

**KLUSTERISASI TOPIK SKRIPSI MAHASISWA PROGRAM
STUDI INFORMATIKA MENGGUNAKAN METODE LDA
(LATENT DIRICHLET ALLOCATION) DAN BERT
(BIDIRECTIONAL ENCODER REPRESENTATIONS
FROM TRANSFORMERS)**

SKRIPSI

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana
Program Studi S1 Informatika



disusun oleh
YAHYA HANDARESTANTO
21.11.4365

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2025

**KLUSTERISASI TOPIK SKRIPSI MAHASISWA PROGRAM
STUDI INFORMATIKA MENGGUNAKAN METODE LDA
(LATENT DIRICHLET ALLOCATION) DAN BERT
(BIDIRECTIONAL ENCODER REPRESENTATIONS
FROM TRANSFORMERS)**

SKRIPSI

untuk memenuhi salah satu syarat mencapai derajat Sarjana

Program Studi *SI* Informatika



disusun oleh

YAHYA HANDARESTANTO

21.11.4365

Kepada

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2025**

HALAMAN PERSETUJUAN

SKRIPSI

**KLUSTERISASI TOPIK SKRIPSI MAHASISWA PROGRAM
STUDI INFORMATIKA MENGGUNAKAN METODE LDA
(LATENT DIRICHLET ALLOCATION) DAN BERT
(BIDIRECTIONAL ENCODER REPRESENTATIONS FROM
TRANSFORMERS)**

yang disusun dan diajukan oleh

Yahya Handarestanto

21.11.4365

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 20 Januari 2025

Dosen Pembimbing,



Hanafi, S.Kom., M.Eng., Ph.D

NIK. 190302024

HALAMAN PENGESAHAN

SKRIPSI

**KLUSTERISASI TOPIK SKRIPSI MAHASISWA PROGRAM
STUDI INFORMATIKA MENGGUNAKAN METODE LDA
(LATENT DIRICHLET ALLOCATION) DAN BERT
(BIDIRECTIONAL ENCODER REPRESENTATIONS FROM
TRANSFORMERS)**

yang disusun dan diajukan oleh

Yahya Handarestanto

21.11.4365

Telah dipertahankan di depan Dewan Pengaji
pada tanggal 20 Januari 2025

Nama Pengaji

Dr. Ferry Wahyu Wibowo, S.Si., M.Cs.
NIK. 190302235

Dina Maulina, S.Kom., M.Kom.
NIK. 190302250

Hanafi, S.Kom., M.Eng., Ph.D
NIK. 190302024

Tanda Tangan



Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal 20 Januari 2025

DEKAN FAKULTAS ILMU KOMPUTER



Hanif Al Fatta, S.Kom., M.Kom., Ph.D.
NIK. 190302096

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

**Nama mahasiswa : Yahya Handarestanto
NIM : 21.11.4365**

Menyatakan bahwa Skripsi dengan judul berikut:

Klusterisasi topik skripsi mahasiswa program studi Informatika menggunakan metode LDA (Latent Dirichlet Allocation) Dan BERT (Bidirectional Encoder Representations from Transformers)

Dosen Pembimbing: Hanafi, S. Kom., M.Eng., Ph.D.

1. Karya tulis ini adalah benar-benar ASLI dan BELUM PERNAH diajukan untuk mendapatkan gelar akademik, baik di Universitas AMIKOM Yogyakarta maupun di Perguruan Tinggi lainnya.
2. Karya tulis ini merupakan gagasan, rumusan dan penelitian SAYA sendiri, tanpa bantuan pihak lain kecuali arahan dari Dosen Pembimbing.
3. Dalam karya tulis ini tidak terdapat karya atau pendapat orang lain, kecuali secara tertulis dengan jelas dicantumkan sebagai acuan dalam naskah dengan disebutkan nama pengarang dan disebutkan dalam Daftar Pustaka pada karya tulis ini.
4. Perangkat lunak yang digunakan dalam penelitian ini sepenuhnya menjadi tanggung jawab SAYA, bukan tanggung jawab Universitas AMIKOM Yogyakarta.
5. Pernyataan ini SAYA buat dengan sesungguhnya, apabila di kemudian hari terdapat penyimpangan dan ketidakbenaran dalam pernyataan ini, maka SAYA bersedia menerima SANKSI AKADEMIK dengan pencabutan gelar yang sudah diperoleh, serta sanksi lainnya sesuai dengan norma yang berlaku di Perguruan Tinggi.

Yogyakarta, 20-01-2025

Yang Menyatakan,



Yahya Handarestanto

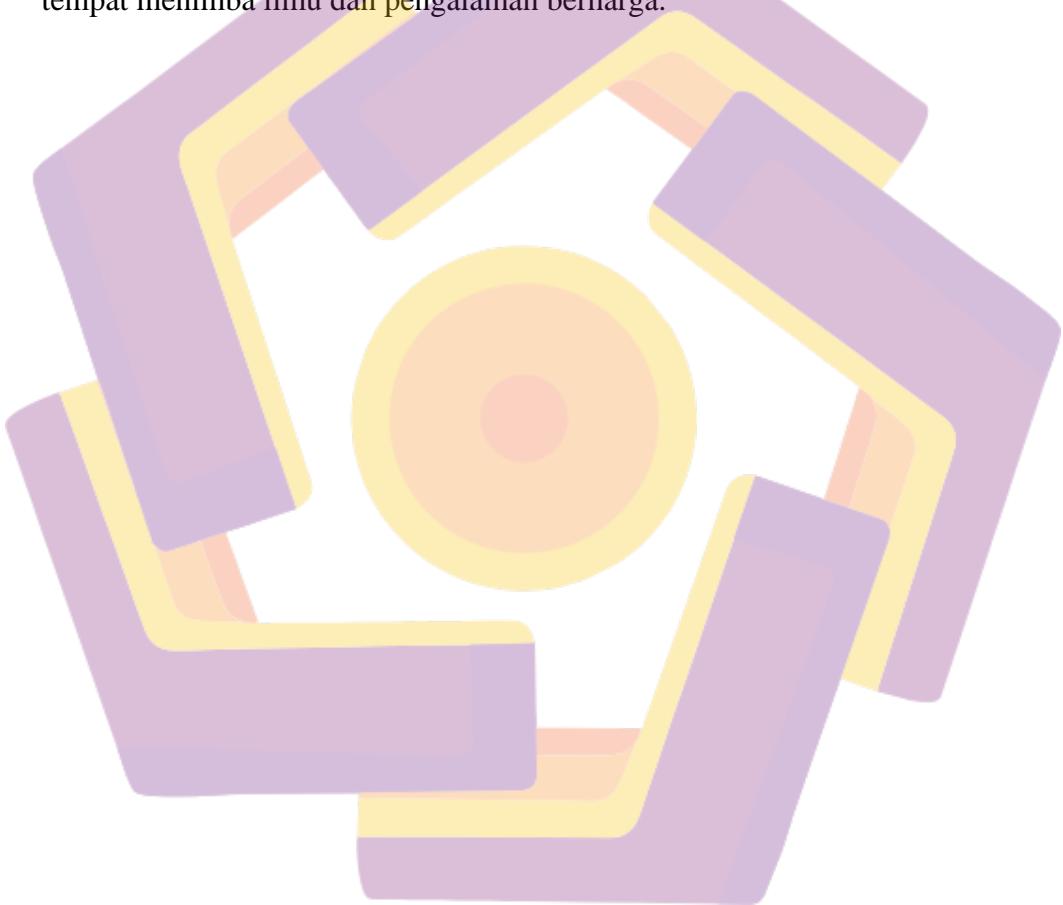
HALAMAN PERSEMBAHAN

Dengan penuh rasa syukur kepada Allah SWT, karya ilmiah ini saya persembahkan kepada:

Kedua orang tua saya tercinta, Bapak Radiyo dan Ibu Kustrini yang telah memberikan dukungan, doa, dan pengorbanan yang tak terhingga.

Para dosen pembimbing yang telah dengan sabar membimbing dan memberikan ilmu yang bermanfaat selama proses penulisan skripsi.

Almamater tercinta, Universitas Amikom Yogyakarta yang telah menjadi tempat menimba ilmu dan pengalaman berharga.



KATA PENGANTAR

Puji syukur penulis panjatkan kepada Allah Yang Maha Esa atas limpahan rahmat, hidayah, dan karunia-Nya sehingga penulis dapat menyelesaikan skripsi yang berjudul "Klusterisasi topik skripsi mahasiswa program studi Informatika menggunakan metode LDA dan BERT" dengan lancar dan tepat waktu. Skripsi ini disusun guna memenuhi salah satu syarat untuk memperoleh gelar S.Kom pada Program Studi S1 Informatika di Universitas Amikom Yogyakarta.

Penulis mengucapkan terima kasih yang sebesar-besarnya kepada:

1. Bapak Hanafi, S.Kom., M.Eng., Ph.D., selaku dosen pembimbing yang telah memberikan bimbingan, arahan, dan motivasi selama proses penyusunan skripsi ini.
2. Tim Dosen Pengaji yang telah memberikan masukan dan saran yang sangat berharga untuk perbaikan skripsi ini.
3. Kedua orang tua yang senantiasa memberikan dukungan moril dan materiil serta doa yang tiada henti.
4. Universitas Amikom Yogyakarta yang telah memberikan kesempatan dan dukungan dalam penelitian ini.

Penulis menyadari bahwa skripsi ini masih jauh dari sempurna. Oleh karena itu, penulis mengharapkan kritik dan saran yang membangun dari berbagai pihak untuk perbaikan di masa mendatang.

Akhir kata, semoga skripsi ini dapat memberikan manfaat bagi pembaca dan pihak-pihak yang berkepentingan.

Yogyakarta, 20 Januari 2025

Penulis

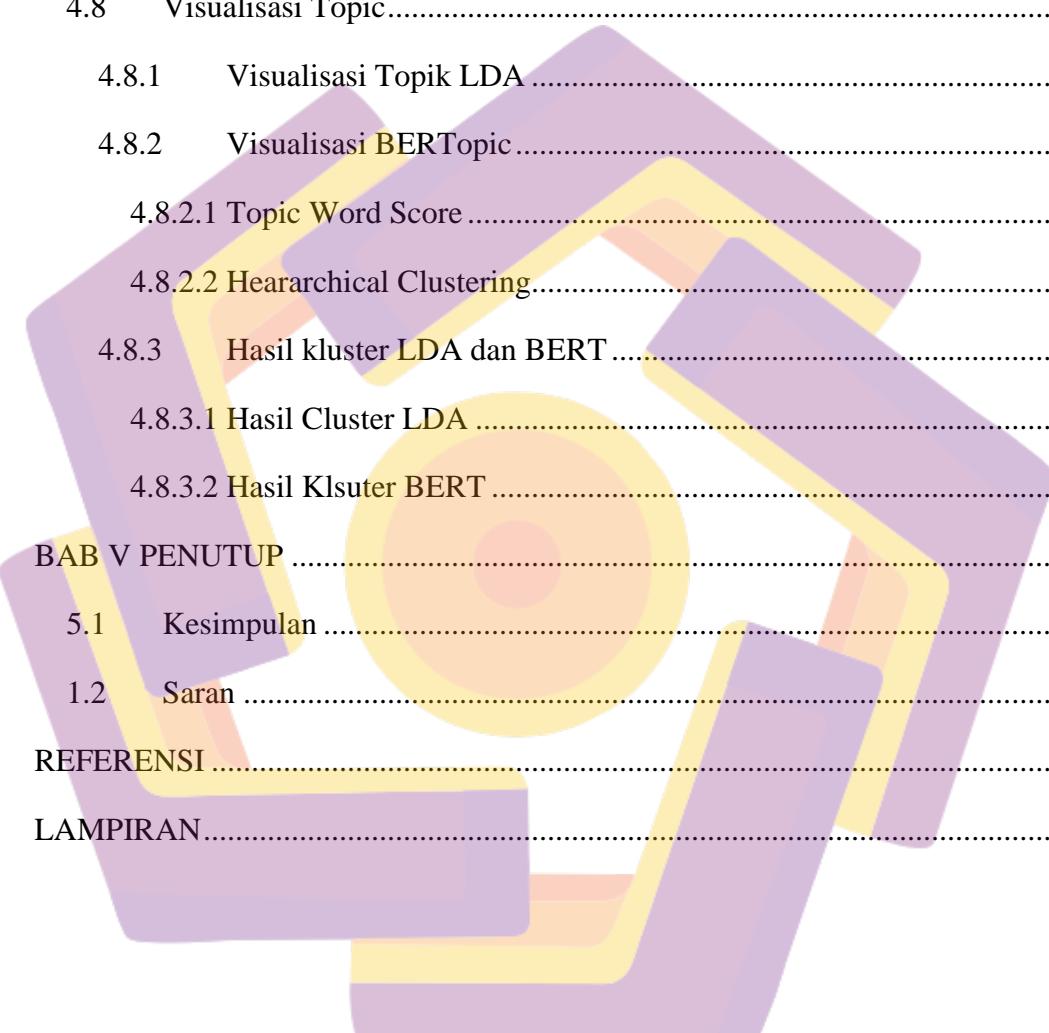
Yahya Handarestanto

DAFTAR ISI

HALAMAN JUDUL	ii
HALAMAN PERSETUJUAN.....	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERNYATAAN KEASLIAN SKRIPSI	iv
HALAMAN PERSEMBERAHAN	v
KATA PENGANTAR	vi
DAFTAR ISI.....	vii
DAFTAR TABEL.....	xi
DAFTAR GAMBAR	xii
DAFTAR LAMPIRAN.....	xiii
DAFTAR LAMBANG DAN SINGKATAN	xiv
DAFTAR ISTILAH.....	xv
INTISARI	xvii
<i>ABSTRACT</i>	xviii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	2
1.3 Batasan Masalah	2
1.4 Tujuan Penelitian	2
1.5 Manfaat Penelitian	3
1.5.1 Manfaat Segi Teoritis Dan Praktis	3
1.5.2 Manfaat Bagi Penelitian Selanjutnya.....	3
1.6 Sistematika Penulisan	3
BAB II TINJAUAN PUSTAKA	5

2.1	Studi Literatur	5
2.2	Dasar Teori.....	10
2.2.1	Google Collab	10
2.2.2	Natural Languange Procesing (NLP).....	11
2.2.3	Text Mining	12
2.2.4	Preprocesing.....	12
2.2.5	Python	13
2.2.6	Bidirectional Encoder Representations from Transformers (BERT)....	14
2.2.7	Packages.....	15
2.2.7.1	Gensim	16
2.2.7.2	Sciket-learn	17
2.2.7.3	NLTK	17
2.2.7.4	PyLDAvis	18
2.2.7.5	Pandas	19
2.2.7.6	BERTopic	19
2.2.7.7	Wordcloud.....	20
2.2.7.8	Seaborn	21
2.2.7.9	Matplotlib.....	21
2.2.8	Term Frequency – Inverse Document Frequency	22
2.2.9	Topik Modeling	24
2.2.10	Latent Dirichlet Allocation (LDA)	25
2.2.11	Topic Coherence	27
2.2.12	Intertopic Distance Map.....	28
2.2.13	Topic Word Score	29
	BAB III METODE PENELITIAN	30

3.1	Objek Penelitian.....	30
3.1.1	Penjelasan Datasheet.....	30
3.2	Skema Penelitian.....	31
3.2.1	Studi Leterature.....	32
3.2.2	Pengumpulan Data	32
3.2.3	Seleksi Data	34
3.2.4	Visualisasi Data	34
3.2.5	Text Preprocesing	34
3.2.6	Fitur Extraction TF-IDF.....	36
3.2.7	Topic Modeling LDA	36
3.2.8	Topic Modeling BERT	39
3.2.9	Evaluasi Model	41
3.3	Alat Dan Bahan.....	44
3.3.1	Alat.....	44
3.3.2	Bahan	44
BAB IV HASIL DAN PEMBAHASAN		46
4.1	Pengumpulan Data	46
4.2	Seleksi Data	47
4.3	Visualisasi Data	48
4.4	Text Preprocesing	49
4.4.1	Cleaning Data.....	49
4.4.2	Lowercase	50
4.4.3	Normalisasi	51
4.4.4	Tokenisasi	52
4.4.5	Stopword Removal.....	53
4.4.6	Lemmatization	55



4.4	Fitur Extraction TF-IDF.....	56
4.5	Topic Modeling LDA	57
4.6	Topic Modeling BERT	59
4.7	Coherence Score	60
4.8	Visualisasi Topic.....	61
4.8.1	Visualisasi Topik LDA	72
4.8.2	Visualisasi BERTopic	72
4.8.2.1	Topic Word Score	72
4.8.2.2	Heararchical Clustering.....	80
4.8.3	Hasil kluster LDA dan BERT	80
4.8.3.1	Hasil Cluster LDA	80
4.8.3.2	Hasil Klsuter BERT	80
BAB V	PENUTUP	84
5.1	Kesimpulan	84
1.2	Saran	84
REFERENSI	86	
LAMPIRAN	91	

DAFTAR TABEL

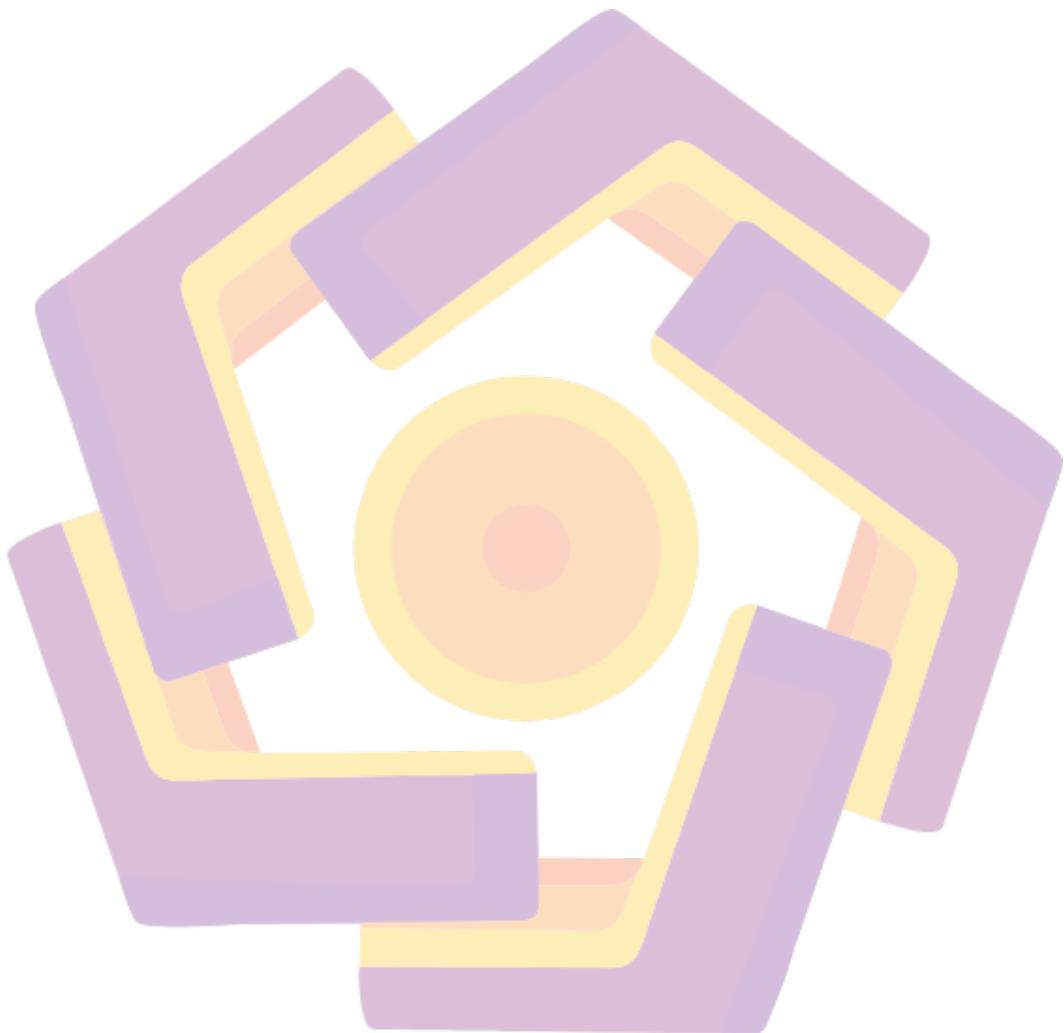
Table 2.1 Keaslian Penelitian	7
Table 3. 1 Penjelasan dari Dataaset	31
Table 4.1 Visualisasi Data Abstract.....	46
Table 4.2 Lanjutan dari table visualisasi data abstrak	47
Table 4.3 Hasil Seleksi Data.....	47
Table 4.4 Contoh Proses Cleaning Data	49
Table 4.5 lanjutan dari Contoh Proses Cleaning Data	50
Table 4.6 Contoh Proses Lowercase	50
Table 4. 7 lanjutan dari contoh proses Lowercase	51
Table 4.8 Contoh Proses Normalisasi Text.....	52
Table 4.9 Contoh proses Tokenisasi	53
Table 4.10 Contoh Proses Stopword Removal	54
Table 4.11 lanjutan dari Contoh Proses Stopword Removal	55
Table 4.12 Contoh Proses lemmatization	55
Table 4.13 Lanjutan dari contoh proses lemmatization	56
Table 4.14 Feature extraction TF-IDF	57
Table 4.15 Topic Modeling LDA	58
Table 4.16 lanjutan dari Topic Modeling LDA	59
Table 4.17 Topic Modeling BERT	60

DAFTAR GAMBAR

Gambar 2.1 Contoh Wordcloud	20
Gambar 2.2 Representasi Graphical Model LDA	26
Gambar 3.1 Gambar Skema Penelitian	32
Gambar 3.2 Proses Pengumpulan Data.....	33
Gambar 3.3 Tahapan Text Preprocesing.....	36
Gambar 3.4 Tahapan Modeling LDA	39
Gambar 3.5 Tahapan Modeling Bert.....	40
Gambar 3.6 Tahapan Evaluasi Model.....	41
Gambar 4.1 Contoh Visualisasi WordCloud	48
Gambar 4.2 Intertopic Distance Map LDA topik kesatu	62
Gambar 4.3 Hasil clusteriasi LDA topik kesatu	63
Gambar 4.4 Intertopic Distance LDA Map topik ke-dua.....	64
Gambar 4.5 Hasil Clusterisasi LDA Topik Kedua	65
Gambar 4.6 Intertopic Distance Map LDA topik ketiga.....	66
Gambar 4.7 Hasil Clusterisasi LDA Topik Ketiga	67
Gambar 4.8 Intertopic Distance Map LDA Topic keempat.....	68
Gambar 4.9 Hasil clusterisasi LDA topic Ke-empat.....	69
Gambar 4.10 Intertopic Distance Map LDA Topik ke lima	70
Gambar 4.11 Hasil Clusterisasi Topik LDA ke-lima.....	71
Gambar 4.12 Hasil Clusterisasi BERT topic 0-1	72
Gambar 4.13 Hasil Clusterisasi BERT topic 2-3	73
Gambar 4.14 Hasil Clusterisasi BERT topic 4 - 5	74
Gambar 4.15 Hasil Clusterisasi BERT topic 6-7	74
Gambar 4. 16 Hasil Clusterisasi BERT topic 8 - 9	75
Gambar 4.17 Hasil Clusterisasi BERT topic 10-11	75
Gambar 4.18 Hasil Clusterisasi BERT topic 12-13	76
Gambar 4. 19 Hasil Clusterisasi BERT topic 14-15	76
Gambar 4. 20 Hasil Heararchical Clustering BERT	78

DAFTAR LAMPIRAN

lampiran 1 link Source Code Github 1 91

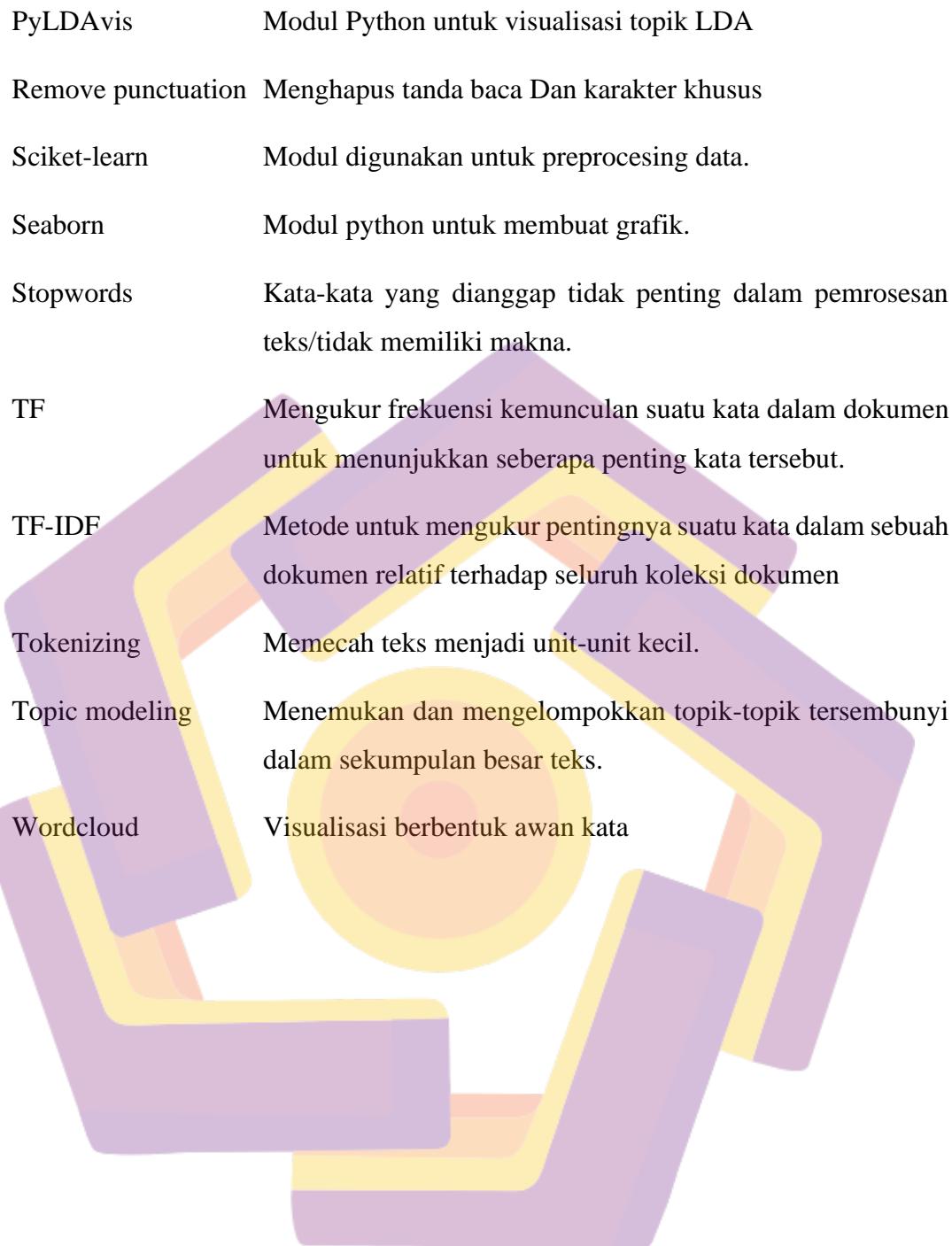


DAFTAR LAMBANG DAN SINGKATAN

AI	Artificial Intelligence
BoW	Bag-Of-Words
CPU	Central Processing Unit
GPU	Graphics Processing Unit
IDF	Inverse Document Frequency
LDA	Latent Dirichlet Allocation
BERT	Bidirectional Encoder Representations from Transformers
NLP	Natural Language Processing
NLTK	Natural Language Toolkit
OOP	Objek-oriented programming
PCA	Principal Component Analysis
TF	Term Frekuensi
TF-IDF	Term Frequency-Inverse Document Frequency
TPU	Tensor Processing Unit
t-SNE	t-distributed Stochastic Neighbor Embedding
URL	Uniform Resource Locator
VAE	Variational Autoencoders
N	Jumlah Kemunculan Kata
θ	Distribusi topik
ω	Probabilitas Kata pada Topik

DAFTAR ISTILAH

Alpha	distribusi topik per dokumen
Beta	distribusi kata per topik
Case folding	Mengubah semua huruf dalam suatu teks menjadi huruf kecil semua
Corpus	kumpulan teks yang digunakan untuk analisis atau pelatihan model.
Gensim	Pustaka python untuk pemodelan topik
IDF	Mengukur seberapa jarang kata tersebut muncul di seluruh dokumen dalam korpus.
Lemmatization	Mengubah kata ke bentuk dasarnya
Machine learning	Cabang dari kecerdasan buatan (Artificial Intelligence/AI) yang memungkinkan komputer untuk belajar dari data dan membuat prediksi atau keputusan tanpa diprogram secara eksplisit.
Matplotlib	Modul untuk visualisasi data secara 2d dan 3d
Meaningless_words	Daftar kata tidak relevan
Missing value	Nilai yang hilang atau kosong
Num Topics	Jumlah topik yang ingin dihasilkan
Overfitting	Model terlalu hafal data pelatihan sehingga sulit menangkap pola umum untuk data lain
Packages	Sekumpulan modul Python
Pandas	Modul Python untuk analisis data.
Preprocessing	Srangkaian langkah untuk mempersiapkan data mentah agar siap digunakan dalam proses pelatihan model.



INTISARI

Penelitian ini bertujuan untuk melakukan pemodelan topik pada kumpulan skripsi mahasiswa Studi Informatika di Universitas Amikom Yogyakarta dengan menggunakan metode Latent Dirichlet Allocation (LDA) dan BERTopic. Analisis dilakukan terhadap abstrak skripsi dari tahun 2020 hingga 2024, yang diolah melalui beberapa tahap preprocessing, seperti penghapusan nilai null, normalisasi teks, tokenisasi, dan penghapusan stopwords. Sebagai langkah awal, analisis TF-IDF dilakukan untuk mengidentifikasi kata-kata kunci yang paling berpengaruh dalam dataset.

Metode LDA digunakan untuk mengidentifikasi topik utama dalam dokumen dengan memandang setiap dokumen sebagai campuran dari beberapa topik tersembunyi. Sementara itu, BERTopic diterapkan dengan pendekatan berbasis BERT untuk menangkap representasi semantik yang lebih dalam dari teks. Evaluasi hasil pemodelan pada penilitian ini dilakukan dengan mencari parameter terbaik yaitu dengan melakukan pencarian skor koherensi tertinggi dengan rentang jumlah topik dari 2 hingga 20. Untuk model LDA menghasilkan skor koherensi sebesar 0.5933 dengan 5 topik, nilai Beta 0.6 dan nilai Alpha nya auto. Dan untuk model BERT/BERTopic menggunakan paraphrase all-MiniLM-L6-v2 menghasilkan skor koherensi 0.7502 dengan ukuran topik 16. Temuan ini menunjukkan bahwa model BERT/BERTopic lebih efektif dalam mengidentifikasi dan mengelompokkan topik pada dataset skripsi mahasiswa Studi Informatika di Universitas Amikom Yogyakarta. Hasil penelitian ini diharapkan dapat memberikan wawasan tentang arah penelitian dan minat mahasiswa, serta menjadi referensi bagi topik penelitian selanjutnya.

Kata kunci: LDA, BERTopic, TF-IDF, Skripsi, Topic modeling

ABSTRACT

This research aims to conduct topic modeling on a collection of Informatics Study student theses at Amikom University Yogyakarta using the Latent Dirichlet Allocation (LDA) and BERTopic methods. Analysis was carried out on thesis abstracts from 2020 to 2024, which were processed through several preprocessing stages, such as removing null values, text normalization, tokenization, and removing stop words. As a first step, TF-IDF analysis was carried out to identify the most influential keywords in the dataset.

The LDA method is used to identify the main topics in documents by viewing each document as a mixture of several hidden topics. Meanwhile, BERTopic is implemented with a BERT-based approach to capture deeper semantic representation of text. Evaluation of the modeling results in this research was carried out by looking for the best parameters, namely by searching for the highest coherence score with a range of topics from 2 to 20. The LDA model produced a coherence score of 0.5933 with 5 topics, a Beta value of 0.6 and an Alpha value of auto. And for the BERT/BERTopic model, using the all-MiniLM-L6-v2 paraphrase produces a coherence score of 0.7502 with a topic size of 16. These findings indicate that the BERT/BERTopic model is more effective in identifying and grouping topics in the thesis dataset of Informatics Study students at Amikom University, Yogyakarta. It is hoped that the results of this research will provide insight into research directions and student interests, as well as become a reference for further research topics.

Keyword: *LDA, BERTopic, TF-IDF, Thesis, Topic modeling*