Hasil Pelabelan Data	40
Tabel 4. 4 Hasil Pembagian Data.....	41
Tabel 4. 5 Konfigurasi Parameter Eksperimen 1: Model IndoBERT dengan Data Normalisasi	41
Tabel 4. 6 Hasil Eksperimen Pertama: Model IndoBERT dengan Data Normalisasi	42
Tabel 4. 7 Konfigurasi Parameter Eksperimen 2: Model IndoBERT dengan Data Non Normalisasi	44
Tabel 4. 8 Hasil Eksperimen Kedua: Model IndoBERT dengan Data Non Normalisasi	45
Tabel 4. 9 Konfigurasi Parameter Eksperimen 3: Model IndoBERTweet dengan Data Normalisasi.....	48
Tabel 4. 10 Hasil Eksperimen Ketiga: Model IndoBERTweet dengan Data Normalisasi	48
Tabel 4. 11 Konfigurasi Parameter Eksperimen 4: Model IndoBERTweet dengan Data Non Normalisasi.....	51

Tabel 4. 12 Hasil Eksperimen Keempat: Model IndoBERTweet dengan Data Non Normalisasi	52
Tabel 4. 13 Hasil Pengujian 1: Model IndoBERT dengan Data Normalisasi.....	55
Tabel 4. 14 Hasil Pengujian 2: Model IndoBERT Data Non Normalisasi	55
Tabel 4. 15 Hasil Pengujian 3: Model IndoBERTweet dengan Data Normalisasi	56
Tabel 4. 16 Hasil Pengujian 4: Model IndoBERTweet dengan Data Non Normalisasi	57
Tabel 4. 17 Perbandingan Hasil Performansi Model IndoBERT dan IndoBERTweet	57
Tabel 4. 18 Hasil Klasifikasi Data Model IndoBERT Terbaik.....	58
Tabel 4. 19 Contoh Hasil Klasifikasi Data.....	58
Tabel 4. 20 Contoh Kalimat Yang Salah Diperbaiki	59

DAFTAR GAMBAR

Gambar 2. 1 Arsitektur Transformers [21]	11
Gambar 2. 2 Arsitektur BERT [23].....	13
Gambar 2. 3 BERT Input Representation [23]	14
Gambar 2. 4 Contoh Intuitif Transfer Learning [25]	15
Gambar 2. 5 Kerangka Pemikiran.....	25
Gambar 3. 1 Alur Penelitian	28
Gambar 3. 2 Proses Implementasi Model IndoBERT dan IndoBERTweet.....	31
Gambar 3. 3 Skema Pengujian.....	36
Gambar 4. 1 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 2e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	42
Gambar 4. 2 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 2e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	43
Gambar 4. 3 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 3e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	43
Gambar 4. 4 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 3e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	43
Gambar 4. 5 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 5e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	44
Gambar 4. 6 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 5e-5 Epoch 2 Model IndoBERT Data Normalisasi.....	44
Gambar 4. 7 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 2e-5 Epoch 2 Model IndoBERT Data Non Normalisasi.....	46
Gambar 4. 8 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 2e-5 Epoch Model IndoBERT Data Non Normalisasi.....	46
Gambar 4. 9 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 3e-5 Epoch 2 Model IndoBERT Data Non Normalisasi.....	46
Gambar 4. 10 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 3e-5 Epoch 2 Model IndoBERT Data Non Normalisasi.....	47
Gambar 4. 11 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 5e-5 Epoch 2 Model IndoBERT Data Non Normalisasi.....	47

Gambar 4. 12 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 5e-5 Epoch 2 Model IndoBERT Data Non Normalisasi.....	47
Gambar 4. 13 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 2e-5 Epoch 2 Model IndoBERTweet Data Normalisasi	49
Gambar 4. 14 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 2e-5 Epoch 2 Model IndoBERTweet Data Normalisasi	49
Gambar 4. 15 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 3e-5 Epoch 2 Model IndoBERTweet Data Normalisasi	50
Gambar 4. 16 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 3e-5 Epoch 2 Model IndoBERTweet Data Normalisasi	50
Gambar 4. 17 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 5e-5 Epoch 2 Model IndoBERTweet Data Normalisasi	50
Gambar 4. 18 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 5e-5 Epoch Model IndoBERTweet Data Normalisasi	51
Gambar 4. 19 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 2e-5 Epoch 2 Model IndoBERTweet Data Non Normalisasi	52
Gambar 4. 20 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 2e-5 Epoch Model IndoBERTweet Data Non Normalisasi	53
Gambar 4. 21 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 3e-5 Epoch 2 Model IndoBERTweet Data Non Normalisasi	53
Gambar 4. 22 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 3e-5 Epoch 2 Model IndoBERTweet Data Non Normalisasi	53
Gambar 4. 23 Nilai Accuracy dan Loss Batch Size 16 dan Learning Rate 5e-5 Epoch 2 Model IndoBERTweet Data Non Normalisasi	54
Gambar 4. 24 Nilai Accuracy dan Loss Batch Size 32 dan Learning Rate 5e-5 Epoch 2 Model IndoBERTweet Data Non Normalisasi	54
Gambar 4. 25 Tampilan Antarmuka Sistem Deteksi Cyberbullying	59

DAFTAR LAMPIRAN

Lampiran 1. Source Code untuk Text Pre-Processing	64
Lampiran 2. Source Code untuk Implementasi Model	68
Lampiran 3. Lembar Konsultasi	80