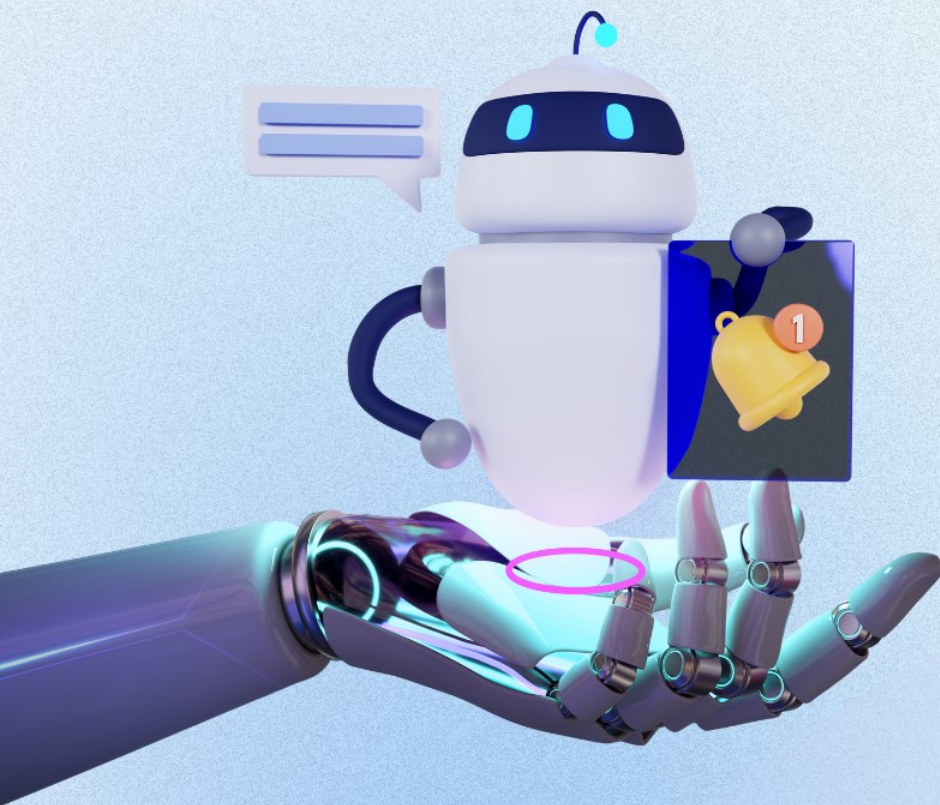


# Pengantar NLP dan Topik Model LDA



Oleh  
Albertus Dwiyoga Widianoro  
Prof. Mustafid  
Prof. Ridwan Sanjaya

# PENGANTAR NLP DAN TOPIK MODEL LDA

Oleh:  
Albertus Dwiyoga Widianoro  
Prof. Mustafid  
Prof. Ridwan Sanjaya

Diterbitkan oleh



2024

# Pengantar NLP dan Topik Model LDA

Oleh:

Albertus Dwiyoga Widianoro

Prof. Mustafid

Prof. Ridwan Sanjaya

Uk. 15,5cm x 23cm (vi + 122hlm)

ISBN: 978-623-10-4853-0

Diterbitkan oleh



Editor: Dr. Usman Ependi, S.Kom., M.Kom

Edisi November 2024

## **Hak Cipta dilindungi undang-undang**

*Dilarang memperbanyak sebagian atau seluruh isi buku ini dalam bentuk apapun, baik secara elektronik maupun mekanik, termasuk memfotokopi, merekam atau dengan menggunakan system penyimpanan, tanpa izin tertulis dari Penulis.*

# KATA PENGANTAR

Dengan bangga kami mempersembahkan buku *Pengantar NLP dan Topik Model LDA*, sebuah karya yang dirancang untuk menjembatani kesenjangan pengetahuan di bidang pemrosesan bahasa alami (Natural Language Processing/NLP) dan pemodelan topik. Dalam era digital yang semakin berkembang pesat, kemampuan untuk memahami dan mengolah data teks menjadi keterampilan yang sangat diperlukan. Buku ini hadir untuk memenuhi kebutuhan pembaca akan pemahaman yang mendalam dan praktis tentang bagaimana teknologi NLP dan model topik seperti Latent Dirichlet Allocation (LDA) bekerja.

Buku ini disusun dengan alur yang sistematis, dimulai dari konsep dasar NLP hingga ke teknik pemodelan topik yang lebih kompleks. Setiap bab telah disusun dengan seksama untuk memberikan wawasan yang tidak hanya mendasar tetapi juga aplikatif, dengan berbagai contoh dan teknik yang relevan dengan kebutuhan analisis data modern. Kami berharap buku ini dapat memberikan manfaat yang luas bagi pembaca, baik yang baru memulai di dunia NLP maupun yang ingin memperdalam pengetahuan tentang pemodelan topik.

Akhir kata, kami mengucapkan terima kasih kepada seluruh pihak yang telah berkontribusi dalam penyusunan buku ini. Semoga buku *Pengantar NLP dan Topik Model LDA* ini dapat menjadi sumber referensi yang berharga dan mendukung perkembangan ilmu pengetahuan di bidang analisis data dan kecerdasan buatan di Indonesia.

Penulis

# DAFTAR ISI

<b>SAMPUL DALAM.....</b>	<b>ii</b>
<b>KATA PENGANTAR.....</b>	<b>iv</b>
<b>DAFTAR ISI.....</b>	<b>v</b>
<b>BAB 1 PENDAHULUAN .....</b>	<b>1</b>
A. Definisi NLP.....	1
B. Sejarah singkat NLP .....	4
C. Relevansi dan penerapan NLP dalam kehidupan .....	5
D. Tujuan.....	7
<b>BAB II DASAR-DASAR NLP .....</b>	<b>9</b>
A. Pengertian dasar tentang bahasa alami .....	9
B. Teknik-teknik dasar dalam NLP .....	11
C. Algoritma dan Model dalam NLP .....	26
1. Naive Bayes .....	29
2. Support Vector Machines (SVM) .....	31
3. Recurrent Neural Networks (RNN) .....	34
4. Convolutional Neural Networks (CNN) .....	37
<b>BAB III PENERAPAN NLP .....</b>	<b>45</b>
A. NLP dalam Industri .....	45
1. Customer Service .....	45
2. Analisis Sentimen .....	46
3. Pencarian Informasi .....	48
4. Fine-tuning Model .....	49
B. NLP dalam Kesehatan .....	51
C. NLP dalam Pendidikan .....	53
1. Evaluasi dan pembelajaran berbasis teks .....	53

2. Analisis plagiarisme .....	55
<b>BAB IV TANTANGAN DAN ISU ETIKA DALAM NLP ...</b>	<b>57</b>
A. Tantangan dalam NLP .....	57
1. Polysemy dan Ambiguitas .....	57
2. Kurangnya data yang berkualitas.....	58
3. Overfitting dan generalisasi .....	59
B. Isu Etika dalam NLP .....	61
1. Privasi dan keamanan data.....	61
2. Bias dalam data dan model .....	62
<b>BAB V TOPIK MODEL.....</b>	<b>64</b>
A. Topik Model LDA .....	64
B. Ruang Lingkup LDA .....	65
C. Tujuan Implementasi LDA .....	66
D. Konsep Dasar Probabilistik .....	67
E. Contoh Konsep Dasar Probabilistik:.....	68
F. Aspek Pemodelan distribusi data memiliki.....	71
G. Contoh Pemodelan Distribusi Data .....	72
<b>BAB VI LATENT DIRICHLET ALLOCATION (LDA) ....</b>	<b>77</b>
A. Prinsip Kerja LDA.....	79
B. Pembentukan Model: .....	80
C. Penyesuaian Parameter .....	81
D. Persamaan dalam Model LDA .....	82
E. Perhitungan Distribusi Posterior .....	82
F. Rumus-rumus LDA .....	84
G. Proses Model LDA .....	85
H. Variabel Tersembunyi: .....	85
1. Distribusi Topik Dokumen $\theta$ .....	85
2. Distribusi Kata dalam Topik $\beta$ .....	86

I. Implementasi LDA.....	86
1. Menghapus tanda baca/huruf kecil .....	96
2. Tokenize words and further clean-up text .....	98
3. Pemodelan Frase: Model Bigram dan Trigram.....	99
4. Remove Stopwords, Make Bigrams and Lemmatize....	101
5. Transformasi data: Korpus dan Kamus .....	102
J. Model Dasar.....	104
K. Penyetelan tuning hyper parameter .....	110
<b>DAFTAR PUSTAKA.....</b>	<b>120</b>

# **BAB 1**

## **PENDAHULUAN**

### **A. Definisi NLP**

Natural Language Processing (NLP) adalah cabang dari kecerdasan buatan yang berfokus pada interaksi antara komputer dan bahasa manusia yang alami. NLP memungkinkan komputer untuk memahami, menganalisis, memanipulasi, dan merespons bahasa manusia. Tujuannya adalah untuk memfasilitasi komunikasi yang mulus antara manusia dan mesin dengan menggunakan bahasa sehari-hari yang kita gunakan.

Definisi NLP mencakup sejumlah teknologi dan algoritma yang memungkinkan komputer untuk:

- **Pemrosesan Teks:** Memahami, menganalisis, dan menghasilkan teks bahasa manusia. Ini melibatkan segala hal mulai dari tokenisasi (memecah teks menjadi unit-unit yang lebih kecil), hingga analisis sintaksis dan semantik.
- **Pemahaman Bahasa:** Memungkinkan komputer untuk memahami konteks dari suatu teks, termasuk arti dari kata-kata, frasa, atau kalimat, serta makna yang lebih luas dalam sebuah konteks.
- **Generasi Bahasa:** Menghasilkan teks yang terstruktur secara gramatikal dan bermakna, seperti dalam pembuatan teks oleh chatbots atau penghasilan konten otomatis.
- **Penerjemahan Bahasa:** Menerjemahkan teks dari satu bahasa ke bahasa lain dengan mempertahankan arti dan konteksnya.

NLP memanfaatkan pendekatan statistik, mesin pembelajaran, dan pemrosesan bahasa alami untuk mencapai tujuannya. Seiring dengan perkembangan teknologi, model NLP yang menggunakan deep learning seperti transformer-based models (misalnya, BERT, GPT) telah menjadi semakin dominan dalam mengatasi banyak tugas NLP kompleks.



Penerapan NLP sangat luas, mulai dari aplikasi sederhana seperti koreksi ejaan hingga sistem yang kompleks seperti analisis sentimen, pengenalan ucapan, dan pembuatan chatbot. Ini juga memiliki peran yang penting dalam berbagai industri seperti kesehatan, keuangan, pendidikan, dan lainnya untuk mengoptimalkan pemrosesan data yang berkaitan dengan bahasa manusia.

NLP merupakan komponen penting dalam berbagai aplikasi perangkat lunak yang kita gunakan dalam kehidupan sehari-hari. Di bagian ini, kami akan memperkenalkan beberapa aplikasi utama dan juga melihat beberapa tugas umum yang akan Anda lihat di berbagai aplikasi NLP.

Aplikasi inti:

- Platform email, seperti Gmail, Outlook, dll., menggunakan NLP secara ekstensif untuk menyediakan serangkaian fitur produk, seperti klasifikasi spam, kotak masuk prioritas, ekstraksi acara kalender, pelengkapan otomatis, dll.
- Asisten berbasis suara, seperti Apple Siri, Google Assistant, Microsoft Cortana, dan Amazon Alexa mengandalkan serangkaian teknik NLP untuk berinteraksi dengan pengguna, memahami perintah pengguna, dan merespons dengan tepat.
- Mesin pencari modern, seperti Google dan Bing, yang merupakan landasan internet saat ini, banyak menggunakan NLP untuk berbagai subtugas, seperti pemahaman kueri, perluasan kueri, menjawab pertanyaan, pengambilan informasi, serta pemeringkatan dan pengelompokan hasil, untuk beberapa nama.
- Layanan terjemahan mesin, seperti Google Translate, Bing Microsoft Translator, dan Amazon Translate semakin banyak digunakan di dunia saat ini untuk menyelesaikan berbagai skenario dan kasus penggunaan bisnis.

Natural Language Processing (NLP) adalah cabang dari kecerdasan buatan yang memfokuskan pada interaksi antara komputer dan bahasa manusia, memungkinkan komputer untuk

memahami, menganalisis, dan merespon bahasa manusia dalam berbagai bentuknya. NLP mencakup sejumlah teknik dan algoritma untuk memproses teks dan ucapan, termasuk pemahaman bahasa, generasi bahasa, dan penerjemahan antarbahasa. Dalam hal ini, NLP berfungsi sebagai jembatan komunikasi antara manusia dan mesin, menggunakan bahasa sehari-hari yang kita gunakan.

Sebagai cabang dari kecerdasan buatan, NLP memungkinkan komputer untuk memahami, menganalisis, memanipulasi, dan merespons bahasa manusia, memfasilitasi komunikasi antara manusia dan mesin menggunakan bahasa alami (Agarwal, 2019). Teknologi NLP mencakup pemrosesan teks, di mana komputer diberdayakan untuk memahami dan menghasilkan teks bahasa manusia, melibatkan proses dari tokenisasi hingga analisis sintaksis dan semantic (Kjell et al., 2023).

Selain itu, pemahaman bahasa memungkinkan komputer untuk memahami konteks dari teks, termasuk makna kata, frasa, atau kalimat, serta makna yang lebih luas dalam konteks tertentu (Basha et al., 2023). Generasi bahasa merupakan aspek lain dari NLP, di mana komputer dapat menghasilkan teks yang terstruktur secara gramatikal dan bermakna, seperti yang digunakan oleh chatbots dan dalam penghasilan konten otomatis. Penerjemahan bahasa, sebagai bagian dari NLP, melibatkan menerjemahkan teks dari satu bahasa ke bahasa lain sambil mempertahankan makna dan konteksnya (Hirschberg & Manning, 2015)

NLP menggunakan pendekatan statistik, mesin pembelajaran, dan deep learning, termasuk model berbasis transformer seperti BERT dan GPT, untuk mengatasi berbagai tugas NLP yang kompleks (Mishra, 2019). Dengan perkembangan teknologi, model NLP yang menggunakan deep learning seperti model berbasis transformer telah menjadi semakin dominan dalam mengatasi banyak tugas NLP yang kompleks, merevolusi cara komputer memahami bahasa manusia.

Berikut adalah kerangka umum untuk sebuah tulisan yang membahas tentang Natural Language Processing (NLP):

## **B. Sejarah singkat NLP**

Sejarah Natural Language Processing (NLP) telah melalui serangkaian perkembangan yang menarik sejak awalnya. Inilah beberapa titik penting dalam sejarah NLP:

Awal Pengembangan tahun 1950-an dan 1960-an: Awalnya, NLP muncul sebagai bagian dari kecerdasan buatan. Pada tahun 1950-an, Alan Turing mengajukan pertanyaan dalam makalahnya yang terkenal, "Apakah mesin bisa berpikir?" yang membuka jalan bagi studi tentang kecerdasan buatan dan pemrosesan bahasa alami<sup>1</sup>. 1954: George Zipf, seorang linguistik, memberikan kontribusi awal dengan hukum Zipf yang menggambarkan distribusi kata dalam bahasa alami. 1950-an dan 1960-an: Pada masa ini, NLP berfokus pada penerjemahan mesin dan pengembangan model untuk pemahaman dan generasi bahasa.

Perkembangan Awal 1970-an hingga 1980-an: Era ini melihat peningkatan dalam penerapan aturan dan pendekatan statistik dalam NLP. Metode-metode seperti penguraian sintaksis berbasis aturan dan model statistik mulai digunakan. 1980-an: Munculnya sistem-sistem seperti SHRDLU (dikembangkan oleh Terry Winograd) yang memungkinkan komunikasi dengan komputer dalam bahasa alami, memberikan dorongan besar pada pengembangan NLP<sup>2</sup>.

Era Statistik dan Mesin Pembelajaran: 1990-an hingga 2000-an: Perkembangan metode-metode statistik semakin mendominasi. Teknik-teknik seperti HMM (Hidden Markov Models) dan algoritma pembelajaran mesin lainnya diperkenalkan untuk tugas-tugas NLP seperti pemodelan bahasa, penerjemahan, dan pengenalan ucapan. Akhir 2000-an: Dengan kemajuan komputasi

---

<sup>1</sup> <https://www.britannica.com/biography/Alan-Turing>

<sup>2</sup> [https://en.wikipedia.org/wiki/Computing\\_Machinery\\_and\\_Intelligence](https://en.wikipedia.org/wiki/Computing_Machinery_and_Intelligence)

dan pendekatan deep learning, NLP mengalami revolusi baru. Model-model seperti Word Embeddings (Word2Vec, GloVe) memungkinkan representasi kata yang lebih baik, sementara neural networks yang lebih dalam meningkatkan kinerja dalam tugas-tugas NLP (Vajjala et al., 2012).

Era Transformer dan Model Bahasa Besar: 2010-an hingga saat ini: Transformer, sebuah arsitektur neural network yang memanfaatkan self-attention mechanism, mengubah lanskap NLP<sup>3</sup>. Model-model besar seperti BERT (Bidirectional Encoder Representations from Transformers) dari Google dan GPT (Generative Pre-Trained Transformer) dari OpenAI mendefinisikan tingkat kinerja baru dalam pemahaman bahasa dan tugas-tugas NLP lainnya (Vaswani et al., 2017).

Seiring dengan perubahan teknologi, dataset yang lebih besar, dan kemajuan dalam algoritma, NLP terus berkembang pesat. Kemampuan untuk memahami, menghasilkan, dan berinteraksi dengan bahasa manusia semakin meningkat, memungkinkan penerapan NLP yang lebih luas dalam berbagai industri dan aplikasi sehari-hari.

### **C. Relevansi dan penerapan NLP dalam kehidupan**

relevansi dan penerapan Natural Language Processing (NLP) dalam kehidupan sehari-hari memiliki dampak yang signifikan dalam berbagai konteks. Berikut adalah beberapa contoh konkret:

1. Asisten Virtual dan Chatbot: Relevansi: Membantu dalam interaksi sehari-hari dengan teknologi. Contoh: Asisten pribadi seperti Siri, Google Assistant, atau Alexa memanfaatkan NLP untuk memahami perintah suara, menjadikan penggunaan perangkat teknologi lebih mudah. Chatbot yang terintegrasi dengan layanan pelanggan

---

<sup>3</sup> <https://blog.research.google/2017/08/transformer-novel-neural-network.html>

online juga menggunakan NLP untuk merespons pertanyaan pelanggan.

2. Pencarian Informasi: Relevansi: Mempercepat akses dan relevansi informasi. Contoh: Mesin pencari seperti Google menggunakan NLP untuk memahami pertanyaan pengguna dan menyajikan hasil pencarian yang relevan. Pemahaman konteks dan arti di balik pertanyaan memungkinkan hasil pencarian yang lebih akurat.
3. Analisis Sentimen Media Sosial: Relevansi: Memahami pandangan publik terhadap suatu topik atau merek. Contoh: Penggunaan NLP dalam analisis sentimen membantu perusahaan memahami umpan balik pelanggan di platform media sosial. Misalnya, menganalisis tweet atau postingan Facebook untuk mengukur reaksi terhadap produk tertentu.
4. Penerjemahan Bahasa: Relevansi: Memfasilitasi komunikasi lintas budaya. Contoh: Google Translate menggunakan teknologi NLP untuk menerjemahkan teks dari satu bahasa ke bahasa lain secara cepat dan akurat. Ini sangat membantu dalam komunikasi global di antara individu yang berbicara bahasa yang berbeda.
5. Otomatisasi Pekerjaan: Relevansi: Meningkatkan efisiensi pekerjaan dan tugas-tugas rutin. Contoh: Penggunaan NLP dalam pengelompokan email, pemrosesan formulir, atau analisis dokumen membantu dalam mengotomatisasi tugas-tugas yang membutuhkan pemahaman bahasa alami.
6. Sistem Pembelajaran Adaptif: Relevansi: Meningkatkan pengalaman belajar yang disesuaikan dengan individu. Contoh: Platform pembelajaran online menggunakan NLP untuk menyediakan materi yang disesuaikan dengan kebutuhan belajar masing-masing siswa. Ini memungkinkan pengalaman belajar yang lebih efektif.
7. Pelayanan Kesehatan yang Lebih Baik: Relevansi: Meningkatkan diagnosis dan pengelolaan informasi medis. Contoh: NLP digunakan dalam analisis rekam medis dan dokumen kesehatan untuk membantu dokter dalam membuat diagnosis yang lebih cepat dan akurat, serta memantau perkembangan pasien.

8. Keamanan dan Penegakan Hukum: Relevansi: Mendeteksi ancaman potensial dan kegiatan kriminal. Contoh: Analisis teks pada platform online untuk mendeteksi ancaman atau kegiatan yang mencurigakan, membantu dalam keamanan siber dan penegakan hukum.

Penerapan NLP dalam kehidupan sehari-hari telah mengubah cara kita berinteraksi dengan teknologi, membawa kemudahan dalam akses informasi, dan meningkatkan efisiensi dalam banyak aspek kehidupan (Tunstall et al., 2022).

## **D. Tujuan**

Tujuan penulisan mengenai Natural Language Processing (NLP) adalah untuk menyampaikan pemahaman yang mendalam tentang konsep, perkembangan, penerapan, dan relevansi NLP dalam berbagai bidang kehidupan. Beberapa tujuan spesifiknya meliputi:

1. Pemahaman Konsep NLP: menjelaskan tentang apa itu NLP, konsep dasar di baliknya, dan bagaimana komputer dapat memahami, memproses, dan menghasilkan bahasa manusia dengan bantuan teknologi.
2. Menjelaskan Teknik dan Algoritma NLP: Memberikan pemahaman tentang teknik-teknik utama dalam NLP seperti tokenisasi, pengenalan entitas, penerjemahan, analisis sentimen, dan model-model seperti transformer-based models. Tujuannya adalah agar pembaca memahami bagaimana NLP diimplementasikan dalam konteks yang berbeda.
3. Penerapan dalam Berbagai Bidang: Menjelaskan penerapan NLP dalam industri, kesehatan, pendidikan, keamanan, dan bidang lainnya. Menyoroti kontribusi NLP dalam mengoptimalkan proses, meningkatkan efisiensi, dan menghadirkan solusi dalam setiap bidang ini.
4. Tantangan dan Isu Etika: Membahas tantangan teknis dalam pengembangan NLP, seperti ambiguitas bahasa dan kurangnya data berkualitas, serta menyoroti isu-isu etika seperti privasi data, bias dalam model, dan tanggung jawab sosial dalam penggunaan teknologi NLP.

5. Edukasi dan Informasi Masyarakat: Memberikan pemahaman yang lebih luas kepada masyarakat umum tentang bagaimana NLP memengaruhi kehidupan sehari-hari mereka, baik dalam penggunaan aplikasi yang dikenal seperti asisten virtual maupun implikasinya dalam industri dan pengembangan teknologi masa depan.
6. Mendorong Pengembangan dan Penelitian Lanjutan: Merangsang minat dan keingintahuan pembaca untuk terlibat dalam pengembangan lebih lanjut dalam bidang NLP, serta mendorong penelitian yang lebih mendalam dalam pengembangan teknologi bahasa.
7. Menggugah Kesadaran akan Potensi dan Tantangan: Memberikan pemahaman yang utuh tentang potensi luar biasa NLP dalam merubah cara kita berinteraksi dengan teknologi, sambil menyadari tantangan teknis dan etika yang terkait.

Tujuan utama dari penulisan tentang NLP adalah untuk memberikan pemahaman menyeluruh tentang konsep ini, memberikan wawasan tentang peran dan pengaruhnya dalam kehidupan sehari-hari, serta menyampaikan tantangan dan peluang yang terkait dengan penggunaan teknologi bahasa ini.

## **BAB II**

### **DASAR-DASAR NLP**

#### **A. Pengertian dasar tentang bahasa alami**

Bahasa alami adalah bahasa yang digunakan oleh manusia untuk berkomunikasi sehari-hari. Bahasa alami bersifat kompleks dan fleksibel, dan dapat digunakan untuk berbagai tujuan, seperti menyampaikan informasi, mengekspresikan emosi, dan membangun hubungan. Ini adalah cara alami di mana manusia menyampaikan ide, emosi, informasi, dan instruksi kepada orang lain menggunakan kata-kata, frasa, kalimat, dan struktur bahasa yang kompleks (Tunstall et al., 2022).

Bahasa alami terdiri dari berbagai komponen, antara lain:

- Morfologi adalah cabang linguistik yang mempelajari bentuk kata. Morfologi meliputi pembentukan kata, seperti proses pengimbuhan, pengulangan, dan penggabungan.
- Sintaksis adalah cabang linguistik yang mempelajari struktur kalimat. Sintaksis meliputi aturan-aturan yang mengatur bagaimana kata-kata dapat digabungkan menjadi kalimat yang bermakna.
- Semantik adalah cabang linguistik yang mempelajari makna kata dan kalimat. Semantik meliputi hubungan antara kata dan dunia nyata.
- Pragmatik adalah cabang linguistik yang mempelajari penggunaan bahasa dalam konteks tertentu. Pragmatik meliputi cara penggunaan bahasa untuk menyampaikan maksud dan tujuan tertentu.

Perbedaan bahasa alami dan bahasa formal: Bahasa alami berbeda dengan bahasa formal, seperti bahasa pemrograman. Bahasa formal memiliki aturan yang ketat dan tidak fleksibel, sedangkan bahasa alami memiliki aturan yang lebih longgar dan fleksibel. Bahasa formal digunakan untuk tujuan-tujuan tertentu, seperti menulis program komputer, sedangkan bahasa alami digunakan



untuk berbagai tujuan, seperti berkomunikasi, mengekspresikan emosi, dan membangun hubungan.

Penerapan bahasa alami: Bahasa alami diterapkan dalam berbagai bidang, antara lain: Bahasa alami digunakan dalam berbagai produk dan layanan teknologi, seperti mesin penerjemah, asisten virtual, dan chatbot. Bahasa alami digunakan untuk membuat konten pembelajaran yang lebih menarik dan interaktif, serta untuk memberikan umpan balik yang lebih personal kepada siswa. Bahasa alami digunakan untuk mendiagnosis penyakit, untuk memberikan layanan kesehatan yang lebih personal, dan untuk meningkatkan penelitian medis. Bahasa alami digunakan untuk meningkatkan layanan pelanggan, untuk meningkatkan pemasaran, dan untuk membuat keputusan bisnis yang lebih baik. Bahasa alami adalah alat komunikasi yang penting bagi manusia. Bahasa alami memiliki berbagai komponen dan digunakan dalam berbagai bidang.

Beberapa aspek penting dalam pengertian dasar tentang bahasa alami meliputi:

- **Kompleksitas dan Struktur:** Bahasa alami memiliki struktur kompleks yang terdiri dari unsur-unsur seperti fonem (unit bunyi), morfem (unit makna terkecil), kata-kata, frasa, kalimat, dan aturan sintaksis yang mempengaruhi arti dari teks atau ucapan.
- **Kekayaan dan Produktivitas:** Bahasa alami memiliki sifat kekayaan yang memungkinkan pembicara untuk menghasilkan kombinasi yang tak terbatas dari kata-kata dan kalimat untuk menyampaikan pesan yang berbeda dalam situasi yang berbeda.
- **Makna dan Konteks:** Makna dalam bahasa alami sangat bergantung pada konteks penggunaannya. Kata atau frase dapat memiliki makna yang berbeda-beda tergantung pada situasi dan bagaimana mereka digunakan dalam kalimat atau percakapan.
- **Ambiguitas:** Bahasa alami sering kali mengandung ambiguitas, di mana satu kata atau frasa dapat memiliki beberapa makna yang berbeda. Ini dapat menjadi

tantangan dalam pemrosesan bahasa alami karena memerlukan pemahaman yang tepat dari konteksnya.

- Perubahan dan Variasi: Bahasa alami terus berkembang dan bervariasi dari waktu ke waktu serta antar komunitas bahasa. Perubahan ini dapat terjadi dalam kosakata, tata bahasa, atau penggunaan kata-kata yang baru.

Dalam konteks Natural Language Processing (NLP), pemahaman tentang bahasa alami menjadi dasar utama. Mesin atau komputer yang memproses bahasa alami harus mampu memahami kompleksitas struktur bahasa, mengatasi ambiguitas, dan memperhitungkan konteks dalam rangka melakukan tugas-tugas seperti pemrosesan teks, analisis sentimen, penerjemahan, dan lainnya. Menciptakan model dan algoritma yang dapat memahami dan memanipulasi bahasa alami dengan baik adalah inti dari pengembangan dalam NLP.

## **B. Teknik-teknik dasar dalam NLP**

Teknik-teknik dasar dalam Natural Language Processing (NLP) adalah fondasi penting untuk memproses, menganalisis, dan memahami bahasa manusia. Berikut adalah beberapa teknik dasar yang sering digunakan dalam NLP (Atkinson-Abutridy, 2022):

- Tokenisasi: Proses memecah teks menjadi unit-unit yang lebih kecil, seperti kata-kata, frasa, atau kalimat. Contoh: Mengubah kalimat menjadi kumpulan kata-kata individual.
- Stopword Removal: Menghapus kata-kata umum yang tidak memberikan nilai tambah dalam analisis teks, seperti "dan", "atau", "di", dll. Dalam analisis sentimen, kata-kata ini sering dihapus karena kurangnya kontribusi terhadap penilaian keseluruhan.
- Stemming dan Lemmatisasi: Stemming: Proses mengubah kata-kata menjadi bentuk dasar atau akar kata dengan menghilangkan imbuhan. Lemmatisasi: Proses mengubah kata-kata ke bentuk dasarnya (kata baku) berdasarkan kamus atau aturan linguistik. Contoh: Mengubah kata-kata

"berlari", "lari", dan "lari-lari" menjadi bentuk dasarnya "lari".

- Part-of-Speech (POS) Tagging: Mengidentifikasi dan menandai jenis kata dalam sebuah kalimat, seperti kata benda, kata kerja, kata sifat, dan lain-lain. Contoh: Menandai kata "makan" sebagai kata kerja dan "rumah" sebagai kata benda.
- Named Entity Recognition (NER): Mendeteksi dan menandai entitas penting dalam teks seperti nama orang, tempat, tanggal, organisasi, dll. Contoh: Mengidentifikasi "Bill Gates" sebagai nama orang atau "Microsoft" sebagai nama perusahaan.
- Word Embeddings dan Word Vectors: Representasi vektor yang menyandikan makna kata-kata dalam bentuk numerik, memungkinkan model untuk memahami hubungan antar kata. Contoh: Model Word2Vec atau GloVe menghasilkan vektor yang merepresentasikan kata-kata dan maknanya dalam ruang vektor.
- Analisis Sintaksis dan Semantik: Sintaksis: Menganalisis struktur gramatikal dari kalimat, seperti hubungan antara kata-kata dalam sebuah kalimat. Semantik: Memahami makna dari kalimat atau teks, termasuk hubungan makna antara kata-kata.
- Analisis Sentimen: Mengidentifikasi, mengekstrak, atau menganalisis sentimen dari teks, seperti apakah sebuah ulasan bersifat positif, negatif, atau netral. Contoh: Menilai apakah sebuah ulasan produk di platform e-commerce adalah positif atau negatif. Teknik-teknik ini membentuk dasar dalam pengolahan bahasa alami. Dalam kombinasi dengan model-model machine learning atau deep learning, teknik-teknik ini memungkinkan mesin untuk memahami, menganalisis, dan menghasilkan teks dalam cara yang semakin mirip dengan cara manusia.

## 1. Tokenisasi

Tokenisasi adalah proses memecah teks menjadi unit-unit yang lebih kecil, seperti kata-kata, frasa, atau kalimat yang disebut token. Setiap token mewakili bagian terpisah dari teks yang dapat

dianggap sebagai unit yang berarti. Proses tokenisasi sangat penting dalam Natural Language Processing (NLP) karena membantu dalam analisis, pemrosesan, dan pemahaman teks oleh mesin.

**Konsep Tokenisasi:** Unit Token: Token bisa berupa kata, frasa, kalimat, atau bahkan karakter tergantung pada kebutuhan analisis atau pemrosesan yang dilakukan. Pemisahan: Teks dapat dipisahkan menjadi token berdasarkan spasi (untuk kata-kata), tanda baca, atau aturan tertentu seperti tokenisasi berdasarkan kata. Pembersihan dan Normalisasi: Tokenisasi dapat melibatkan pembersihan teks dari karakter khusus, tanda baca, dan normalisasi huruf menjadi huruf kecil atau huruf besar untuk konsistensi. Berikut teknik untuk Melakukan Tokenisasi:

**Tokenisasi Berdasarkan Spasi:**

Pemisahan teks menjadi token berdasarkan spasi antara kata-kata.

Contoh: "Saya sedang belajar NLP." akan menjadi token: ["Saya", "sedang", "belajar", "NLP", "."]

**Tokenisasi Berdasarkan Kata:**

Memisahkan teks berdasarkan aturan kata, mengabaikan tanda baca atau spasi.

Contoh: "Dia tidak suka berjalan-jalan." akan menjadi token: ["Dia", "tidak", "suka", "berjalan-jalan", "."]

**Tokenisasi Berdasarkan Frasa atau Kalimat:**

Memisahkan teks menjadi token berdasarkan frasa atau kalimat.

Contoh: "Saya belajar NLP. Ini menarik!" akan menjadi token kalimat: ["Saya belajar NLP.", "Ini menarik!"]

**Tokenisasi dengan Penggunaan Algoritma Khusus:**

Menggunakan aturan linguistik atau algoritma yang lebih kompleks untuk memisahkan teks menjadi token.

Contoh: Algoritma seperti WordPunctTokenizer dalam Python yang membagi teks menjadi kata-kata dan tanda baca sebagai token terpisah. Langkah-langkah dalam Proses Tokenisasi:

- a) Pemisahan Teks: Pisahkan teks menjadi unit-unit yang relevan berdasarkan jenis token yang diinginkan (kata, frasa, kalimat).
- b) Pembersihan: Hilangkan karakter khusus, normalisasikan huruf, dan lakukan pre-processing lainnya jika diperlukan.
- c) Representasi dalam Bentuk Token: Hasil tokenisasi diwakili dalam bentuk daftar atau struktur data lainnya yang menyimpan token-token tersebut.
- d) Pentingnya Tokenisasi: Memungkinkan mesin untuk memproses teks dalam format yang dapat dipahami dan diolah. Dasar untuk langkah-langkah pemrosesan NLP lainnya seperti analisis sintaksis, pembangunan model, atau analisis sentimen. Membantu dalam mempersiapkan data untuk berbagai tugas NLP seperti machine translation, analisis teks, dan lainnya. Tokenisasi merupakan langkah penting dalam pemrosesan teks dalam NLP yang membantu dalam memahami dan memanipulasi bahasa manusia secara efisien oleh komputer atau mesin.

## 2. Stopword removal

Stopword removal adalah proses menghilangkan kata-kata umum yang tidak memberikan kontribusi signifikan terhadap makna dalam analisis teks. Kata-kata semacam ini, seperti "dan", "atau", "di", "yang", dan lainnya, sering muncul dalam teks tetapi cenderung tidak memiliki nilai informasi yang besar dalam pemrosesan atau analisis teks. Proses penghapusan stopwords sangat penting dalam tahap pra-pemrosesan dalam NLP untuk meningkatkan kualitas analisis dan model yang dibangun. Berikut adalah konsep Stopword Removal:

- a) Kata-kata Umum: Stopwords adalah kata-kata umum yang sering ditemukan dalam bahasa namun jarang membawa makna khusus atau signifikan dalam analisis teks.
- b) Mengurangi Noise: Penghapusan stopwords membantu mengurangi noise dalam data teks. Ini memungkinkan

fokus pada kata-kata yang lebih penting dalam pemrosesan.

- c) Pembersihan Teks: Proses ini melibatkan menghapus stopwords dari teks, meninggalkan hanya kata-kata yang dianggap lebih relevan dalam analisis.

### 3. Stemming dan Lemmatisasi

Stemming dan lemmatisasi adalah dua teknik dalam pemrosesan bahasa alami yang bertujuan untuk mengubah kata-kata menjadi bentuk dasar atau kata baku agar lebih mudah dianalisis atau dipahami oleh mesin.

Stemming:

Konsep: Stemming adalah proses menghilangkan imbuhan atau akhiran kata untuk mendapatkan bentuk dasar atau stem dari sebuah kata. Tujuan: Mengurangi kata-kata ke bentuk dasarnya sehingga kata-kata yang memiliki akar kata yang sama dapat dianggap sebagai bentuk yang sama. Teknik: Stemming menggunakan aturan sederhana untuk menghapus imbuhan kata. Namun, hasil stemming tidak selalu merupakan kata yang benar dalam bahasa yang sesungguhnya. Contoh: "Berlari", "lari", "lari-lari" akan diubah menjadi bentuk dasarnya "lar".

Lemmatisasi:

Konsep: Lemmatisasi adalah proses mengubah kata-kata menjadi bentuk dasar atau kata baku berdasarkan kamus atau aturan linguistik. Tujuan: Menghasilkan kata yang merupakan bentuk dasar kata dalam bahasa yang tepat. Teknik: Lemmatisasi menggunakan kamus kata-kata yang telah didefinisikan dan aturan linguistik untuk mengubah kata-kata menjadi bentuk dasar. Contoh: "Berlari", "lari", "lari-lari" akan diubah menjadi bentuk dasarnya "lari".

Perbandingan Antara Stemming dan Lemmatisasi:

Stemming: Lebih sederhana karena hanya menghapus imbuhan untuk mendapatkan akar kata yang mungkin tidak selalu benar.

Lemmatisasi: Lebih kompleks karena memerlukan pengetahuan tentang kamus dan struktur bahasa untuk mengubah kata menjadi bentuk yang benar secara linguistik.

Langkah-langkah dalam Stemming dan Lemmatisasi:

Tokenisasi: Pemecahan teks menjadi token (kata-kata).

Proses Stemming atau Lemmatisasi: Penggunaan algoritma atau aturan linguistik untuk mengubah kata-kata menjadi bentuk dasar. Pentingnya Stemming dan Lemmatisasi: Normalisasi Teks: Membantu dalam normalisasi teks untuk analisis yang lebih akurat. Reduksi Redundansi: Mengurangi redundansi kata-kata yang memiliki akar yang sama. Baik stemming maupun lemmatisasi digunakan untuk menyederhanakan kata-kata menjadi bentuk dasar atau kata baku untuk memfasilitasi analisis teks dalam NLP. Meskipun tidak sempurna, keduanya membantu dalam memproses dan memahami teks dalam analisis bahasa alami.

#### 4. Part-of-Speech (POS) Tagging

Part-of-Speech (POS) tagging adalah proses yang dilakukan dalam Natural Language Processing (NLP) untuk mengidentifikasi jenis kata dalam sebuah kalimat, seperti kata benda, kata kerja, kata sifat, kata tanya, dan lainnya. Ini penting dalam pemahaman makna dan struktur kalimat dalam bahasa alami.

Konsep POS Tagging:

Jenis Kata: Setiap kata dalam sebuah kalimat memiliki peran atau fungsi tertentu dalam kalimat tersebut, misalnya sebagai subjek, predikat, atau objek.

Tujuan: POS tagging bertujuan untuk mengidentifikasi peran atau fungsi setiap kata dalam kalimat untuk memahami struktur dan makna kalimat.

Teknik untuk Melakukan POS Tagging:

Menggunakan Model Statistik: Memanfaatkan model statistik seperti Hidden Markov Models (HMM) atau Conditional Random

Fields (CRF) untuk memprediksi jenis kata berdasarkan konteks kalimat.

Pemanfaatan Kamus Kata: Menggunakan kamus kata-kata yang sudah diberi label jenis kata untuk mengidentifikasi tipe kata.

Penggunaan Algoritma Berbasis Aturan: Penggunaan aturan linguistik dan struktur bahasa untuk menentukan jenis kata dalam konteks kalimat.

Langkah-langkah dalam Proses POS Tagging:

Tokenisasi: Pemisahan kata-kata dalam kalimat menjadi token.

Ekstraksi Fitur: Mendapatkan fitur-fitur dari kata-kata yang dapat membantu dalam prediksi jenis kata, seperti kata sebelumnya atau kata-kata yang terkait dalam kalimat.

Pemodelan: Menggunakan model (statistik atau berbasis aturan) untuk memprediksi jenis kata untuk setiap token dalam kalimat.

Labeling: Memberikan label atau tag untuk setiap kata dalam kalimat berdasarkan jenis kata yang diprediksi.

Contoh POS Tagging:

Dalam kalimat "Ani membaca buku di perpustakaan," hasil POS taggingnya mungkin seperti:

"Ani" -> Noun (Kata benda)

"membaca" -> Verb (Kata kerja)

"buku" -> Noun (Kata benda)

"di" -> Preposition (Kata depan)

"perpustakaan" -> Noun (Kata benda)

Pentingnya POS Tagging:

Pemahaman Struktur Kalimat: Memahami struktur kalimat dan hubungan antara kata-kata dalam konteks kalimat.

Pemrosesan Bahasa yang Lebih Lanjut: Membantu dalam berbagai tugas NLP seperti analisis sintaksis, parsing, dan terjemahan.

POS Tagging membantu dalam memahami makna dan struktur kalimat, yang merupakan langkah penting dalam analisis dan



pemrosesan teks dalam NLP. Ini memungkinkan mesin untuk memahami peran dan hubungan antar kata-kata dalam konteks kalimat.

## 5. Named Entity Recognition (NER)

Named Entity Recognition (NER) adalah proses dalam Natural Language Processing (NLP) yang bertujuan untuk mengidentifikasi dan menandai entitas penting dalam teks seperti nama orang, tempat, tanggal, organisasi, dan lainnya. Tujuannya adalah untuk mengenali dan mengekstrak informasi yang relevan dari teks yang dapat dianggap sebagai entitas (Li et al., 2022).

NER adalah teknik penting dalam NLP karena memungkinkan pengenalan dan penandaan entitas penting dalam teks, yang mendukung berbagai tugas pemrosesan bahasa alami dan analisis teks yang lebih lanjut. Named Entity Recognition (NER) merupakan salah satu tugas dasar dalam Natural Language Processing (NLP). NER bertujuan untuk mengidentifikasi dan mengklasifikasikan entitas bernama dalam teks, seperti nama orang, nama organisasi, lokasi, tanggal, dan waktu.

Pengenalan entitas dinamis telah berkembang secara signifikan dalam dekade terakhir, dengan penelitian terbaru yang semakin banyak mengadopsi pembelajaran mendalam, pembelajaran transfer, basis pengetahuan, dan metode lainnya. Penelitian NER untuk bahasa-bahasa sumber daya rendah juga meningkat pesat (Sun et al., 2018). Salah satu tantangan utama dalam NER adalah mengatasi entitas bersarang dan tumpang tindih, di mana token yang sama bisa menjadi bagian dari lebih dari satu kategori entitas. Strategi baru yang diusulkan melibatkan memformalkan tugas ekstraksi entitas sebagai tugas pemahaman bacaan berbasis kueri, di mana tugas mengekstraksi entitas dengan PER diformalkan sebagai menjawab pertanyaan "orang mana yang disebutkan dalam teks?" (Meng et al., 2019).

Dalam hal pengembangan model NER, telah diperkenalkan S-NER, model NER berbasis rentang yang terlebih dahulu membagi teks mentah menjadi rentang teks dan menganggapnya sebagai

kandidat entitas. Ini kemudian langsung memperoleh jenis rentang dengan melakukan klasifikasi jenis entitas pada representasi semantik rentang, yang menghilangkan kebutuhan akan ketergantungan label (Yu et al., 2022). Metode yang diusulkan ini telah menunjukkan peningkatan signifikan dalam performa NER, terutama dalam menangani jumlah contoh pendukung yang rendah, yang menyoroti pentingnya mengadaptasi strategi NER untuk domain dan tantangan baru (Ziyadi et al., 2020). Dengan terus berkembangnya metode dan pendekatan baru dalam NER, penelitian di bidang ini tetap menjadi area yang aktif dan penting dalam pemrosesan bahasa alami, memberikan kontribusi signifikan terhadap berbagai aplikasi NLP yang bergantung pada pemahaman teks yang akurat dan mendalam (Hannon et al., 2024)

**Konsep Named Entity Recognition (NER):**

**Entitas Nama:** Dalam sebuah teks, entitas nama adalah segmen yang merujuk kepada orang, tempat, tanggal, organisasi, dan entitas penting lainnya yang memiliki makna spesifik dalam konteks yang diberikan. Tujuannya untuk Mengidentifikasi dan menandai entitas penting ini dalam teks untuk menggali informasi yang berguna. Teknik untuk Melakukan Named Entity Recognition (NER), **Pemodelan Berbasis Aturan:** Menggunakan aturan linguistik atau pola tertentu untuk mengidentifikasi entitas nama dalam teks.

**Pemanfaatan Machine Learning:** Menerapkan pendekatan machine learning seperti Conditional Random Fields (CRF) atau deep learning menggunakan model seperti Recurrent Neural Networks (RNN) atau Transformer untuk mengidentifikasi entitas.

**Langkah-langkah dalam Proses Named Entity Recognition (NER):**

1. Tokenisasi: Pemisahan teks menjadi token atau kata-kata.
2. POS Tagging: Identifikasi jenis kata-kata dalam teks menggunakan Part-of-Speech tagging.
3. Feature Extraction: Ekstraksi fitur yang relevan, seperti kata sebelumnya, jenis kata, atau pola tertentu dalam teks.

4. Pemodelan dan Prediksi: Menggunakan model machine learning atau aturan tertentu untuk memprediksi dan menandai entitas dalam teks.

Contoh Named Entity Recognition (NER):

Dalam kalimat "Mark Zuckerberg mendirikan Facebook pada tahun 2004 di Harvard University," hasil NER-nya mungkin seperti:

"Mark Zuckerberg" -> Nama Orang (Person)

"Facebook" -> Organisasi (Organization)

"2004" -> Tanggal (Date)

"Harvard University" -> Lokasi (Location)

Pentingnya Named Entity Recognition (NER):

Ekstraksi Informasi: Membantu dalam mengekstrak informasi penting dari teks seperti nama, tempat, atau tanggal.

Analisis Data: Memungkinkan pemrosesan lebih lanjut untuk analisis data seperti analisis sentimen, klasifikasi teks, dan lainnya.

Entitas bernama adalah frasa benda (noun phrase) yang memiliki tipe spesifik. Misalnya, "John Doe" adalah nama orang, "Microsoft" adalah nama organisasi, "Jakarta" adalah lokasi, "2023-12-13" adalah tanggal, dan "10:00" adalah waktu. Teknik untuk melakukan NER. Ada dua pendekatan utama untuk melakukan NER, yaitu pendekatan berbasis aturan dan Pendekatan berbasis pembelajaran mesin.

1. Pendekatan berbasis aturan: Pendekatan berbasis aturan menggunakan serangkaian aturan untuk mengidentifikasi dan mengklasifikasikan entitas bernama. Aturan-aturan ini biasanya dibuat oleh ahli bahasa berdasarkan pengetahuan mereka tentang bahasa.
2. Keuntungan dari pendekatan berbasis aturan: Efektif untuk kasus-kasus yang sederhana dan terdefinisi dengan baik. Dapat digunakan untuk bahasa yang tidak memiliki data pelatihan yang besar
3. Kerugian dari pendekatan berbasis aturan: Sulit untuk membuat aturan yang lengkap dan akurat untuk semua kasus.

Sulit untuk mengadaptasi aturan untuk bahasa yang baru atau berubah.

Pendekatan berbasis pembelajaran mesin. Pendekatan berbasis pembelajaran mesin menggunakan model pembelajaran mesin untuk mengidentifikasi dan mengklasifikasikan entitas bernama. Model pembelajaran mesin dilatih pada data pelatihan yang berisi contoh entitas bernama. Keuntungan dari pendekatan berbasis pembelajaran mesin: Dapat menangani kasus-kasus yang kompleks dan tidak terdefinisi dengan baik. Dapat beradaptasi dengan data pelatihan yang baru atau berubah.

Kerugian dari pendekatan berbasis pembelajaran mesin: Membutuhkan data pelatihan yang besar. Dapat menghasilkan hasil yang tidak akurat jika data pelatihan tidak representative. Kombinasi pendekatan berbasis aturan dan pembelajaran mesin. Pendekatan berbasis aturan dan pembelajaran mesin dapat dikombinasikan untuk meningkatkan akurasi NER. Misalnya, pendekatan berbasis aturan dapat digunakan untuk mengidentifikasi entitas bernama yang umum, dan pendekatan berbasis pembelajaran mesin dapat digunakan untuk mengidentifikasi entitas bernama yang tidak umum. Pencarian informasi. NER dapat digunakan untuk mengidentifikasi entitas bernama dalam dokumen teks. Misalnya, NER dapat digunakan untuk mengidentifikasi nama orang, nama organisasi, dan lokasi dalam dokumen berita. Pemrosesan bahasa alami. NER dapat digunakan untuk berbagai tugas pemrosesan bahasa alami, seperti:

\* \*\*Peringkasan teks\*\*

\* \*\*Pertanyaan jawab\*\*

\* \*\*Pemahaman bahasa alami\*\*

Aplikasi bisnis: NER dapat digunakan untuk berbagai aplikasi bisnis, seperti:

\* \*\*Analisis sentimen\*\*

\* \*\*Pencegahan penipuan\*\*

\* \*\*Pemasarkan\*\*

Named Entity Recognition (NER) adalah tugas dasar dalam Natural Language Processing (NLP). NER bertujuan untuk mengidentifikasi dan mengklasifikasikan entitas bernama dalam teks. Ada dua pendekatan utama untuk melakukan NER, yaitu: Pendekatan berbasis aturan dan Pendekatan berbasis pembelajaran mesin. Pendekatan berbasis aturan efektif untuk kasus-kasus yang sederhana dan terdefinisi dengan baik, sedangkan pendekatan berbasis pembelajaran mesin dapat menangani kasus-kasus yang kompleks dan tidak terdefinisi dengan baik. Pendekatan berbasis aturan dan pembelajaran mesin dapat dikombinasikan untuk meningkatkan akurasi NER. NER memiliki berbagai penerapan, antara lain: Pencarian informasi. Pemrosesan bahasa alami. Aplikasi bisnis

## 6. Word Embeddings dan Word Vectors

Word Embeddings dan Word Vectors adalah teknik penting dalam Natural Language Processing (NLP) yang digunakan untuk merepresentasikan kata-kata dalam bentuk vektor numerik dalam ruang dimensi yang lebih rendah. Representasi ini memungkinkan mesin untuk memahami dan memanipulasi makna kata-kata dalam pemrosesan bahasa alami.

Word embeddings dan word vectors dalam Natural Language Processing (NLP) memiliki peran penting karena kemampuan mereka untuk mengkodekan hubungan antar kata dalam ruang vektor. Hal ini bermanfaat untuk berbagai tugas pemrosesan bahasa, dari komponen dalam sistem NLP hingga alat untuk analisis linguistik dalam studi bahasa dan literatur. Menginterpretasikan embeddings dan memahami hubungan gramatikal dan semantik yang dikodekan di dalamnya berguna namun menantang. Visualisasi dapat membantu dalam interpretasi embeddings tersebut (Heimerl & Gleicher, 2018).

Word embeddings merupakan teknik pembelajaran fitur yang memetakan kata-kata dari kosakata ke dalam vektor bilangan riil dalam ruang berdimensi rendah. Dengan memanfaatkan korpus teks tanpa label yang besar, representasi ruang kontinu ini dapat

dihitung untuk menangkap informasi sintaktis dan semantik tentang kata-kata. Ketika digunakan sebagai representasi input dasar, word embeddings telah terbukti menjadi aset besar untuk berbagai tugas NLP. Teknik-teknik terkini untuk mendapatkan word embeddings sebagian besar berbasis pada model bahasa neural network (NNLM), di mana vektor kata diinisialisasi secara acak dan kemudian dilatih untuk memprediksi konteks di mana kata-kata yang bersangkutan cenderung muncul (Gavhane et al., 2022).

Dalam pemrosesan Bahasa Gujarati, yang merupakan bahasa dengan sumber daya rendah, word2vec dan fastText merupakan beberapa teknik word embeddings yang paling umum. Sementara banyak pekerjaan telah dilakukan untuk mendapatkan embeddings dalam bahasa dengan sumber daya kaya seperti Bahasa Inggris, masih ada pekerjaan yang harus dilakukan untuk bahasa dengan sumber daya rendah. Fokus pada pengembangan vektor kata untuk bahasa Gujarati dan penyiapan dataset tes analogi untuk mengevaluasi akurasi embeddings yang diperoleh telah dilakukan. Kinerja model juga dibandingkan dengan model Gujarati pra-latih yang sudah tersedia (Joshi et al., 2019).

Word embeddings (representasi vektor kata yang didistribusikan) telah menjadi komponen penting dalam banyak tugas pemrosesan bahasa alami (NLP) seperti terjemahan mesin, analisis sentimen, analogi kata, pengenalan entitas bernama, dan kesamaan kata. Meskipun demikian, pekerjaan terkini hanya menyediakan vektor kata untuk bahasa Hausa yang dilatih menggunakan fastText, terdiri dari hanya beberapa vektor kata. Penelitian ini menyajikan model embeddings kata menggunakan model Continuous Bag of Words (CBoW) dan Skip Gram dari Word2Vec. Model-model ini, hauWE (Hausa Words Embedding), lebih besar dan lebih baik dari model sebelumnya, membuatnya lebih berguna dalam tugas-tugas NLP. Untuk membandingkan model, mereka digunakan untuk memprediksi 10 kata yang paling mirip dengan 30 kata Hausa yang dipilih secara acak. hauWE CBoW dengan akurasi prediksi 88,7% dan hauWE SG dengan 79,3% jauh melampaui performa model [1] dengan 22,3% (Abdulummin & Galadanci, 2019).

Konsep Word Embeddings dan Word Vectors:

1. **Representasi Numerik:** Kata-kata dalam teks diubah menjadi vektor numerik dalam ruang dimensi yang lebih rendah.
2. **Makna dan Hubungan:** Vektor kata-kata yang serupa atau terkait secara semantik ditempatkan lebih dekat satu sama lain dalam ruang vektor.

Teknik untuk Melakukan Word Embeddings:

**Word2Vec:** Salah satu teknik paling terkenal yang mempelajari representasi vektor kata-kata dengan memanfaatkan jaringan saraf tiruan. **GloVe** (Global Vectors for Word Representation): Teknik lain yang menggabungkan informasi dari matriks co-occurrence kata-kata dalam corpus untuk menghasilkan representasi vektor kata. Langkah-langkah dalam Proses Word Embeddings:

1. Pra-Pemrosesan: Pra-pemrosesan teks seperti tokenisasi, penghapusan stopwords, stemming, atau lemmatisasi.
2. Pembuatan Model: Pembuatan model Word2Vec atau GloVe dengan menggunakan data teks yang besar.
3. Pelatihan Model: Melatih model pada teks yang digunakan untuk membuat representasi vektor kata-kata.
4. Contoh Word Embeddings: Misalkan representasi vektor untuk kata-kata "king" dan "queen". Dalam ruang vektor, mereka mungkin memiliki hubungan yang serupa dengan kata "royal" atau "throne" karena keterkaitan semantiknya.
5. Pentingnya Word Embeddings: Semantik yang Lebih Dalam: Memungkinkan mesin untuk memahami hubungan dan makna antara kata-kata dalam konteks.
6. Pemrosesan Bahasa yang Lebih Baik: Meningkatkan kinerja model dalam berbagai tugas NLP seperti analisis sentimen, penerjemahan, dan klasifikasi teks.

Word Embeddings menjadi kunci dalam NLP karena memungkinkan representasi kata-kata dalam ruang vektor numerik yang memperhitungkan hubungan semantik dan makna. Representasi ini memperkaya pemahaman mesin terhadap bahasa manusia dan mendukung kinerja model dalam berbagai tugas NLP.

Ada berbagai teknik yang dapat digunakan untuk melakukan word embeddings dan word vectors. Beberapa teknik yang umum digunakan adalah: Skip-gram, Continuous Bag-of-Words (CBOW), Glove, Word2Vec

Skip-gram adalah teknik yang menggunakan model pembelajaran mesin untuk memprediksi kata-kata di sekitar kata tertentu. Misalnya, model skip-gram akan dilatih untuk memprediksi kata "makan" jika kata "nasi" muncul di sekitarnya. CBOW adalah teknik yang menggunakan model pembelajaran mesin untuk memprediksi kata tertentu berdasarkan kata-kata di sekitarnya. Misalnya, model CBOW akan dilatih untuk memprediksi kata "nasi" jika kata "makan" dan "ayam" muncul di sekitarnya. Glove adalah teknik yang menggunakan metode statistik untuk menghitung representasi vektor dari kata-kata. Metode Glove menghitung representasi vektor dari kata-kata berdasarkan frekuensi kemunculan kata-kata dalam dokumen teks.

Word2Vec adalah teknik yang menggabungkan teknik skip-gram dan CBOW. Teknik Word2Vec dapat menghasilkan representasi vektor dari kata-kata yang lebih akurat daripada teknik skip-gram atau CBOW saja. Word embeddings dan word vectors adalah representasi vektor dari kata-kata dalam bahasa alami. Word embeddings dan word vectors dapat digunakan untuk berbagai tugas pemrosesan bahasa alami. Ada berbagai teknik yang dapat digunakan untuk melakukan word embeddings dan word vectors. Beberapa teknik yang umum digunakan adalah skip-gram, CBOW, Glove, dan Word2Vec.



## C. Algoritma dan Model dalam NLP

Algoritma dan model dalam NLP dapat diklasifikasikan menjadi dua kategori utama, yaitu: Algoritma berbasis aturan dan Algoritma berbasis pembelajaran mesin.

Algoritma berbasis aturan

Algoritma berbasis aturan menggunakan serangkaian aturan untuk menyelesaikan tugas NLP. Aturan-aturan ini biasanya dibuat oleh ahli bahasa berdasarkan pengetahuan mereka tentang bahasa. Beberapa contoh algoritma berbasis aturan dalam NLP adalah:

1. Algoritma part-of-speech tagging
2. Algoritma named entity recognition
3. Algoritma sentiment analysis
4. Algoritma berbasis pembelajaran mesin

Algoritma berbasis pembelajaran mesin menggunakan model pembelajaran mesin untuk menyelesaikan tugas NLP. Model pembelajaran mesin dilatih pada data pelatihan yang berisi contoh input dan output. Algoritma machine translation, Algoritma summarization, Algoritma question answering. Selain itu, ada beberapa algoritma dan model NLP yang khusus untuk tugas tertentu, seperti:

Algoritma speech recognition, Algoritma natural language generation, Algoritma natural language understanding. Berikut adalah penjelasan lebih rinci tentang beberapa algoritma dan model NLP yang umum digunakan: Algoritma part-of-speech tagging, Algoritma part-of-speech tagging (POS tagging) adalah algoritma yang digunakan untuk menentukan kelas kata (part-of-speech) dari setiap kata dalam kalimat. Kelas kata menentukan fungsi kata dalam kalimat. Beberapa contoh kelas kata adalah:

Nama (noun)

Kata kerja (verb)

Kata sifat (adjective)

Kata keterangan (adverb)

Kata ganti (pronoun)

Algoritma POS tagging dapat digunakan untuk berbagai tugas NLP, seperti:

Pencarian informasi, Pemahaman bahasa alami, Pertanyaan jawab, Algoritma named entity recognition. Algoritma named entity recognition (NER) adalah algoritma yang digunakan untuk mengidentifikasi dan mengklasifikasikan entitas bernama dalam teks. Entitas bernama adalah frasa benda (noun phrase) yang memiliki tipe spesifik, seperti nama orang, nama organisasi, lokasi, tanggal, dan waktu. Beberapa contoh entitas bernama adalah:

John Doe (nama orang)  
Microsoft (nama organisasi)  
Jakarta (lokasi)  
2023-12-13 (tanggal)  
10:00 (waktu)

Algoritma NER dapat digunakan untuk berbagai tugas NLP, seperti: Pencarian informasi, Pemahaman bahasa alami, Aplikasi bisnis. Algoritma machine translation: Algoritma machine translation (MT) adalah algoritma yang digunakan untuk menerjemahkan teks dari satu bahasa ke bahasa lain. Algoritma MT menggunakan model pembelajaran mesin untuk mempelajari hubungan antara kata dan frasa dalam dua bahasa. Algoritma MT dapat digunakan untuk berbagai keperluan, seperti: Penerjemahan dokumen, Penerjemahan percakapan, Penerjemahan situs web

Algoritma summarization: Algoritma summarization adalah algoritma yang digunakan untuk membuat ringkasan dari teks. Algoritma summarization menggunakan model pembelajaran mesin untuk mengidentifikasi informasi yang penting dalam teks dan untuk menyusun informasi tersebut menjadi ringkasan yang ringkas dan informatif. Algoritma summarization dapat digunakan untuk berbagai keperluan, seperti: Pencarian informasi, Pemahaman bahasa alami, Layanan pelanggan

Algoritma question answering: Algoritma question answering (QA) adalah algoritma yang digunakan untuk menjawab pertanyaan tentang teks. Algoritma QA menggunakan model pembelajaran mesin untuk memahami pertanyaan dan untuk mencari jawaban yang relevan dalam teks. Algoritma QA dapat digunakan untuk berbagai keperluan, seperti: Pendidikan, Pencarian informasi, Layanan pelanggan. Terdapat banyak algoritma dan model yang digunakan dalam Natural Language Processing (NLP) untuk berbagai tugas analisis bahasa. Berikut adalah beberapa di antaranya:

**Naive Bayes:** Digunakan dalam klasifikasi teks atau analisis sentimen berdasarkan probabilitas dan teorema Bayes. **Support Vector Machines (SVM):** Model pembelajaran yang digunakan untuk klasifikasi teks dan pembedaan antara kelas-kelas yang berbeda dalam data teks. **Hidden Markov Models (HMM):** Digunakan dalam pemodelan urutan kata-kata atau kata-kata tersembunyi dalam konteks seperti tugas NER atau POS tagging. **Conditional Random Fields (CRF):** Algoritma yang digunakan dalam tugas NER, POS tagging, atau labeling urutan berdasarkan ketergantungan kondisional antara elemen dalam urutan.

Model dalam NLP:

**Transformer:** Model yang revolusioner dalam NLP, seperti BERT, GPT (Generative Pre-trained Transformer), dan lainnya. Mereka menggunakan self-attention mechanism untuk memproses teks dan menghasilkan representasi yang lebih baik. **Recurrent Neural Networks (RNN):** Model yang memproses data urutan seperti teks dan menghasilkan output berdasarkan pemahaman konteks sebelumnya. Variannya seperti LSTM dan GRU sering digunakan dalam NLP. **Word2Vec:** Model yang menghasilkan representasi vektor kata-kata dalam ruang vektor berdasarkan kemunculan kata-kata dalam konteks tertentu. **GloVe (Global Vectors for Word Representation):** Model yang menghasilkan representasi vektor kata-kata berdasarkan matriks co-occurrence kata-kata dalam teks. **BERT (Bidirectional**

**Encoder Representations from Transformers):** Model yang menggunakan transformer dan pelatihan unsupervised learning pada teks besar untuk memahami konteks dan menyediakan representasi yang lebih baik untuk kata-kata. **CNN (Convolutional Neural Networks):** Meskipun awalnya digunakan dalam pengolahan citra, varian CNN juga digunakan dalam NLP untuk tugas seperti klasifikasi teks dan analisis sentiment (Tunstall et al., 2022).

Pendekatan Hybrid:

Beberapa model NLP menggabungkan berbagai elemen dari model-model yang berbeda atau menggunakan strategi ensemble untuk meningkatkan kinerja dan kemampuan dalam pemrosesan bahasa alami. Setiap algoritma dan model memiliki keunggulan dan kelemahan tertentu, dan pemilihan yang tepat tergantung pada tugas spesifik dalam NLP yang akan dijalankan serta ketersediaan data yang tersedia. Kombinasi model dan pendekatan tertentu sering kali memberikan kinerja yang lebih baik dalam berbagai konteks pengolahan bahasa alami.

## 1. Naive Bayes

Naive Bayes merupakan algoritma klasifikasi yang menggunakan Teorema Bayes untuk menentukan probabilitas kelas suatu sampel data berdasarkan fitur-fitur yang diamati. Meskipun sering digunakan dalam konteks klasifikasi teks atau analisis sentimen, Naive Bayes juga diterapkan dalam berbagai masalah klasifikasi di bidang Machine Learning.

Konsep Naive Bayes: Teorema Bayes: Mendasarkan diri pada Teorema Bayes yang menyatakan hubungan antara probabilitas posterior (kemungkinan kejadian setelah melihat data baru) dengan probabilitas prior (kejadian sebelum melihat data baru) dan likelihood (kemungkinan data yang diamati jika kelas tertentu).

Klasifikasi dengan Probabilitas: Naive Bayes memprediksi kelas suatu sampel data berdasarkan perhitungan probabilitas kelas tersebut dengan mengasumsikan independensi antara fitur-fitur yang diamati (yang sering disebut "naive" karena asumsi ini).

Proses Naive Bayes dalam Klasifikasi Teks atau Analisis Sentimen: Pra-Pemrosesan: Tahap ini melibatkan tokenisasi teks, penghapusan stopwords, dan representasi fitur-fitur teks dalam bentuk vektor (misalnya, TF-IDF atau Bag-of-Words). Penghitungan Probabilitas: Menghitung probabilitas masing-masing kelas (positif, negatif, atau kelas lainnya) berdasarkan kemunculan kata-kata atau fitur-fitur yang diamati dalam data pelatihan. Teorema Bayes: Memanfaatkan Teorema Bayes untuk menghitung probabilitas posterior dari kelas berdasarkan fitur-fitur yang diamati dalam data uji. Klasifikasi: Memilih kelas dengan probabilitas posterior tertinggi sebagai prediksi kelas untuk sampel data yang diamati.

Keunggulan dan Keterbatasan Naive Bayes: Keunggulan: Cepat dalam pembelajaran dan prediksi, bahkan dengan dataset yang besar. Efektif dalam klasifikasi teks dengan fitur yang besar. Keterbatasan: Asumsi naif tentang independensi fitur bisa tidak realistis dalam konteks nyata. Kinerjanya dapat terpengaruh jika ada ketergantungan antara fitur-fitur yang diamati.

Aplikasi Naive Bayes dalam NLP: Analisis Sentimen: Klasifikasi teks berdasarkan sentimen (positif, negatif, atau netral). Klasifikasi Teks: Pengelompokan teks ke dalam kategori tertentu seperti klasifikasi berita, spam detection, dan lainnya. Naive Bayes, meskipun memiliki asumsi naif tentang independensi fitur, tetap menjadi salah satu algoritma klasifikasi yang cukup populer dalam NLP karena kemampuannya dalam menangani klasifikasi teks dengan baik, terutama ketika dataset besar dan fitur-fitur teks yang kompleks terlibat.

## 2. Support Vector Machines (SVM)

Support Vector Machines (SVM) merupakan salah satu model pembelajaran yang umum digunakan dalam klasifikasi, termasuk dalam konteks analisis teks dalam Natural Language Processing (NLP). SVM digunakan untuk memisahkan dan mengklasifikasikan data ke dalam kelas-kelas yang berbeda dengan mencari hyperplane terbaik yang memisahkan antara kelas-kelas tersebut di dalam ruang fitur.

Konsep SVM:

**Hyperplane:** SVM mencari hyperplane (bidang dalam kasus dua dimensi) yang memisahkan data ke dalam kelas-kelas yang berbeda. Pemisahan ini dilakukan sedemikian rupa sehingga jarak (marginal) antara hyperplane dan titik-titik data (yang disebut support vectors) dari kedua kelas adalah maksimum.

**Pemisahan Non-linear:** SVM dapat mengatasi masalah pemisahan yang tidak linier dengan menggunakan fungsi kernel yang dapat mentransformasi data ke dalam dimensi yang lebih tinggi, memungkinkan pembuatan hyperplane yang lebih kompleks untuk pemisahan yang lebih baik.

Proses SVM dalam Klasifikasi Teks:

1. **Pra-Pemrosesan:** Mirip dengan langkah pra-pemrosesan untuk model lainnya dalam NLP, seperti tokenisasi, penghapusan stopwords, dan pembuatan vektor fitur seperti TF-IDF atau Bag-of-Words.
2. **Pemilihan Hyperplane:** SVM akan mencari hyperplane terbaik yang memisahkan antara kelas-kelas dalam ruang fitur berdasarkan vektor fitur teks.
3. **Penentuan Margin Terbesar:** SVM berusaha menemukan hyperplane yang memiliki margin terbesar, yaitu jarak terbesar antara support vectors dan hyperplane tersebut.
4. **Penentuan Kelas:** Setelah mendapatkan hyperplane terbaik, SVM dapat mengklasifikasikan data baru ke dalam kelas yang sesuai berdasarkan posisi relatifnya terhadap hyperplane.

Keunggulan dan Keterbatasan SVM: **Keunggulan:** Efektif dalam ruang fitur berdimensi tinggi, mampu menangani dataset yang kompleks, dan dapat berkinerja baik bahkan dengan jumlah fitur yang lebih besar. **Keterbatasan:** Rentan terhadap overfitting jika hyperparameter tidak disesuaikan dengan baik. Membutuhkan pemrosesan yang memakan waktu untuk pemilihan parameter yang tepat.

Aplikasi SVM dalam NLP yaitu **Klasifikasi Teks:** Klasifikasi dokumen, analisis sentimen, deteksi spam, kategorisasi teks, dan lainnya. **Pengelompokan Teks:** SVM digunakan untuk mengelompokkan teks-teks dengan karakteristik yang serupa. SVM merupakan algoritma yang kuat dalam klasifikasi teks dan telah banyak digunakan dalam berbagai tugas NLP karena kemampuannya dalam menangani pemisahan antar kelas dalam ruang fitur dengan baik, bahkan pada dataset dengan dimensi yang tinggi atau kompleks. SVM (Support Vector Machine) adalah algoritma klasifikasi yang bertujuan untuk menemukan hyperplane terbaik yang memisahkan antara kelas-kelas dalam data dengan margin terbesar. Dalam kasus klasifikasi biner (dua kelas), rumus untuk SVM dengan hyperplane linier dapat dijelaskan sebagai berikut:

Hyperplane Linier:

Misalkan kita memiliki data pelatihan dengan vektor fitur  $x_i$  yang terdiri dari  $n$  fitur, dan label kelasnya adalah  $y_i$  (dengan  $y_i=1$  atau  $y_i=-1$ ). Fungsi hipotesis untuk SVM dapat ditulis sebagai:

$$f(x)=w^T x+b$$

di mana:

$w$  adalah vektor bobot normal ke hyperplane,

$x$  adalah vektor fitur,

$b$  adalah bias.

Persamaan untuk Hyperplane:

Hyperplane dipilih untuk memiliki margin terbesar antara kelas yang dipisahkan. Jarak dari titik data ke hyperplane adalah  $\frac{1}{||w||}$ .

Maksimalkan margin dengan meminimalkan  $\|w\|$  (norma Euclidean dari  $w$ ) yang setara dengan meminimalkan  $1/2\|w\|^2$ , dengan mempertimbangkan pembatasan:  $y_i(w^T x_i + b) \geq 1$  untuk semua titik data yang merupakan support vectors. Fungsi Objektif (Objective Function): Objektif utama dari SVM adalah meminimalkan fungsi objektif berikut:

$$\min_{w,b} 1/2\|w\|^2$$

dengan pembatasan:

$$y_i(w^T x_i + b) \geq 1$$

### **Fungsi Objektif (Objective Function):**

Objektif utama dari SVM adalah meminimalkan fungsi objektif berikut:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

dengan pembatasan:

$$y_i(w^T x_i + b) \geq 1$$

### **Metode Lagrange:**

Dalam praktiknya, konsep Lagrange multipliers digunakan untuk menyelesaikan masalah optimasi SVM dengan pembatasan ini.

Menggunakan metode Lagrange, kita mendapatkan fungsi Lagrange yang akan dioptimalkan untuk mencari nilai minimum:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1]$$

di mana  $\alpha$  adalah Lagrange multipliers.

Solusi SVM ditemukan dengan mencari nilai minimum dari fungsi Lagrange tersebut. Solusi ini menghasilkan hyperplane terbaik yang memisahkan antara kelas dengan margin terbesar di antara support vectors.



Rumus-rumus ini mencerminkan dasar dari bagaimana SVM bekerja dalam menemukan hyperplane terbaik untuk pemisahan kelas dalam data. Dalam prakteknya, untuk kasus-kasus yang lebih kompleks, seperti SVM non-linier, digunakan kernel untuk memungkinkan pemisahan yang lebih kompleks dalam ruang fitur yang lebih tinggi.

### **3. Recurrent Neural Networks (RNN)**

Jaringan Saraf Rekurensial (RNN) adalah jenis arsitektur jaringan saraf yang dirancang khusus untuk memproses data urutan atau data yang terstruktur secara sekuensial, seperti teks, audio, atau data deret waktu. Konsep dasar RNN: Memori Jangka Pendek: RNN memiliki kemampuan untuk menyimpan informasi dalam keadaan internal atau "memori jangka pendek". Ini memungkinkannya untuk mengingat informasi sebelumnya saat memproses elemen berikutnya dalam urutan. Keterkaitan Antar Elemen: Setiap elemen dalam urutan diolah secara berurutan, dan informasi dari elemen sebelumnya digunakan untuk memproses elemen berikutnya.

Arsitektur RNN: RNN memiliki struktur yang mengizinkan informasi untuk mengalir mundur melalui jaringan, memungkinkan koneksi siklus atau rekurensi. Setiap langkah waktu (time step) dalam urutan diproses oleh lapisan yang sama dari jaringan. Jenis-jenis RNN: One-to-One: Input tunggal menghasilkan output tunggal, seperti dalam jaringan saraf biasa (feedforward). One-to-Many: Satu input menghasilkan serangkaian output, seperti menghasilkan kalimat dari gambar tunggal. Many-to-One: Serangkaian input menghasilkan output tunggal, seperti klasifikasi teks dari urutan kata-kata. Many-to-Many: Urutan input dihubungkan dengan urutan output, seperti terjemahan mesin atau POS tagging.

**Keunggulan RNN:**

**Penanganan Data Urutan:** Cocok untuk data yang memiliki struktur sekuensial, seperti teks, audio, atau data deret waktu.  
**Memori Jangka Pendek:** Kemampuan untuk "mengingat"

informasi dari langkah-langkah sebelumnya dalam urutan. Keterbatasan RNN: Masalah Pelatihan: Rentan terhadap masalah menghilangnya atau meledaknya gradien, yang mempengaruhi kemampuannya untuk memahami hubungan jarak jauh dalam urutan panjang. Komputasi yang Lambat: Keterkaitan antar elemen membuat proses komputasi RNN cenderung lambat dalam pengolahan data yang panjang.

Penggunaan RNN dalam NLP: RNN sering digunakan dalam NLP untuk tugas-tugas seperti: Analisis Sentimen: Mengklasifikasikan sentimen dalam teks. Penerjemahan Mesin: Menerjemahkan teks dari satu bahasa ke bahasa lain. Generasi Teks: Menghasilkan teks yang baru, seperti pembuatan cerita atau artikel. RNN adalah alat yang berguna dalam NLP karena kemampuannya untuk memproses data teks secara sekuensial dan mempertahankan memori jangka pendek, yang memungkinkannya untuk menghadapi tugas-tugas kompleks dalam analisis bahasa alami. Tahapan dalam Recurrent Neural Networks (RNN) melibatkan proses yang berurutan dari pengolahan data urutan di setiap langkah waktu (time step). Berikut adalah tahapan-tahapan utama dalam RNN:

#### 1. Input Data:

**Data Urutan:** Seperti teks, data deret waktu, atau data yang memiliki struktur sekuensial. Misalnya, urutan kata-kata dalam kalimat.

#### 2. Representasi Data:

**Tokenisasi:** Pemecahan data urutan menjadi unit-unit terpisah (token), misalnya, pemecahan kalimat menjadi kata-kata atau karakter. **Pembuatan Vektor Fitur:** Representasi data dalam bentuk vektor fitur seperti Bag-of-Words atau word embeddings (misalnya, Word2Vec, GloVe).

#### 3. Pembuatan Arsitektur RNN:

**Inisialisasi:** Pembuatan model RNN dengan menentukan jumlah neuron, lapisan, dan jenis RNN yang akan digunakan (misalnya, SimpleRNN, LSTM, atau GRU). **Pengaturan Hyperparameter:**

Penyesuaian hyperparameter seperti jumlah neuron, jumlah lapisan, learning rate, dan jenis fungsi aktivasi.

#### 4. Pelatihan (Training):

**Pengolahan Data Langkah per Langkah:** Pengolahan data langkah per langkah (time step) dalam urutan. **Pembelajaran Bobot:** Pembelajaran parameter bobot dalam jaringan dengan menggunakan algoritma seperti backpropagation melalui waktu (Backpropagation Through Time - BPTT).

#### 5. Validasi dan Evaluasi:

**Validasi Model:** Menggunakan data validasi untuk mengevaluasi performa model dan menyesuaikan hyperparameter jika diperlukan. **Evaluasi Performa:** Menggunakan metrik evaluasi seperti akurasi, F1-score, atau perplexity untuk mengevaluasi kinerja model.

#### 6. Prediksi dan Penggunaan Model:

**Prediksi:** Menggunakan model yang telah dilatih untuk membuat prediksi pada data baru atau untuk tugas yang ditentukan, seperti klasifikasi, generasi teks, atau penerjemahan.

#### 7. Penyesuaian dan Peningkatan Model:

**Optimisasi:** Peningkatan model dengan penyesuaian hyperparameter, pemilihan arsitektur yang lebih kompleks, atau teknik regularisasi.

Catatan: **Backpropagation Through Time (BPTT):** Penting dalam RNN karena memungkinkan propagasi gradien dari langkah waktu ke langkah waktu, membantu dalam pembelajaran parameter bobot. **Overfitting:** RNN rentan terhadap overfitting pada data urutan yang panjang. Oleh karena itu, pemilihan model yang tepat dan teknik regularisasi seperti dropout atau batch normalization dapat membantu mengurangi masalah ini. Tahapan-tahapan ini membentuk proses umum dalam penggunaan dan pengolahan data menggunakan Recurrent Neural Networks dalam berbagai tugas dalam pemrosesan bahasa alami dan pengenalan pola dalam data sekuensial.

#### 4. Convolutional Neural Networks (CNN)

Convolutional Neural Networks (CNN) adalah jenis arsitektur jaringan saraf yang dirancang khusus untuk pemrosesan dan pengenalan pola dalam data grid, seperti gambar. CNN sangat efektif dalam mengekstraksi fitur-fitur spasial dari data dan telah menjadi alat yang kuat dalam pengolahan citra dan pengenalan pola.

**Konsep dasar CNN:** Konvolusi: CNN menggunakan operasi konvolusi untuk mengekstraksi fitur-fitur dari gambar. Ini melibatkan pergerakan filter (kernel) ke seluruh gambar untuk mendeteksi pola-pola visual seperti tepi, sudut, atau tekstur.

Pooling: Operasi pooling seperti max pooling digunakan untuk mereduksi dimensi dari fitur yang diekstraksi, mempertahankan informasi penting dan mengurangi kompleksitas model.

**Struktur CNN:**

Convolutional Layer: Lapisan konvolusi terdiri dari filter yang memindai gambar untuk mengekstraksi fitur-fitur. Pooling Layer: Lapisan pooling yang mengurangi dimensi spasial dari fitur yang diekstraksi oleh lapisan konvolusi. Fully Connected Layer: Lapisan-lapisan terhubung sepenuhnya yang menggabungkan fitur-fitur yang diekstraksi untuk klasifikasi akhir.

**Keunggulan CNN:**

Ekstraksi Fitur Otomatis: Kemampuan untuk secara otomatis mengekstraksi fitur-fitur hierarkis dari data gambar. Invariansi Spatial: Invariansi terhadap pergeseran dan transformasi kecil dalam gambar. Aplikasi CNN dalam Computer Vision: Klasifikasi Gambar: Mengenali objek atau kelas dari gambar. Deteksi Objek: Menemukan dan menandai lokasi objek dalam gambar. Segmentasi Gambar: Memisahkan objek dari latar belakang.

**Perkembangan Terkini:** Transfer Learning: Pemanfaatan model-model yang sudah dilatih sebelumnya untuk tugas-tugas spesifik dalam gambar atau domain lain. Architectural Advancements: Pengembangan arsitektur yang lebih kompleks seperti ResNet,

Inception, atau EfficientNet untuk kinerja yang lebih baik dalam pengenalan gambar.

CNN telah menjadi landasan dalam bidang Computer Vision dan telah membawa kemajuan signifikan dalam berbagai aplikasi, mulai dari pengenalan wajah hingga mobil otonom. Kemampuannya dalam mengekstraksi fitur dari gambar dan merangkai informasi spasial membuatnya sangat efektif dalam memahami data visual. Convolutional Neural Networks (CNN) pada awalnya dikembangkan untuk memproses data grid seperti gambar. Namun, beberapa penelitian terbaru telah mencoba menerapkan konsep dasar CNN dalam pemrosesan data teks dengan perubahan dalam representasi dan pemrosesan.

Penerapan CNN untuk data teks melibatkan transformasi data teks menjadi matriks atau tensor yang dapat diproses oleh CNN. Berikut adalah cara umum dalam menggunakan CNN untuk pemrosesan data teks:

#### 1. Representasi Data Teks:

Word Embeddings: Mengubah kata-kata dalam teks menjadi vektor numerik, seperti Word2Vec, GloVe, atau FastText. One-Hot Encoding: Representasi biner dari kata-kata dalam bentuk matriks yang besar (sering kali digunakan dalam kasus yang lebih sederhana).

#### 2. Convolutional Layer:

Convolution: Konvolusi diterapkan pada representasi vektor kata atau matriks berukuran kelompok kata-kata yang disebut filter (kernel) untuk mengekstraksi fitur-fitur dari urutan kata-kata. Misalnya, filter 1 akan mendeteksi fitur-fitur seperti kata-kata yang berdekatan, filter 2 mungkin mendeteksi pola khusus lainnya.

#### 3. Max Pooling Layer:

Pooling: Operasi pooling (misalnya, max pooling) digunakan untuk mereduksi dimensi vektor hasil konvolusi, mempertahankan fitur-fitur penting sambil mengurangi kompleksitas.

#### 4. Fully Connected Layer:

Layer Terhubung Penuh: Fitur-fitur yang diekstraksi oleh CNN kemudian digabungkan ke dalam lapisan-lapisan terhubung penuh untuk klasifikasi akhir atau tugas-tugas lainnya seperti analisis sentimen atau klasifikasi teks.

Keunggulan dan Penggunaan dalam NLP: Ekstraksi Fitur: CNN dapat mempelajari representasi hierarkis dari teks, mengenali pola seperti kata-kata yang berdekatan, frasa, atau makna tertentu dalam teks. Analisis Sentimen: Penerapan CNN pada data teks untuk analisis sentimen atau klasifikasi teks. Pemrosesan Bahasa: Meskipun RNN dan Transformer lebih umum dalam NLP, CNN telah digunakan dalam tugas-tugas seperti pemrosesan bahasa dan generasi teks. Catatan Penting: Ukuran Jendela (Window Size): Penting untuk memilih ukuran jendela yang sesuai untuk filter dalam CNN agar dapat menangkap pola yang relevan dalam teks.

Penggunaan Bersama dengan Model Lain: Pada beberapa kasus, CNN digunakan bersama dengan arsitektur lain seperti LSTM atau dilakukan fine-tuning pada model pre-trained untuk kinerja yang lebih baik. Meskipun awalnya CNN dirancang untuk data grid seperti gambar, adaptasi terbaru telah memperluas penggunaannya dalam pemrosesan data teks dengan mengubah representasi teks menjadi format yang dapat diproses oleh CNN. Meskipun RNN dan Transformer masih dominan dalam NLP, CNN tetap menjadi area penelitian yang menarik dalam pemrosesan bahasa alami. Tahapan Convolutional Neural Networks (CNN) untuk pengolahan data teks melibatkan serangkaian langkah yang mirip dengan penggunaan CNN pada data gambar, namun dengan representasi data teks yang berbeda. Berikut tahapan-tahapan utama:

##### 1. Preprocessing Data Teks:

Tokenisasi: Pemecahan teks menjadi token, seperti kata-kata atau karakter. Pembuatan Sequence: Membentuk urutan dari token-token yang terbagi dalam suatu urutan (kalimat atau teks).

## 2. Representasi Data Teks:

**Word Embeddings:** Mengonversi token-token teks menjadi vektor numerik menggunakan Word Embeddings seperti Word2Vec, GloVe, atau FastText. Ini mengubah kata-kata menjadi representasi vektor dalam ruang dimensi yang lebih kecil.

**One-Hot Encoding:** Representasi biner dari kata-kata dalam bentuk matriks yang besar. Sering kali digunakan dalam eksperimen awal atau dalam kasus dengan dataset yang lebih sederhana.

## 3. Convolutional Layer:

**Convolution:** Lapisan konvolusi diterapkan pada representasi vektor kata atau matriks berukuran kelompok kata-kata yang disebut filter (kernel) untuk mengekstraksi fitur-fitur dari urutan kata-kata. Filter akan "memindai" sekuens dan menemukan pola yang relevan.

## 4. Pooling Layer:

**Pooling:** Operasi pooling (misalnya, max pooling) digunakan untuk mereduksi dimensi vektor hasil konvolusi, mempertahankan fitur-fitur penting sambil mengurangi kompleksitas dan ukuran data.

## 5. Fully Connected Layer:

**Layer Terhubung Penuh:** Fitur-fitur yang diekstraksi oleh CNN kemudian digabungkan ke dalam lapisan-lapisan terhubung penuh untuk klasifikasi akhir atau tugas lain seperti analisis sentimen, klasifikasi teks, atau generasi teks.

## 6. Output Layer:

**Layer Output:** Lapisan terakhir yang menghasilkan output yang cocok dengan tugas spesifik yang dijalankan, seperti klasifikasi kategori teks atau nilai sentimen.

**Catatan Penting: Hyperparameter Tuning:** Pemilihan filter size, jumlah filter, tingkat dropout, dan learning rate adalah kunci untuk meningkatkan performa CNN pada data teks.

Penyesuaian dengan Arsitektur: Beberapa penelitian menggunakan CNN bersama dengan arsitektur lain seperti LSTM atau menggunakan fine-tuning pada model pre-trained untuk hasil yang lebih baik. Dimensi Data: Representasi vektor kata-kata dari Word Embeddings atau matriks one-hot encoding mempengaruhi dimensi data masukan dan proses konvolusi. Penerapan CNN pada data teks mengharuskan perubahan representasi teks menjadi format yang dapat diproses oleh CNN, dan meskipun CNN lebih sering digunakan dalam pengolahan gambar, terdapat penelitian dan aplikasi yang menarik dalam penggunaannya pada data teks.

#### 5. Transformer-based models (seperti BERT, GPT, dll.)

Transformer-based models adalah model pembelajaran mesin yang menggunakan arsitektur Transformer. Transformer adalah arsitektur jaringan saraf tiruan yang dirancang untuk menangani data berurutan, seperti teks, audio, dan video. Transformer memiliki kemampuan untuk memahami hubungan antara input dan output, sehingga membuatnya lebih cocok untuk tugas-tugas yang membutuhkan pemahaman kontekstual.

Model berbasis Transformer merupakan pendekatan revolusioner dalam pemrosesan bahasa alami yang menghilangkan ketergantungan pada urutan (sequence dependency) dan menggabungkan mekanisme self-attention untuk memahami konteks dari kata-kata atau token dalam teks. Ini memungkinkan model untuk memperoleh pemahaman yang lebih baik tentang hubungan antara kata-kata dalam teks. Konsep dasar Transformer: Self-Attention Mechanism: Transformer menggunakan self-attention untuk menimbang hubungan antara semua token dalam kalimat secara sekaligus. Ini memungkinkan model untuk memberikan bobot yang tepat untuk setiap token berdasarkan hubungannya dengan token lain dalam kalimat tersebut. Multi-Head Attention: Transformer memiliki beberapa head attention yang beroperasi secara independen, memungkinkan model untuk mempelajari hubungan yang lebih kompleks di berbagai sudut pandang.



Encoder-Decoder Architecture: Transformer umumnya terdiri dari blok encoder untuk memproses input dan blok decoder untuk menghasilkan output dalam tugas seperti penerjemahan bahasa.

Komponen-komponen Transformer: Positional Encoding: Karena Transformer tidak mempertahankan urutan token seperti RNN, positional encoding diperkenalkan untuk memperkenalkan informasi urutan ke dalam representasi vektor token. Encoder Layers: Setiap layer dalam blok encoder memiliki modul self-attention dan modul feed-forward neural network. Decoder Layers: Mirip dengan encoder, blok decoder memiliki self-attention, ditambah attention terhadap output dari encoder (untuk tugas penerjemahan).

Keunggulan Transformer: Paralelisme yang Lebih Baik: Transformer memungkinkan perhitungan paralel yang lebih efisien dibandingkan dengan RNN atau LSTM, mempercepat pelatihan model. Kemampuan untuk Memahami Konteks yang Lebih Luas: Mekanisme self-attention memungkinkan model untuk memahami hubungan antara kata-kata yang jauh dalam teks.

Aplikasi Transformer-based Models: Penerjemahan Bahasa: Model Transformer seperti GPT (Generative Pre-trained Transformer) dan BERT telah digunakan untuk penerjemahan bahasa dan pemahaman teks yang lebih baik. Generasi Teks: Transformer juga digunakan untuk menghasilkan teks yang lebih lancar dan realistis dalam tugas generasi teks. Analisis Sentimen: Penggunaan Transformer dalam tugas analisis sentimen telah meningkatkan akurasi dan pemahaman konteks sentimen dalam teks.

Contoh Model Transformer: BERT (Bidirectional Encoder Representations from Transformers): Model yang dilatih secara pre-trained untuk pemahaman konteks bahasa. GPT (Generative Pre-Trained Transformer): Model yang berfokus pada generasi teks yang lebih baik berdasarkan pemahaman konteks bahasa. Transformer-based models telah mengubah lanskap dalam

pemrosesan bahasa alami dengan meningkatkan kemampuan memahami konteks dan urutan kata-kata dalam teks, dan mereka terus menjadi fokus penelitian dan pengembangan dalam NLP.

**Arsitektur Transformer:** Arsitektur Transformer terdiri dari dua bagian utama, yaitu encoder dan decoder. Encoder bertanggung jawab untuk menganalisis input, sedangkan decoder bertanggung jawab untuk menghasilkan output. Encoder terdiri dari beberapa layer self-attention. Self-attention adalah teknik yang digunakan untuk menghitung hubungan antara setiap input. Self-attention memungkinkan Transformer untuk memahami hubungan antara input yang berdekatan dan input yang jauh.

Decoder juga terdiri dari beberapa layer self-attention. Selain itu, decoder juga memiliki layer feedforward. Feedforward adalah teknik yang digunakan untuk mengubah representasi vektor dari input.

**Keunggulan Transformer-based models.** Transformer-based models memiliki beberapa keunggulan dibandingkan dengan model-model sebelumnya, seperti RNN dan CNN. Keunggulan-keunggulan tersebut antara lain:

1. Kemampuan untuk memahami hubungan: Transformer dapat memahami hubungan antara input dan output, sehingga membuatnya lebih cocok untuk tugas-tugas yang membutuhkan pemahaman kontekstual.
2. Efisiensi: Transformer lebih efisien daripada RNN dan CNN, terutama untuk tugas-tugas yang membutuhkan pemahaman jangka panjang.
3. Keakuratan: Transformer dapat menghasilkan akurasi yang lebih tinggi daripada RNN dan CNN untuk berbagai tugas, seperti terjemahan mesin, pengenalan bahasa alami, dan klasifikasi teks.

4. Aplikasi Transformer-based models Transformer-based models telah digunakan untuk berbagai tugas pemrosesan data berurutan, seperti:
5. Terjemahan mesin: Transformer telah menjadi standar de facto untuk terjemahan mesin. Transformer dapat menghasilkan terjemahan yang lebih akurat dan alami daripada model-model sebelumnya.
6. Pemahaman bahasa alami: Transformer telah digunakan untuk berbagai tugas pemrosesan bahasa alami, seperti pengenalan entitas, klasifikasi teks, dan sentiment analysis.
7. Pengenalan suara: Transformer telah digunakan untuk meningkatkan akurasi pengenalan suara.
8. Komposisi musik: Transformer telah digunakan untuk menghasilkan musik yang mirip dengan musik yang sudah ada.

Transformer-based models adalah model pembelajaran mesin yang kuat dan serbaguna yang dapat digunakan untuk berbagai tugas pemrosesan data berurutan. Transformer telah menjadi standar de facto untuk beberapa tugas, seperti terjemahan mesin.

# BAB III

## PENERAPAN NLP

### A. NLP dalam Industri

#### 1. Customer Service

Penerapan Natural Language Processing (NLP) dalam layanan pelanggan (Customer Service) telah menjadi kunci dalam meningkatkan pengalaman pelanggan, efisiensi operasional, dan pemahaman yang lebih baik terhadap kebutuhan pelanggan. Berikut adalah rincian tentang bagaimana NLP digunakan dalam layanan pelanggan:

1. Chatbot dan Asisten Virtual: **Otomatisasi Respon:** Chatbot menggunakan NLP untuk memahami pertanyaan atau masalah pelanggan dan memberikan respons yang relevan secara otomatis. Mereka dapat membantu dalam menjawab pertanyaan umum, memberikan informasi produk, atau menyelesaikan masalah tertentu. **Pemahaman Bahasa Alami:** Melalui NLP, chatbot dapat memahami pertanyaan dalam bahasa alami, bahkan dengan variasi atau frasa yang berbeda.
2. Analisis Sentimen: **Pemantauan Sentimen:** NLP digunakan untuk menganalisis sentimen dari ulasan atau feedback pelanggan di media sosial, forum, atau platform lainnya. Ini membantu perusahaan memahami perasaan pelanggan terhadap produk atau layanan mereka. **Deteksi Masalah:** Dengan menganalisis sentimen, perusahaan dapat mendeteksi masalah atau keluhan yang muncul dari pelanggan secara cepat dan meresponnya dengan tepat waktu.
3. Analisis Percakapan Pelanggan: **Pemrosesan Transkripsi:** NLP digunakan untuk menganalisis percakapan telepon, chat, atau email dengan pelanggan untuk mengekstrak informasi penting seperti masalah umum, kebutuhan, atau keluhan.

- Peningkatan Layanan:** Analisis NLP pada percakapan pelanggan dapat membantu dalam mengidentifikasi area di mana layanan dapat ditingkatkan atau masalah yang perlu diselesaikan.
4. **Klasifikasi dan Pemrosesan Permintaan Pelanggan:**  
**Klasifikasi Permintaan:** NLP digunakan untuk mengklasifikasi permintaan pelanggan ke dalam kategori yang tepat. Ini membantu dalam menentukan prioritas dan menanggapi dengan lebih cepat.  
**Pemrosesan Otomatis:** Dengan pemahaman NLP terhadap permintaan pelanggan, beberapa tugas dapat diproses secara otomatis, seperti pembuatan tiket layanan atau pengiriman pesan balasan awal.
  5. **Personalisasi Layanan:**  
**Analisis Riwayat:** NLP membantu menganalisis riwayat interaksi pelanggan untuk memahami preferensi, kebutuhan, dan pola perilaku. Ini memungkinkan perusahaan untuk memberikan layanan yang lebih personal dan relevan.
  6. **Implementasi NLP di Platform Layanan Pelanggan:**  
**Integrasi dalam CRM:** Integrasi NLP dalam perangkat lunak manajemen hubungan pelanggan (CRM) membantu dalam pemrosesan dan pengelolaan data pelanggan untuk memberikan layanan yang lebih baik.  
**Penggunaan API:** Penggunaan API (Application Programming Interface) NLP dari penyedia layanan dapat memperkaya fungsionalitas platform layanan pelanggan dengan kemampuan bahasa alami. Penerapan NLP dalam layanan pelanggan membantu perusahaan dalam memahami dan merespons kebutuhan pelanggan dengan lebih efisien, meningkatkan interaksi, dan memberikan pengalaman pelanggan yang lebih baik secara keseluruhan.

## 2. Analisis Sentimen

Penerapan Natural Language Processing (NLP) dalam analisis sentimen bertujuan untuk memahami dan mengekstrak sentimen, opini, atau perasaan dari teks yang dihasilkan oleh pengguna,

konsumen, atau pemangku kepentingan. Berikut adalah rincian tentang bagaimana NLP digunakan dalam analisis sentimen:

1. **Preprocessing Data Teks:** Tokenisasi: Pemecahan teks menjadi token, seperti kata-kata, frasa, atau karakter. Stopword Removal: Penghapusan kata-kata umum yang tidak memberikan makna penting dalam analisis sentimen. Stemming atau Lemmatisasi: Normalisasi kata-kata menjadi bentuk dasar mereka untuk mengurangi variasi dalam teks.
2. **Representasi Data Teks:** Word Embeddings: Mengubah teks menjadi vektor numerik menggunakan teknik Word Embeddings seperti Word2Vec, GloVe, atau FastText. TF-IDF (Term Frequency-Inverse Document Frequency): Menghitung bobot kata-kata dalam teks berdasarkan frekuensi kemunculan kata-kata tersebut dalam dokumen dan seberapa umum kata-kata tersebut dalam korpus keseluruhan.
3. **Analisis Sentimen:** Pendekatan Supervised Learning: Menggunakan metode klasifikasi (misalnya, Support Vector Machines, Naive Bayes, atau Neural Networks) yang dilatih dengan data yang dilabeli untuk mengklasifikasikan teks ke dalam kategori sentimen tertentu (positif, negatif, atau netral). Unsupervised Learning: Menggunakan teknik Clustering atau analisis topik untuk mengelompokkan teks ke dalam kelompok sentimen berdasarkan kesamaan topik atau karakteristik.
4. **Emotion Analysis:** Deteksi Emosi: Menerapkan NLP untuk mengidentifikasi emosi atau perasaan tertentu dalam teks seperti kegembiraan, kemarahan, atau kecemasan.
5. **Aspect-Based Sentiment Analysis:** Analisis Berbasis Aspek: Memahami sentimen terkait dengan aspek-aspek tertentu dalam teks, seperti produk dalam ulasan, fitur spesifik, atau layanan yang disediakan.
6. **Pengembangan Model Sentiment Analysis:** Fine-tuning Model Pre-trained: Menggunakan model yang telah dilatih sebelumnya dalam bahasa alami (seperti BERT, GPT, atau Transformer) dan menyesuaikannya dengan tugas analisis sentimen tertentu.

7. **Evaluasi Model: Menggunakan Metrics:** Menggunakan metrik evaluasi seperti akurasi, F1-score, atau Confusion Matrix untuk mengukur kinerja model dalam memprediksi sentimen dengan benar.

Penerapan NLP dalam analisis sentimen memungkinkan perusahaan untuk memahami perasaan pelanggan, umpan balik produk, atau sentimen pasar secara luas. Hal ini membantu dalam pengambilan keputusan yang lebih baik, penyesuaian strategi bisnis, dan meningkatkan interaksi dengan pelanggan berdasarkan pemahaman yang lebih baik tentang sentimen yang terkandung dalam teks.

### 3. Pencarian Informasi

Penerapan Natural Language Processing (NLP) dalam pencarian informasi membantu dalam pemrosesan, pemahaman, dan relevansi pencarian terhadap teks yang dimasukkan pengguna. Berikut adalah rincian tentang bagaimana NLP digunakan dalam pencarian informasi:

1. **Query Understanding: Analisis Pencarian:** NLP digunakan untuk memahami query atau pertanyaan pengguna yang dimasukkan ke dalam mesin pencarian, memecahnya menjadi token dan mengidentifikasi kata kunci penting. **Pemrosesan Bahasa Alami:** Memahami makna atau intent di balik query, termasuk penanganan variasi frasa atau pertanyaan yang mirip namun memiliki struktur yang berbeda.
2. **Pengindeksan Informasi: Tokenisasi dan Representasi:** Dokumen-dokumen atau konten yang akan diindeks dianalisis menggunakan NLP untuk tokenisasi, mengubah teks menjadi representasi vektor, dan menghitung bobot kata-kata (TF-IDF) untuk membangun indeks yang mempercepat proses pencarian. **Entity Recognition:** Mengidentifikasi entitas seperti nama orang, tempat, atau organisasi dalam teks untuk meningkatkan akurasi pencarian.
3. **Relevansi Pencarian: Matching dan Ranking:** NLP digunakan untuk mencocokkan query pengguna dengan dokumen yang

relevan dan memberi peringkat pada hasil pencarian berdasarkan relevansi. **Pemahaman Konteks:** Memperhitungkan konteks dalam pencarian, memastikan hasil yang diberikan sesuai dengan kebutuhan pengguna.

4. Pencarian Semantik: **Analisis Semantik:** Menggunakan NLP untuk memahami arti sebenarnya dari query atau dokumen, bukan hanya kata-kata yang digunakan, melainkan juga konsep yang terkandung dalam teks. **Pemrosesan Teks yang Lebih Lanjut:** Penggunaan teknik seperti Word Embeddings atau Transformer untuk pemahaman yang lebih dalam tentang hubungan antar kata-kata atau makna di balik teks.
5. Personalisasi Pencarian: **Pemahaman User Intent:** NLP membantu dalam memahami intent atau tujuan pengguna yang berbeda, memungkinkan sistem untuk memberikan hasil yang lebih relevan berdasarkan histori pencarian atau profil pengguna.
6. Evaluasi dan Peningkatan Sistem: **Analisis Feedback:** Menggunakan NLP untuk menganalisis umpan balik pengguna terhadap hasil pencarian untuk meningkatkan relevansi dan kualitas hasil.

#### 4. Fine-tuning Model

Meningkatkan model pencarian berbasis pada informasi dari evaluasi dan umpan balik untuk meningkatkan performa dan akurasi. Penerapan NLP dalam pencarian informasi membantu dalam meningkatkan akurasi, relevansi, dan kecepatan pencarian, memastikan bahwa pengguna mendapatkan informasi yang mereka cari dengan lebih efisien dan sesuai dengan kebutuhan mereka.

1. Chatbots dan Virtual Assistants: Penerapan Natural Language Processing (NLP) dalam Chatbots dan Virtual Assistants memungkinkan sistem untuk memahami, memproses, dan merespons bahasa manusia secara efektif. Berikut adalah rincian tentang bagaimana NLP digunakan dalam Chatbots dan Virtual Assistants:
2. Pengenalan dan Pemahaman Bahasa Manusia: **Pemrosesan Bahasa Alami:** NLP digunakan untuk memahami perintah,



- pertanyaan, atau masukan pengguna dalam bahasa alami. Ini melibatkan pemecahan kalimat, tokenisasi, dan pemahaman intent di balik permintaan. **Entity Recognition:** Identifikasi entitas seperti nama orang, lokasi, tanggal, atau objek tertentu dalam teks untuk memberikan respon yang lebih tepat.
3. Pembangunan Chatbots yang Responsif: **Generasi Respon:** NLP membantu dalam menghasilkan respon yang relevan dan kontekstual berdasarkan pemahaman terhadap permintaan pengguna. Ini melibatkan pembuatan respon yang sesuai dengan konteks, bahasa yang ramah, dan pilihan kata yang tepat. **Personalisasi Respon:** Sistem dapat menyesuaikan respon berdasarkan informasi pengguna atau sejarah interaksi sebelumnya.
  4. Pengelolaan Dialog dan Konteks: **Memahami Konteks:** NLP membantu dalam mempertahankan konteks percakapan dan memastikan Chatbot atau Virtual Assistant dapat mengingat percakapan sebelumnya untuk memberikan respon yang lebih baik. **Dialog State Management:** Manajemen status percakapan yang memungkinkan sistem untuk menanggapi permintaan yang berurutan atau berkelanjutan.
  5. Integrasi dengan Pengetahuan dan Informasi Tambahan: **Akses ke Informasi:** NLP memungkinkan Chatbot untuk mengakses basis pengetahuan, database, atau sumber informasi lainnya untuk memberikan jawaban yang lebih lengkap dan akurat.
  6. Evaluasi dan Peningkatan Kualitas Respon: **Analisis Sentimen:** Menggunakan NLP untuk memahami sentimen pengguna terhadap respon yang diberikan oleh Chatbot dan mengadaptasi respons berdasarkan umpan balik. **Peningkatan Model:** Penggunaan umpan balik pengguna dan analisis performa untuk meningkatkan model Chatbot, termasuk fine-tuning model berbasis NLP.
  7. Pengembangan Multilingual Chatbots: **Penerapan Bahasa Lain:** NLP digunakan untuk mendukung chatbot dalam bahasa yang berbeda, memungkinkan sistem untuk beroperasi dalam lingkungan multilingual. Penerapan NLP dalam Chatbots dan Virtual Assistants membantu dalam memberikan pengalaman interaktif yang lebih manusiawi, efisien, dan responsif bagi pengguna. Kemampuan sistem untuk memahami bahasa

manusia secara alami merupakan inti dari efektivitas Chatbot dalam memberikan layanan yang berguna dan informatif kepada pengguna.

## **B. NLP dalam Kesehatan**

1. Analisis data medis: Penerapan Natural Language Processing (NLP) dalam analisis data medis bertujuan untuk memahami, mengekstrak, dan mengelola informasi dari catatan medis, laporan laboratorium, dokumen kesehatan, atau literatur medis. Berikut adalah rincian tentang bagaimana NLP digunakan dalam analisis data medis:
2. Pemrosesan Catatan Medis: **Ekstraksi Informasi**: NLP digunakan untuk mengekstrak informasi klinis seperti diagnosis, tindakan medis, gejala, atau riwayat penyakit dari catatan medis yang sering kali terstruktur atau tidak terstruktur. **Named Entity Recognition (NER)**: Mengidentifikasi entitas medis seperti nama pasien, dokter, obat-obatan, atau prosedur medis dalam catatan medis.
3. Klasifikasi dan Analisis Teks Medis: **Diagnosis Otomatis**: Penggunaan NLP dalam klasifikasi teks medis untuk mendukung diagnosa otomatis berdasarkan informasi yang terkandung dalam catatan medis. **Analisis Sentimen Kesehatan**: Menganalisis catatan medis untuk mengevaluasi sentimen pasien terhadap pengalaman perawatan atau prosedur medis tertentu.
4. Penelitian dan Literatur Medis: **Literature Review**: NLP digunakan untuk memproses dan menganalisis literatur medis yang luas, membantu peneliti untuk mendapatkan wawasan dari artikel dan penelitian terbaru dalam bidang kesehatan. **Pengelompokan Tema**: Mengelompokkan artikel atau makalah medis berdasarkan tema tertentu menggunakan teknik pengelompokan topik.
5. Pemahaman Bahasa Kesehatan: **Terminologi Medis**: Memahami istilah medis yang kompleks dan terminologi khusus yang digunakan dalam catatan medis atau literatur medis. **Analisis Percakapan Medis**: Memahami percakapan

- antara dokter dan pasien dalam rekaman medis untuk meningkatkan pemahaman terhadap situasi kesehatan pasien.
6. **Prediksi Penyakit dan Perawatan: Prediksi Risiko:** Menggunakan NLP untuk mengidentifikasi faktor risiko atau prediksi perjalanan penyakit berdasarkan informasi yang terdapat dalam catatan medis.
  7. **Privasi dan Keamanan Data: Anonimisasi Data:** Penggunaan NLP dalam menghapus atau mengaburkan informasi identitas pribadi dari catatan medis untuk menjaga keamanan dan privasi data.
  8. **Pengembangan Sistem Berbasis NLP: Sistem Dukungan Keputusan:** Membangun sistem NLP yang mendukung pengambilan keputusan klinis atau memberikan saran terhadap perawatan medis.

Penerapan NLP dalam analisis data medis membuka potensi besar untuk meningkatkan pengelolaan data kesehatan, penelitian medis, pelayanan kesehatan, dan pengembangan sistem yang mendukung pengambilan keputusan klinis yang lebih baik. Ini juga memainkan peran penting dalam meningkatkan efisiensi, akurasi, dan pemahaman terhadap informasi kesehatan yang terkandung dalam dokumen medis.

Penerapan Natural Language Processing (NLP) dalam pengenalan entitas medis (NER - Named Entity Recognition) adalah tentang mengidentifikasi dan mengekstrak entitas spesifik dalam teks medis seperti nama pasien, nama dokter, jenis penyakit, obat-obatan, prosedur medis, tanggal, dan informasi penting lainnya. Berikut adalah rincian tentang penerapan NLP dalam NER untuk pengenalan entitas medis:

1. **Pemrosesan Teks Medis: Tokenisasi:** Memecah teks medis menjadi token (kata-kata, frasa, atau bagian-bagian lain) untuk analisis lebih lanjut. **Stopword Removal:** Penghapusan kata-kata umum yang tidak relevan dalam teks medis.
2. **Penggunaan Model NLP untuk Pengenalan Entitas Medis: Model Berbasis Aturan (Rule-Based):** Menerapkan aturan linguistik atau peraturan manual untuk mengidentifikasi entitas medis. Contohnya, pengenalan nama orang berdasarkan pola penulisan nama manusia. **Machine Learning-Based Models:**

- Penggunaan algoritma pembelajaran mesin seperti Conditional Random Fields (CRFs), Support Vector Machines (SVM), atau Deep Learning (misalnya, LSTM atau Transformer) yang telah dilatih pada data yang dilabeli untuk mengenali entitas medis.
3. Pengembangan Anotasi dan Dataset: **Anotasi Manual:** Menandai atau memberi label entitas medis dalam teks medis oleh ahli manusia untuk membuat dataset pelatihan yang dilabeli.
  4. Feature Engineering: **Penggunaan Fitur:** Pemilihan fitur yang relevan seperti kata-kata sekitar, morfologi kata, atau konteks untuk membantu model dalam mengenali entitas medis dengan lebih akurat.
  5. Evaluasi dan Peningkatan Model: **Cross-Validation:** Menggunakan teknik cross-validation untuk mengukur kinerja model dalam mengenali entitas medis dan menghindari overfitting. **Fine-Tuning Model:** Memperbarui atau menyesuaikan model NER berdasarkan umpan balik dari evaluasi hasil model untuk meningkatkan akurasi.
  6. Penerapan dalam Aplikasi Kesehatan: **Sistem Manajemen Kesehatan Elektronik:** Menggunakan NER untuk mengekstrak dan mengelola informasi penting dalam catatan medis elektronik untuk memfasilitasi pencarian, analisis, dan perawatan pasien. **Penelitian Klinis:** Penerapan NER dalam analisis literatur medis untuk mengidentifikasi informasi penting dalam artikel penelitian atau makalah medis. Penerapan NLP dalam NER untuk pengenalan entitas medis merupakan langkah penting dalam pengelolaan data kesehatan, penelitian medis, dan perawatan pasien yang memungkinkan pengambilan informasi yang lebih cepat dan akurat dari teks medis yang besar dan kompleks.

## C. NLP dalam Pendidikan

### 1. Evaluasi dan pembelajaran berbasis teks

Penerapan Natural Language Processing (NLP) dalam evaluasi dan pembelajaran berbasis teks melibatkan analisis teks untuk mengukur kinerja, meningkatkan pemahaman, dan

mengembangkan sistem yang mendukung pendidikan dan evaluasi. Berikut adalah rincian tentang bagaimana NLP digunakan dalam konteks ini:

1. **Penilaian Otomatis: Analisis Jawaban Siswa:** Penggunaan NLP dalam mengevaluasi jawaban siswa dalam bentuk teks, mengidentifikasi kesalahan atau kekurangan dalam jawaban mereka, serta memberikan umpan balik yang sesuai. **Penilaian Tugas:** NLP digunakan untuk memberikan penilaian otomatis terhadap tugas yang mencakup teks, seperti esai, tugas menulis, atau penugasan proyek.
2. **Analisis Sentimen dan Partisipasi: Pemantauan Sentimen:** Menganalisis sentimen dari diskusi kelas, tanggapan siswa, atau umpan balik untuk memahami perasaan dan tingkat partisipasi.
3. **Pemahaman Konten: Ringkasan Otomatis:** Menggunakan NLP untuk merangkum teks panjang, seperti materi pelajaran atau artikel, agar lebih mudah dipahami oleh siswa. **Pemahaman Materi Pelajaran:** Penerapan NLP untuk memahami pertanyaan siswa, memberikan informasi tambahan, atau menjelaskan konsep yang rumit dalam teks.
4. **Peningkatan Pengalaman Belajar: Personalisasi Pembelajaran:** Memanfaatkan NLP untuk mempersonalisasi pengalaman belajar siswa berdasarkan kemajuan mereka, preferensi, dan kebutuhan individu. **Rekomendasi Konten:** Menyediakan rekomendasi materi pembelajaran berdasarkan minat dan kemajuan siswa.
5. **Analisis Diskusi Kelas dan Forum: Analisis Diskusi:** Menggunakan NLP untuk menganalisis percakapan atau diskusi dalam forum online atau kelas virtual guna mengidentifikasi topik populer, pola partisipasi, atau pemahaman umum.
6. **Pengembangan Sistem Pendidikan: Sistem Tutor Cerdas:** Membangun sistem tutor yang menggunakan NLP untuk memahami kebutuhan siswa, memberikan bantuan, dan menyesuaikan pembelajaran. **Pengembangan Platform Pembelajaran:** Integrasi NLP dalam platform pembelajaran daring untuk meningkatkan interaksi, pembelajaran adaptif, dan evaluasi.

Penerapan NLP dalam evaluasi dan pembelajaran berbasis teks memberikan peluang untuk meningkatkan efisiensi dalam penilaian, personalisasi pembelajaran, dan pemahaman konten secara lebih baik. Hal ini juga mendukung perkembangan sistem pendidikan yang lebih adaptif, interaktif, dan responsif terhadap kebutuhan individual siswa.

## 2. Analisis plagiarism

Penerapan Natural Language Processing (NLP) dalam analisis plagiarisme adalah tentang penggunaan teknologi bahasa alami untuk mendeteksi dan menganalisis kesamaan atau plagiarisme antara dokumen atau teks. Berikut adalah rincian tentang bagaimana NLP diterapkan dalam konteks ini:

1. **Preprocessing Teks: Tokenisasi dan Representasi:** Mengubah teks ke dalam representasi numerik atau token untuk analisis lebih lanjut. **Pembersihan Teks:** Membersihkan teks dari informasi yang tidak relevan, seperti tanda baca, karakter khusus, atau formatting.
2. **Penggunaan Model NLP: Model Berbasis Aturan (Rule-Based):** Menerapkan aturan linguistik atau logika untuk mendeteksi kesamaan teks, terutama dalam dokumen panjang atau struktur kompleks. **Machine Learning-Based Models:** Penggunaan algoritma pembelajaran mesin (misalnya, Decision Trees, Support Vector Machines, atau Neural Networks) yang dilatih dengan data yang dilabeli untuk mengidentifikasi pola plagiarisme.
3. **Analisis Struktural Teks: Alignment dan Similarity Detection:** Menggunakan NLP untuk menemukan kesamaan atau kemiripan antara teks, baik dalam frasa, kalimat, paragraf, atau struktur keseluruhan. **Deteksi Plagiarisme Paragraf atau Dokumen:** Menggunakan teknik seperti cosine similarity atau Levenshtein distance untuk mendeteksi plagiarisme pada tingkat dokumen atau paragraf.
4. **Analisis Konten dan Semantik: Kesamaan Makna:** Menggunakan NLP untuk memahami makna di balik kata-kata atau frasa, bukan hanya kesamaan kata. **Analisis Semantik:**

Mengidentifikasi makna dan inti dari teks untuk membandingkan konten secara lebih mendalam.

5. Pengembangan Algoritma Khusus: **Fine-Tuning Model:** Meningkatkan model NLP untuk deteksi plagiarisme, memastikan sensitivitas yang tinggi dan akurasi dalam mengidentifikasi kesamaan teks.
6. Evaluasi Hasil: **Pengukuran Similaritas:** Mengukur tingkat kesamaan atau plagiarisme antara teks berdasarkan hasil analisis NLP. **Umpan Balik dan Peningkatan:** Menggunakan umpan balik untuk meningkatkan algoritma deteksi plagiarisme dan meningkatkan ketepatan serta ketelitian hasil.

Penerapan NLP dalam analisis plagiarisme memiliki peran penting dalam memastikan keaslian dan integritas karya tulis. Ini membantu institusi pendidikan, editor, peneliti, atau penerbit untuk mengidentifikasi plagiarisme dengan lebih efisien dan akurat, serta menjaga kejujuran dalam publikasi karya.

# **BAB IV**

## **TANTANGAN DAN ISU ETIKA DALAM NLP**

### **A. Tantangan dalam NLP**

#### **1. Polysemy dan Ambiguitas**

Dalam Natural Language Processing (NLP), terdapat sejumlah tantangan dan isu etika yang berkaitan dengan polisemi (polysemy) dan ambiguitas dalam bahasa:

1. Polisemi dan Ambiguitas: Arti Ganda: Polisemi merujuk pada kata-kata atau frasa yang memiliki lebih dari satu arti yang sah. Ambiguitas mencakup situasi di mana kalimat atau teks memiliki arti yang tidak jelas atau lebih dari satu arti yang memungkinkan.
2. Kesulitan dalam Pemahaman Konteks: Konteks yang Membingungkan: Penafsiran kata-kata atau frasa terkadang tergantung pada konteksnya. Meskipun NLP mungkin mampu mengenali variasi kata, memahami makna sebenarnya dalam konteks yang tepat bisa menjadi sulit.
3. Tantangan dalam Pemrosesan Bahasa Alami: Disambiguasi: NLP harus mampu untuk mengatasi polisemi dan ambiguitas, memutuskan makna yang benar berdasarkan konteks yang diberikan. Ini sering kali menantang karena bahasa manusia penuh dengan kompleksitas dan nuansa.
4. Isu Etika dalam Penggunaan Teknologi: Bias dan Penafsiran yang Salah: Ketika NLP menghadapi polisemi dan ambiguitas, ada risiko penafsiran yang salah atau bias dalam analisis, yang dapat memengaruhi hasil dan keputusan yang dibuat oleh sistem berbasis NLP.
5. Dampak pada Aplikasi NLP: Ketepatan dalam Pemrosesan Teks: Polisemi dan ambiguitas dapat mengganggu ketepatan hasil analisis NLP, misalnya dalam kasus klasifikasi teks atau pemahaman konten.



6. Perlunya Penanganan yang Lebih Cermat: Peningkatan Algoritma: Diperlukan pengembangan algoritma NLP yang lebih canggih untuk menangani polisemi dan ambiguitas secara lebih efektif.

Isu Etika Terkait: **Transparansi dan Akuntabilitas:** Penerapan NLP yang tidak mempertimbangkan polisemi atau ambiguitas bisa mengakibatkan kesalahan atau interpretasi yang salah, yang bisa menjadi isu etika jika hal itu memengaruhi keputusan penting atau menimbulkan bias.

**Privasi dan Keamanan:** Ketika NLP digunakan dalam aplikasi yang melibatkan data sensitif, risiko salah tafsir atau manipulasi akibat polisemi atau ambiguitas bisa memengaruhi privasi atau keamanan data. Penanganan polisemi dan ambiguitas dalam NLP adalah tantangan penting karena memengaruhi tingkat akurasi, keandalan, dan interpretasi yang tepat dari teks dalam konteks yang berbeda. Isu etika terkait juga harus dipertimbangkan secara cermat untuk memastikan penggunaan NLP yang bertanggung jawab dan tepat.

## 2. Kurangnya data yang berkualitas

Salah satu tantangan utama dalam Natural Language Processing (NLP) adalah kurangnya data yang berkualitas. Hal ini bisa menjadi hambatan serius dalam mengembangkan model NLP yang baik. Berikut adalah penjelasan lebih rinci:

1. Keterbatasan Dataset: **Kurangnya Volume Data:** Dalam beberapa kasus, dataset yang tersedia untuk pelatihan model NLP bisa sangat terbatas, terutama untuk bahasa yang kurang umum atau domain tertentu seperti medis atau hukum. **Kualitas Data yang Buruk:** Data yang tidak terstruktur, tidak terlabeli dengan baik, atau tidak terkumpul dengan baik dapat menghambat kemampuan model untuk belajar dengan baik.
2. Tantangan Variasi Bahasa: **Variasi Dialek dan Gaya Bahasa:** Bahasa manusia sangat bervariasi, termasuk penggunaan dialek, slang, atau variasi dalam gaya penulisan. Ini membutuhkan data yang representatif untuk melatih model

- agar mampu memahami variasi bahasa tersebut. **Bahasa yang Kurang Dikenal:** Bahasa-bahasa yang kurang umum atau kuno sering kali memiliki keterbatasan dalam data, membuat pengembangan model NLP yang akurat menjadi sulit.
3. Biaya dan Waktu Pengumpulan Data: **Biaya Pengumpulan Data:** Mengumpulkan dataset yang besar dan berkualitas memerlukan sumber daya yang signifikan, baik itu biaya maupun waktu. **Keterbatasan Waktu:** Dalam beberapa kasus, pembangunan model NLP yang baik membutuhkan waktu yang lama karena proses pengumpulan, pembersihan, dan anotasi data yang memadai.
  4. Ketergantungan pada Data Label: **Ketergantungan pada Data yang Dilabeli:** Model pembelajaran mesin sering kali memerlukan data yang sudah dilabeli dengan benar untuk melatih dan memvalidasi kinerja. Kurangnya data yang dilabeli bisa menjadi kendala. Cara Mengatasinya: **Augmentasi Data:** Menciptakan data tambahan dari data yang ada dengan teknik seperti penggandaan, translasi, atau penggabungan untuk meningkatkan jumlah dan variasi data. **Transfer Learning:** Memanfaatkan model yang sudah dilatih pada data yang besar (pre-trained models) dan menyesuaikannya dengan data yang tersedia dalam domain atau bahasa tertentu. **Collaborative Efforts:** Kerja sama dan pertukaran dataset antara lembaga, peneliti, atau komunitas dapat membantu mengatasi keterbatasan data.

Kurangnya data yang berkualitas bisa menjadi tantangan utama dalam pengembangan model NLP yang akurat dan andal. Strategi pengumpulan data yang cerdas, teknik augmentasi data, dan pemanfaatan model yang sudah dilatih dapat membantu mengatasi sebagian dari kendala ini dalam mengembangkan model NLP yang lebih baik.

### 3. Overfitting dan generalisasi

Dalam konteks Natural Language Processing (NLP), overfitting dan generalisasi adalah dua konsep penting yang memengaruhi

kualitas dan performa model yang dikembangkan untuk pemrosesan bahasa alami.

1. **Overfitting:** Definisi Overfitting: Overfitting terjadi ketika model terlalu "memorize" data pelatihan dan tidak mampu melakukan generalisasi dengan baik pada data baru atau data yang belum pernah dilihat sebelumnya. Penyebab Overfitting: Overfitting seringkali terjadi ketika model terlalu kompleks atau memiliki kapasitas yang berlebihan untuk mempelajari detail-detail kecil yang sebenarnya bersifat acak atau tidak relevan dalam data. Indikasi Overfitting: Biasanya, tanda-tanda overfitting termasuk performa model yang sangat baik pada data pelatihan tetapi performa yang buruk pada data validasi atau data uji. Strategi Penanggulangan: Menggunakan teknik regularisasi seperti dropout, pengurangan kompleksitas model, atau menggunakan teknik validasi silang untuk memvalidasi performa model.
2. **Generalisasi:** Definisi Generalisasi: Generalisasi adalah kemampuan model untuk mengadopsi pola yang ditemukan dari data pelatihan dan menerapkannya dengan baik pada data baru atau data yang belum dilihat sebelumnya. Penyebab Generalisasi: Model yang mampu menangkap pola yang umum dan relevan dari data pelatihan tanpa terlalu fokus pada detail yang mungkin bersifat acak. Indikasi Generalisasi: Model yang baik dalam generalisasi akan menunjukkan performa yang konsisten pada data yang tidak pernah dilihat selama pelatihan. Strategi Peningkatan Generalisasi: Menggunakan teknik penambahan data, pengaturan yang tepat terkait kompleksitas model, dan menggunakan metode regularisasi yang tepat untuk mencegah overfitting.

**Relevansi dalam NLP:** Dalam NLP, overfitting bisa terjadi saat model NLP terlalu "memorize" teks pelatihan dengan sangat baik, tetapi tidak bisa menerapkan pemahaman yang diperolehnya pada teks baru yang berbeda. Generalisasi yang baik dalam NLP menunjukkan kemampuan model untuk memahami bahasa secara umum tanpa terlalu terkait dengan detail-detail yang mungkin tidak relevan atau acak dalam teks. Tantangan dalam NLP adalah membangun model yang memiliki tingkat generalisasi yang tinggi

sehingga dapat memproses dan memahami beragam jenis teks dengan akurat, bahkan teks yang belum pernah dilihat sebelumnya. Memahami konsep overfitting dan generalisasi penting dalam mengembangkan model NLP yang handal dan efektif dalam memahami, memproses, dan menghasilkan hasil yang akurat dari teks dalam berbagai konteks dan jenis data.

## **B. Isu Etika dalam NLP**

### **1. Privasi dan keamanan data**

Isu etika privasi dan keamanan data dalam Natural Language Processing (NLP) menjadi sangat penting karena penggunaan data teks yang melibatkan informasi pribadi atau sensitif dari individu atau kelompok. Berikut adalah beberapa poin terkait isu etika ini:

1. **Privasi Data: Penggunaan Informasi Pribadi:** Penggunaan data teks yang mengandung informasi pribadi seperti riwayat medis, percakapan pribadi, atau informasi identitas individu menimbulkan kekhawatiran privasi. **Risiko Identifikasi:** Proses analisis NLP yang tidak memadai bisa mengungkap informasi sensitif yang dapat mengidentifikasi individu, bahkan jika nama tidak disebutkan.
2. **Keamanan Data: Kekhawatiran Keamanan:** Data teks yang disimpan, diproses, atau ditransmisikan dalam sistem NLP rentan terhadap ancaman keamanan seperti peretasan atau akses tidak sah. **Risiko Penyalahgunaan Informasi:** Data teks yang tidak terlindungi dapat disalahgunakan untuk tujuan jahat seperti penipuan, pencurian identitas, atau penargetan individu.
3. **Isu Etika Terkait: Transparansi Penggunaan Data:** Pentingnya memberikan informasi kepada pengguna terkait bagaimana data mereka digunakan dalam sistem NLP dan untuk tujuan apa. **Konsentisasi dan Izin:** Menghargai hak privasi dan mendapatkan izin atau persetujuan dari individu sebelum menggunakan atau memproses data teks mereka. **Pemulihan dan Hapus Data:** Menciptakan mekanisme untuk menghapus atau memulihkan data teks secara efektif jika diminta oleh individu terkait hak privasi mereka.

4. Penerapan Etika dalam Pengembangan Model NLP: Pengembangan Model yang Bertanggung Jawab: Pentingnya membangun model NLP dengan mempertimbangkan prinsip-prinsip privasi dan keamanan, serta memastikan bahwa data sensitif diperlakukan dengan hati-hati.

Enkripsi dan Perlindungan Data: Menggunakan teknologi enkripsi dan pengamanan data yang tepat untuk melindungi informasi yang disimpan dan diproses oleh sistem NLP. Kerangka Regulasi: Perlunya kerangka regulasi yang kuat untuk mengatur penggunaan data teks dalam NLP, memastikan perlindungan yang tepat terhadap privasi dan keamanan. Menyadari dan mempertimbangkan isu privasi dan keamanan data dalam pengembangan dan penerapan teknologi NLP sangat penting untuk memastikan penggunaan yang etis, aman, dan bertanggung jawab dari informasi teks yang sensitif atau pribadi.

## **2. Bias dalam data dan model**

Isu etika tentang bias dalam data dan model dalam Natural Language Processing (NLP) mengacu pada ketidakseimbangan atau distorsi dalam data serta model yang dapat menyebabkan hasil yang tidak adil atau tidak representatif. Berikut adalah beberapa poin terkait isu etika ini:

1. Bias dalam Data: Ketidakseimbangan Representasi: Data yang digunakan untuk melatih model NLP mungkin tidak mencerminkan keberagaman masyarakat, menyebabkan kurangnya representasi yang merata dari berbagai kelompok atau perspektif. Replikasi Bias Manusia: Data teks bisa mencerminkan bias yang ada dalam masyarakat, seperti gender, ras, atau kecenderungan budaya, yang dapat tercermin dalam model NLP.
2. Bias dalam Model: Pengambilan Keputusan Tidak Adil: Model NLP yang dikenai bias dalam data latihnya dapat menghasilkan keputusan atau penilaian yang tidak adil atau diskriminatif. Perpetuasi Bias: Model NLP yang belajar dari data yang sudah terbias cenderung memperkuat atau memperpanjang bias tersebut dalam hasilnya.

3. Isu Etika Terkait: Keadilan dan Kesetaraan: Model NLP yang bias dapat memberikan keputusan atau prediksi yang tidak adil, memengaruhi kesetaraan akses atau perlakuan yang adil. Transparansi dan Akuntabilitas: Perlunya transparansi dalam proses pembangunan model dan pengambilan keputusan untuk memahami dan memeriksa bias yang ada. Keragaman dan Representasi: Pentingnya memastikan keberagaman dan representasi yang adil dari berbagai perspektif dalam data dan model NLP.
4. Penanganan Bias dalam NLP: Pembersihan Data: Identifikasi, analisis, dan pembersihan data yang memiliki bias yang tidak diinginkan. Pengaturan Model: Menerapkan teknik seperti debiasing atau fine-tuning untuk mengurangi atau menghilangkan bias yang ditemukan dalam model. Monitoring dan Evaluasi Berkelanjutan: Melakukan evaluasi berkelanjutan terhadap model untuk mengidentifikasi dan mengatasi bias yang baru muncul. Pendekatan yang Berbasis Etika: Menggunakan pendekatan yang berbasis etika dalam pengembangan model untuk memastikan keadilan, transparansi, dan kesetaraan. Mengatasi isu bias dalam data dan model NLP sangat penting untuk memastikan bahwa teknologi ini diterapkan secara adil, transparan, dan menghormati keberagaman serta hak asasi manusia. Hal ini membantu mencegah model NLP menyebarkan atau memperkuat ketidaksetaraan yang ada dalam masyarakat.

## **BAB V**

### **TOPIK MODEL**

#### **A. Topik Model LDA**

Dalam dunia yang dipenuhi dengan ledakan informasi, pengelolaan dan pemahaman terhadap teks telah menjadi tantangan besar. Bagaimana kita bisa mengurai ratusan, bahkan ribuan dokumen, untuk menemukan pola dan tema yang tersembunyi di dalamnya? Inilah di mana Model Latent Dirichlet Allocation (LDA) memasuki panggung sebagai alat penting dalam pemrosesan teks dan analisis topik.

LDA, yang merupakan singkatan dari Latent Dirichlet Allocation, adalah sebuah model probabilistik yang memungkinkan kita untuk mengidentifikasi topik-topik yang tersembunyi di dalam sebuah koleksi besar dokumen. Konsep utama di balik LDA adalah ide bahwa setiap dokumen dalam koleksi tersebut merupakan kombinasi dari beberapa topik, sementara setiap topik sendiri adalah distribusi probabilitas atas sekelompok kata-kata.

Mengapa LDA penting? Alat ini memungkinkan kita untuk menjelajahi teks dengan cara yang tidak terlalu langsung, dengan mengidentifikasi hubungan dan tema yang ada di antara kumpulan kata-kata yang mungkin tidak terlihat pada pandangan pertama. Dengan kemampuannya untuk menemukan pola tersembunyi, LDA telah diterapkan dalam berbagai bidang mulai dari analisis sentimen hingga klasifikasi dokumen, serta pemahaman yang lebih dalam terhadap tren dan opini dalam teks yang besar.

Namun, seperti halnya alat analisis lainnya, LDA memiliki kelebihan dan batasannya sendiri. Penggunaannya yang efektif membutuhkan pemahaman yang baik akan parameter, proses preprocessing data yang teliti, serta interpretasi hasil yang tepat. Meskipun memberikan wawasan yang kuat, model ini juga

memerlukan penggunaan yang bijaksana dan penyesuaian yang cermat sesuai dengan konteks aplikasinya.

Dalam buku ini, kami akan membawa Anda melalui perjalanan mendalam dari dasar-dasar probabilitas hingga implementasi praktis dari Model LDA. Kami akan membahas teori di balik model ini, langkah-langkah untuk menerapkannya secara efektif, strategi evaluasi, dan juga memperlihatkan berbagai studi kasus yang memperlihatkan aplikasi nyata dari model ini. Semoga buku ini membantu Anda memahami, menerapkan, dan mengambil manfaat dari kekuatan analisis teks yang ditawarkan oleh Model Latent Dirichlet Allocation. Selamat menikmati perjalanan Anda dalam mempelajari model yang luar biasa ini. Ruang lingkup dan tujuan dalam penggunaan Model Latent Dirichlet Allocation (LDA) sangat penting untuk memberikan pemahaman yang jelas kepada pembaca tentang apa yang dapat dicapai dengan model ini dan bagaimana mereka bisa menerapkannya secara praktis. Berikut penjelasan mengenai ruang lingkup dan tujuan LDA:

## **B. Ruang Lingkup LDA**

Ruang lingkup LDA meliputi pemahaman tentang bagaimana model ini digunakan untuk menganalisis teks secara probabilistik. Dalam penggunaannya, LDA membantu dalam:

1. **Penemuan Topik Tersembunyi:** LDA membantu mengidentifikasi pola dan topik yang tersembunyi di dalam kumpulan dokumen, memungkinkan kita untuk mengetahui topik apa saja yang sedang dibahas.
2. **Representasi Dokumen:** Model ini memungkinkan dokumen direpresentasikan sebagai distribusi topik, memberikan cara yang kuat untuk melihat bagaimana dokumen terkait dengan topik-topik tertentu.
3. **Analisis Sentimen dan Klasifikasi Dokumen:** Dengan memahami topik utama dalam dokumen, LDA dapat digunakan untuk menganalisis sentimen, mengelompokkan dokumen ke dalam kategori tertentu, atau bahkan membantu dalam pemrosesan pencarian.



### **C. Tujuan Implementasi LDA**

**Pemahaman Teori Probabilistik di Balik LDA:** Tujuan pertama adalah memberikan pemahaman yang kuat tentang dasar-dasar probabilistik yang mendasari model LDA sehingga pembaca dapat mengerti alasan di balik proses dan hasilnya. **Implementasi Praktis dengan Alat yang Tersedia:** Buku ini bertujuan untuk membantu pembaca dalam mengimplementasikan LDA dengan alat dan bahasa pemrograman yang umum digunakan seperti Python, R, atau bahasa lainnya yang mendukung analisis teks.

**Strategi Preprocessing dan Evaluasi yang Efektif:** Penting bagi pembaca untuk memahami langkah-langkah pra-pemrosesan data yang diperlukan sebelum menerapkan LDA, serta cara melakukan evaluasi yang tepat terhadap model yang telah dibangun. **Studi Kasus dan Contoh yang Nyata:** Buku ini akan memaparkan studi kasus yang bervariasi dan contoh penggunaan nyata LDA di berbagai bidang agar pembaca mendapatkan gambaran yang jelas tentang cara praktis dalam menerapkan model ini. Dengan memahami ruang lingkup dan tujuan penggunaan LDA, diharapkan pembaca dapat merencanakan, mengimplementasikan, dan mengevaluasi model ini secara efektif untuk kebutuhan analisis teks mereka.

Dalam dunia pembelajaran mesin, model probabilistik menjadi fondasi yang kuat untuk pemahaman dan analisis data. Dasar-dasar model probabilistik mengacu pada representasi matematis dari ketidakpastian dalam suatu sistem. Konsep ini melibatkan probabilitas sebagai alat utama untuk menggambarkan ketidakpastian dalam data. Dalam konteks model probabilistik, variabel yang diamati diasumsikan memiliki distribusi probabilitas tertentu yang menentukan kemungkinan nilai-nilai yang mungkin mereka miliki.

Model probabilistik menawarkan pendekatan kuat untuk memahami dan memodelkan data yang kompleks. Dengan memperhitungkan distribusi probabilitas dari berbagai variabel dan parameter, model-model ini memungkinkan penanganan

ketidakpastian dengan cara yang sistematis. Penggunaannya yang luas meliputi pembelajaran mesin, di mana model probabilistik digunakan untuk membuat prediksi yang bergantung pada distribusi probabilitas, bukan hanya untuk memberikan hasil biner atau deterministik.

Dasar-dasar model probabilistik melibatkan konsep teoritis seperti distribusi probabilitas, fungsi likelihood, teori keputusan, dan inferensi statistik. Melalui representasi matematis yang rumit namun sistematis, model probabilistik memungkinkan kita untuk mengeksplorasi dan menganalisis data secara lebih mendalam. Dalam konteks pembelajaran mesin, ini memungkinkan pengembangan model yang mampu mengidentifikasi pola, menarik kesimpulan, dan membuat keputusan berdasarkan analisis statistik yang kuat. Dengan demikian, pemahaman yang kuat tentang dasar-dasar model probabilistik menjadi krusial dalam menjelajahi dan menerapkan teknik-teknik analisis data yang lebih canggih.

#### **D. Konsep Dasar Probabilistik**

Konsep Dasar Probabilistik merujuk pada teori dan prinsip yang mendasari penggunaan probabilitas dalam pemodelan fenomena yang tidak pasti. Probabilitas adalah ukuran untuk mengukur seberapa mungkin suatu peristiwa akan terjadi, dan konsep dasar probabilistik digunakan untuk menggambarkan ketidakpastian dalam berbagai situasi.

Di dalamnya terdapat beberapa konsep utama, salah satunya adalah Distribusi Probabilitas. Ini merujuk pada cara peristiwa acak atau variabel acak tersebar di berbagai nilai dengan berbagai kemungkinan. Distribusi probabilitas memungkinkan kita untuk menggambarkan peluang masing-masing nilai yang mungkin diambil oleh variabel acak.

Selain itu, konsep dasar probabilistik juga mencakup Fungsi Likelihood. Ini menggambarkan seberapa baik suatu model statistik cocok dengan data yang diamati. Fungsi Likelihood

menjadi dasar bagi banyak metode estimasi parameter dalam statistika, dan digunakan untuk mengukur seberapa mungkin parameter model yang diestimasi memproduksi data yang diamati. Konsep dasar ini juga mencakup Teori Keputusan, yang berkaitan dengan cara kita membuat keputusan dalam kondisi ketidakpastian. Teori Keputusan berusaha untuk menggabungkan aspek keuntungan (reward) dan risiko dalam pengambilan keputusan dengan mempertimbangkan probabilitas dan dampak dari pilihan yang dibuat.

Inferensi Statistik juga merupakan bagian penting dari konsep dasar probabilistik. Ini berkaitan dengan proses membuat kesimpulan atau generalisasi tentang populasi atau fenomena berdasarkan data yang hanya diambil dari sampel. Dengan menggunakan prinsip-prinsip dasar probabilitas, inferensi statistik memungkinkan kita untuk melakukan generalisasi yang masuk akal dari data sampel ke populasi yang lebih besar.

Keseluruhan, Konsep Dasar Probabilistik menyediakan kerangka kerja matematis dan konseptual yang penting dalam memahami dan menerapkan teori probabilitas. Ini membantu kita untuk memodelkan dan memahami fenomena kompleks dengan menggambarkan ketidakpastian, membuat keputusan dalam kondisi tidak pasti, serta membuat inferensi yang dapat dipercaya berdasarkan data yang terbatas.

### **E. Contoh Konsep Dasar Probabilistik:**

Berikut adalah beberapa contoh Konsep Dasar Probabilistik:

1. **Distribusi Probabilitas:** Misalkan Anda melempar koin. Kemungkinan hasilnya adalah gambar (heads) atau angka (tails), di mana masing-masing hasil memiliki probabilitas 0.5 (asumsi koin yang adil). Distribusi probabilitas ini membantu menggambarkan kemungkinan hasil yang mungkin dari eksperimen acak ini.
2. **Fungsi Likelihood:** Bayangkan Anda memiliki data pengamatan tentang tinggi badan orang-orang di suatu populasi. Dengan menggunakan model statistik, Anda ingin

menemukan distribusi tinggi badan yang paling mungkin mewakili data yang diamati. Dalam konteks ini, fungsi likelihood membantu mengukur seberapa baik model distribusi tinggi badan ini cocok dengan data yang ada.

3. **Teori Keputusan:** Anda berada di supermarket dan ingin memilih antara dua merek produk dengan harga yang berbeda. Anda tidak yakin kualitas produk mana yang lebih baik. Melalui teori keputusan, Anda mempertimbangkan kemungkinan manfaat dari masing-masing pilihan berdasarkan harga dan probabilitas bahwa salah satu merek produk lebih baik daripada yang lain.
4. **Inferensi Statistik:** Anda ingin mengetahui rata-rata waktu yang dibutuhkan seseorang untuk menyelesaikan tes tertentu. Anda hanya memiliki data waktu yang diperlukan oleh sekelompok sampel orang. Dengan menggunakan inferensi statistik, Anda dapat membuat perkiraan rata-rata waktu yang diperlukan oleh seluruh populasi berdasarkan data sampel ini.

Semua contoh di atas menunjukkan penerapan Konsep Dasar Probabilistik dalam berbagai konteks, mulai dari eksperimen acak hingga pengambilan keputusan dan estimasi parameter dari data terbatas. Konsep-konsep ini membantu dalam menggambarkan, memahami, dan membuat prediksi dalam situasi di mana terdapat ketidakpastian atau variasi dalam hasil yang mungkin terjadi.

Model Probabilistik dalam Pembelajaran Mesin mengacu pada pendekatan di mana model statistik menggunakan konsep probabilitas untuk memodelkan dan mengevaluasi data. Ini adalah salah satu pendekatan yang sangat berguna dalam membuat estimasi, klasifikasi, dan prediksi berdasarkan data yang tidak pasti.

Dalam pembelajaran mesin, model-model ini memungkinkan kita untuk menggabungkan informasi dari data yang diberikan dengan ketidakpastian yang melekat pada proses pengambilan keputusan. Beberapa konsep utama dalam model probabilistik pembelajaran mesin termasuk:

1. **Probabilitas sebagai Landasan Utama:** Model probabilistik menggunakan probabilitas sebagai dasar untuk memahami dan memodelkan ketidakpastian dalam data. Mereka mengekspresikan hubungan antara input dan output dengan distribusi probabilitas, yang membantu dalam mengukur ketidakpastian dalam prediksi.
2. **Pemodelan Distribusi Data:** Model probabilistik mampu memodelkan distribusi data yang kompleks, memberikan cara yang lebih fleksibel untuk menggambarkan keragaman dan kompleksitas data dalam pembelajaran mesin.
3. **Estimasi Parameter dengan Maksimum Likelihood atau Metode Bayesian:** Model ini sering kali menggunakan metode maksimum likelihood atau pendekatan Bayesian untuk mengestimasi parameter dari data yang diamati. Dengan cara ini, mereka dapat menyesuaikan model mereka dengan data yang ada dan menghasilkan prediksi yang lebih akurat.
4. **Penggunaan dalam Klasifikasi dan Regresi:** Model-model ini digunakan untuk klasifikasi, regresi, atau tugas-tugas pembelajaran mesin lainnya. Mereka mampu memberikan prediksi dengan menghasilkan distribusi probabilitas atas output yang mungkin, bukan hanya memberikan label atau nilai tunggal.
5. **Interpretasi yang Lebih Mudah:** Dalam beberapa kasus, model probabilistik dapat memberikan interpretasi yang lebih intuitif atas hasilnya. Mereka memungkinkan kita untuk memahami seberapa yakin model terhadap prediksi yang dibuatnya.
6. **Dalam keseluruhan, Model Probabilistik** dalam Pembelajaran Mesin membantu dalam mengatasi ketidakpastian yang melekat dalam data, memungkinkan model untuk membuat prediksi yang lebih cermat, dan memberikan cara yang lebih terstruktur untuk memahami distribusi data yang kompleks. Ini adalah pendekatan yang kuat dalam konteks di mana informasi probabilistik diperlukan untuk membuat keputusan yang cerdas dan akurat.
7. **Pemodelan Distribusi Data** merujuk pada upaya untuk menggambarkan atau memahami bagaimana data yang diamati tersebar atau didistribusikan di dalam ruang sampel. Ini adalah

konsep yang penting dalam statistika dan pembelajaran mesin karena membantu dalam memahami sifat-sifat data dan mencari model yang cocok untuk menjelaskan data tersebut.

Dalam pemodelan distribusi data, kita mencari fungsi matematis yang paling sesuai untuk menggambarkan sebaran atau distribusi data yang diamati. Fungsi ini sering kali didasarkan pada sejumlah parameter yang kemudian akan diestimasi dari data yang tersedia. Beberapa distribusi probabilitas yang sering digunakan untuk memodelkan data meliputi distribusi normal (Gaussian), distribusi binomial, distribusi Poisson, distribusi eksponensial, dan banyak lagi.

## **F. Aspek Pemodelan distribusi data memiliki**

1. **Deskripsi Distribusi Data:** Melalui pemodelan distribusi, kita bisa mendapatkan gambaran yang jelas tentang sebaran data. Misalnya, jika data terdistribusi normal, kita dapat menggunakan parameter rata-rata dan deviasi standar untuk mendeskripsikan distribusi tersebut.
2. **Prediksi dan Estimasi:** Dengan mengetahui distribusi data, kita dapat membuat prediksi atau estimasi terkait nilai-nilai yang mungkin dari data yang baru. Misalnya, dalam prediksi, jika kita mengetahui distribusi data yang ada, kita dapat membuat perkiraan tentang nilai yang paling mungkin terjadi.
3. **Pemilihan Model yang Sesuai:** Pemodelan distribusi membantu kita memilih model yang paling cocok untuk data yang kita hadapi. Ini membantu dalam pembuatan model yang lebih akurat dan representatif terhadap data yang sebenarnya.
4. **Analisis Statistik Lanjutan:** Distribusi data juga menjadi dasar untuk banyak analisis statistik lanjutan. Misalnya, dalam inferensi statistik, pemodelan distribusi menjadi kunci dalam membuat asumsi tentang distribusi data sampel terhadap populasi yang lebih besar.

Dengan kata lain, pemodelan distribusi data adalah usaha untuk menemukan atau menyesuaikan model matematika yang paling cocok untuk menjelaskan cara data tersebar, sehingga membantu

dalam analisis, prediksi, dan pengambilan keputusan yang berkaitan dengan data tersebut.

## **G. Contoh Pemodelan Distribusi Data**

### **Berikut adalah beberapa contoh Pemodelan Distribusi Data:**

**Distribusi Normal (Gaussian):** Contoh: Misalkan Anda memiliki ketinggian orang-orang dalam sebuah populasi. Jika data ketinggian tersebut terdistribusi secara mendekati kurva normal (bell curve) dengan rata-rata 170 cm dan deviasi standar 10 cm, Anda dapat menggunakan distribusi normal untuk memodelkan sebaran ketinggian tersebut. Dengan model ini, Anda dapat memprediksi seberapa mungkin orang memiliki ketinggian tertentu di dalam populasi berdasarkan karakteristik distribusi tersebut.

**Distribusi Binomial:** Contoh: Bayangkan Anda melakukan serangkaian uji coba di mana setiap uji coba memiliki dua hasil mungkin: sukses atau gagal. Misalnya, Anda melempar koin 10 kali dan mencatat berapa kali hasilnya adalah gambar (heads). Jika Anda ingin memodelkan distribusi jumlah gambar yang mungkin muncul dari 10 lemparan tersebut, Anda bisa menggunakan distribusi binomial. Dengan ini, Anda bisa memprediksi probabilitas munculnya sejumlah gambar tertentu dalam serangkaian lemparan koin.

**Distribusi Poisson:** Contoh: Anda mengamati jumlah kendaraan yang melewati suatu titik dalam satu jam di suatu jalan raya yang jarang dilewati. Anda mencatat rata-rata lima kendaraan per jam. Distribusi Poisson bisa digunakan untuk memodelkan sebaran jumlah kendaraan yang melewati titik tersebut dalam interval waktu tertentu. Dengan model ini, Anda dapat memperkirakan probabilitas munculnya sejumlah kendaraan dalam interval waktu yang telah ditentukan.

**Distribusi Eksponensial:** Contoh: Bayangkan Anda ingin memodelkan waktu antara kedatangan pelanggan ke suatu layanan perbankan. Jika waktu antara kedatangan pelanggan terdistribusi

eksponensial dengan rata-rata 5 menit, Anda dapat menggunakan distribusi eksponensial untuk memodelkan interval waktu antara kedatangan pelanggan ke lokasi tersebut. Dengan ini, Anda dapat membuat perkiraan tentang waktu yang diharapkan untuk kedatangan pelanggan berikutnya.

Contoh-contoh di atas mengilustrasikan cara pemodelan distribusi data digunakan dalam berbagai konteks untuk menggambarkan cara data tersebar, memberikan prediksi, serta membantu dalam pengambilan keputusan berdasarkan karakteristik distribusi tersebut. Berikut adalah contoh matematis dari beberapa distribusi data yang umum digunakan dalam pemodelan statistik:

**Distribusi Normal (Gaussian):** Distribusi Normal didefinisikan oleh fungsi densitas probabilitas (probability density function, PDF)<sup>45</sup>:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Di sini,  $\mu$  adalah nilai rata-rata,  $\sigma$  adalah deviasi standar, dan  $x$  adalah variabel acak yang diukur.

**Distribusi Binomial:** Distribusi Binomial menggambarkan probabilitas  $p$  sukses atau  $1-p$  gagal dalam  $n$  uji coba independen. Fungsi mass probabilitas (probability mass function, PMF) untuk distribusi binomial diberikan oleh:

$$P(X = k) = \binom{n}{k} \cdot p^k \cdot (1 - p)^{n-k}$$

Di sini,  $k$  adalah jumlah sukses yang diharapkan dalam  $n$  uji coba,  $p$  adalah probabilitas sukses dalam satu uji coba, dan  $X$  adalah variabel acak yang menggambarkan jumlah sukses.

---

<sup>4</sup> [https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution)

<sup>5</sup> <https://itl.nist.gov/div898/handbook/eda/section3/eda3661.htm>



**Distribusi Poisson:** Distribusi Poisson menggambarkan jumlah peristiwa yang terjadi dalam Distribusi Poisson adalah distribusi probabilitas diskrit yang menggambarkan jumlah peristiwa yang terjadi dalam suatu interval waktu atau ruang tertentu, ketika peristiwa-peristiwa tersebut terjadi dengan tingkat kejadian yang konstan dan secara independen dari waktu sebelumnya. Distribusi ini sering digunakan dalam berbagai bidang seperti ilmu statistik, matematika, ilmu sosial, dan lainnya untuk memodelkan peristiwa yang jarang terjadi namun memiliki tingkat kejadian yang stabil. Rumus Distribusi Poisson adalah:

$$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$P(X=k)$  adalah probabilitas bahwa terjadi  $k$  peristiwa dalam interval waktu atau ruang yang diberikan.  $e$  adalah konstanta Euler (sekitar 2.71828).

$\lambda$  adalah tingkat kejadian rata-rata per interval waktu atau ruang. Ini bisa dianggap sebagai rata-rata jumlah peristiwa yang terjadi.  $k$  adalah jumlah peristiwa yang ingin dihitung probabilitasnya.  $k!$  adalah faktorial dari  $k$  (produk dari semua bilangan bulat positif kurang dari atau sama dengan  $k$ ).

Misalnya, jika kita memiliki situasi di mana rata-rata jumlah mobil yang melewati suatu jalan dalam satu jam adalah 5, kita dapat menggunakan distribusi Poisson untuk menghitung probabilitas bahwa tepat 3 mobil akan melewati jalan dalam waktu satu jam:

$$P(X = 3) = \frac{e^{-5} 5^3}{3!}$$

Ini akan memberikan probabilitas bahwa tepat 3 mobil akan melewati jalan dalam interval waktu satu jam, berdasarkan asumsi tingkat kejadian rata-rata sebanyak 5 mobil per jam.

Distribusi Poisson berguna dalam memodelkan peristiwa-peristiwa yang jarang terjadi namun memiliki distribusi kejadian yang terukur.

### **Distribusi Eksponensial:**

Distribusi Eksponensial menggambarkan waktu antara peristiwa-peristiwa yang terjadi secara acak. Distribusi Eksponensial adalah distribusi probabilitas yang digunakan untuk memodelkan waktu antara peristiwa-peristiwa yang terjadi secara acak dan independen dalam suatu proses yang memiliki tingkat kejadian konstan. Ini sering digunakan dalam analisis waktu tunggu di berbagai bidang seperti ilmu statistik, ilmu komputer, sistem antrian, dan lainnya. Rumus Distribusi Eksponensial adalah sebagai berikut:

$$f(x; \lambda) = \lambda e^{-\lambda x}$$

$f(x; \lambda)$  adalah fungsi kepadatan probabilitas (PDF) dari variabel acak

$x$  dengan parameter

$\lambda$ , yang menyatakan tingkat kejadian.

$\lambda$  adalah tingkat kejadian yang merupakan invers dari rata-rata peristiwa yang terjadi per unit waktu. Semakin besar nilai

$\lambda$ , semakin cepat peristiwa-peristiwa terjadi.

$e$  adalah konstanta Euler (sekitar 2.71828).

$x$  adalah waktu tunggu atau interval waktu antara peristiwa-peristiwa.

Misalnya, jika kita ingin menghitung probabilitas bahwa waktu antara dua kejadian (misalnya, kedatangan dua kendaraan pada suatu titik dalam sistem transportasi) berada pada interval waktu tertentu, kita dapat menggunakan distribusi Eksponensial. Jika tingkat kedatangan rata-rata kendaraan adalah 4 per jam ( $\lambda=4$ ), maka probabilitas bahwa waktu antara kedatangan dua kendaraan adalah lebih dari 15 menit ( $x > \frac{15}{60}$  jam) adalah:

$$P\left(x, \frac{15}{60}\right) = \int_{\frac{15}{60}}^{\infty} \lambda e^{-\lambda x} dx$$

Distribusi Eksponensial juga sering digunakan dalam model antrian untuk memprediksi waktu tunggu dalam antrian atau interval antara kedatangan pelanggan dalam sistem layanan. Ini membantu dalam analisis kinerja sistem di mana waktu antara peristiwa memegang peranan penting.

## **BAB VI**

# **LATENT DIRICHLET ALLOCATION (LDA)**

Dalam dunia yang semakin dipenuhi oleh volume besar informasi teks, pengelolaan dan pemahaman terhadap konten tersebut telah menjadi tantangan yang semakin kompleks. Salah satu alat yang paling penting dalam menganalisis data teks secara menyeluruh adalah Model Latent Dirichlet Allocation (LDA) (Blei et al., 2003). Konsep ini, yang diadaptasi dari bidang statistik dan pembelajaran mesin, memungkinkan kita untuk mengurai struktur tersembunyi dari dokumen-dokumen yang kompleks, mengidentifikasi pola-pola yang tak terlihat pada pandangan pertama. Sejarah dan latar belakang Model Latent Dirichlet Allocation (LDA) berasal dari dunia ilmu komputer, statistik, dan pengolahan bahasa alami. Model ini diperkenalkan pertama kali oleh David Blei, Andrew Ng, dan Michael Jordan pada tahun 2003 melalui makalah penelitian yang diterbitkan dalam jurnal ilmiah "Journal of Machine Learning Research".

Latar belakang LDA berakar dari upaya untuk menemukan cara efektif untuk mengatasi kompleksitas dalam analisis teks. Sebelum LDA, memahami dan mengelompokkan dokumen-dokumen berdasarkan topik atau pola yang tersembunyi dalam jumlah yang besar merupakan tantangan besar. LDA diciptakan sebagai jawaban untuk mengatasi masalah ini, dengan tujuan memberikan metode yang lebih sistematis dan terstruktur untuk mengekstraksi topik tersembunyi dari kumpulan dokumen yang besar.

Pada dasarnya, LDA diilhami oleh konsep tentang bagaimana dokumen-dokumen terbentuk. Model ini mengasumsikan bahwa dokumen-dokumen dibangun dari sejumlah topik yang tersembunyi, dan setiap kata dalam dokumen tersebut berasal dari salah satu dari topik-topik ini. LDA menggunakan pendekatan probabilistik untuk mengekstraksi distribusi topik dari kumpulan dokumen dan mengidentifikasi pola yang mendasarinya.

Sejak diperkenalkan, LDA telah menjadi salah satu alat yang sangat populer dalam analisis teks, pengelompokan dokumen, sistem rekomendasi, dan pemrosesan bahasa alami. Penggunaannya telah meluas di berbagai bidang seperti ilmu sosial, ekonomi, biomedis, dan lainnya, karena kemampuannya dalam mengurai dan memahami konten teks yang kompleks menjadi topik-topik yang lebih terdefinisi.

Latent Dirichlet Allocation (LDA) adalah model topik probabilistik yang sangat digunakan dalam pemrosesan bahasa alami (NLP) karena kemampuannya untuk mengidentifikasi struktur semantik dalam kumpulan teks besar. LDA beroperasi dengan mengasumsikan bahwa setiap dokumen adalah campuran dari sejumlah topik, dan setiap topik diwakili sebagai distribusi atas kata-kata. Hal ini memungkinkan LDA untuk mengekstraksi dan memahami topik-topik yang mendasari dalam kumpulan data teks tanpa perlu label atau anotasi manual, membuatnya sangat berguna untuk berbagai aplikasi analisis teks.

Pentingnya LDA dalam NLP sangat signifikan karena kemampuannya untuk menangani masalah skala besar, seperti yang ditemukan dalam analisis Big Data. Model ini membantu 'membuka' dan membuat koneksi percakapan laten yang sebelumnya tidak terlihat dalam korpus teks yang luas, seperti profil, thread diskusi, forum, dan media sosial lainnya. LDA membantu dalam mengidentifikasi hubungan yang belum diketahui sebelumnya dan menyediakan wawasan yang lebih dalam tentang struktur semantik data teks (Gross & Murthy, 2014).

Dalam konteks NLP, LDA sering digunakan untuk meningkatkan aplikasi seperti klasifikasi dokumen, pengelompokan teks, dan sistem rekomendasi. Model ini menawarkan kerangka kerja yang kuat untuk memahami dan mengelola variabilitas semantik dan sintaktik dalam teks. Dengan memetakan dokumen ke dalam ruang topik, LDA memfasilitasi pengurangan dimensi yang efektif dan interpretasi semantik yang kaya, yang sangat penting dalam tugas pemahaman teks dan pengambilan informasi (Wang et al.,

2012). Selain itu, variasi LDA, seperti LDA semi-supervised, telah dikembangkan untuk menggabungkan pengetahuan yang diawasi ke dalam prosedur pembelajaran, memungkinkan penggunaan label terawasi untuk memandu pemodelan topik dan meningkatkan akurasi klasifikasi dokumen. Ini menunjukkan fleksibilitas dan kemampuan adaptasi LDA untuk memenuhi kebutuhan spesifik dari berbagai tugas NLP (Wang et al., 2012).

Secara keseluruhan, LDA merupakan alat yang sangat berharga dalam kotak alat NLP, memberikan wawasan mendalam tentang struktur semantik yang kompleks dari teks dan memfasilitasi pengembangan aplikasi pemrosesan teks yang canggih.

## A. Prinsip Kerja LDA

Model Latent Dirichlet Allocation (LDA) bekerja dengan cara mengasumsikan bahwa setiap dokumen dalam kumpulan dokumen dibentuk oleh kombinasi dari beberapa topik, dan setiap kata dalam dokumen berasal dari salah satu dari topik-topik tersebut. Prinsip kerja LDA secara rinci dapat dijabarkan sebagai berikut:

**Inisialisasi Awal:** LDA dimulai dengan tahap inisialisasi di mana setiap kata dalam setiap dokumen ditugaskan secara acak ke salah satu dari sejumlah topik yang telah ditentukan. Awalnya, distribusi kata dalam dokumen ditetapkan secara acak.

**Iterasi Estimasi:** Model melakukan iterasi untuk menyesuaikan distribusi topik di setiap dokumen dan distribusi kata di setiap topik. Dalam setiap iterasi, LDA mencoba untuk memperbaiki penugasan kata-kata ke topik-topik berdasarkan dua hal utama:

1. Perhitungan Proporsi Topik dalam Dokumen: Model memperkirakan seberapa banyak setiap topik mempengaruhi setiap dokumen. Ini dilakukan dengan menghitung proporsi atau distribusi probabilitas dari setiap topik dalam setiap dokumen.
2. Perhitungan Proporsi Kata dalam Topik: LDA juga memperkirakan seberapa banyak setiap kata terkait dengan setiap topik. Ini dilakukan dengan menghitung proporsi atau distribusi probabilitas dari setiap kata dalam setiap topik.

**Update Parameter:** Setelah iterasi yang berulang, model memperbarui parameter-parameternya untuk memperbaiki estimasi proporsi kata dalam topik dan proporsi topik dalam dokumen.

**Penentuan Topik:** Setelah proses iterasi selesai, LDA menghasilkan distribusi topik yang diperkirakan untuk setiap dokumen dan distribusi kata yang diperkirakan untuk setiap topik. Dengan hasil ini, kita dapat melihat topik-topik yang mendominasi setiap dokumen dan kata-kata yang paling terkait dengan masing-masing topik.

Prinsip utama di balik LDA adalah bagaimana model mencoba untuk memperbaiki estimasi awal terkait dengan bagaimana kata-kata terdistribusi di antara topik-topik dan bagaimana topik-topik didistribusikan di antara dokumen-dokumen. Tujuannya adalah untuk menemukan pola yang tersembunyi dalam dokumen-dokumen dan menghasilkan representasi yang lebih terstruktur dan informatif tentang topik-topik yang ada dalam kumpulan dokumen tersebut. Implementasi matematis dari Model Latent Dirichlet Allocation (LDA) melibatkan langkah-langkah yang kompleks dalam memodelkan distribusi kata-kata di dalam dokumen dan distribusi topik di dalam kumpulan dokumen. Di bawah ini adalah detail langkah-langkah implementasi matematis LDA:

## **B. Pembentukan Model:**

**Variabel Laten:** LDA melibatkan variabel laten (tersembunyi), termasuk variabel topik dan variabel distribusi topik pada dokumen-dokumen.

**Parameter Model:** Parameter yang diperlukan meliputi jumlah topik yang diinginkan ( $K$ ), distribusi Dirichlet untuk topik dalam dokumen ( $\alpha$ ), dan distribusi Dirichlet untuk kata dalam topik ( $\beta$ ).  
**Representasi Dokumen:** Dokumen direpresentasikan dalam bentuk matriks di mana setiap baris mewakili sebuah dokumen, dan setiap kolom mewakili jumlah kata dalam kosa kata yang

digunakan. Nilai di dalam matriks ini mewakili frekuensi kemunculan kata dalam dokumen tersebut.

**Proses Estimasi dan Iterasi:** Iterasi dimulai dengan menginisialisasi secara acak nilai-nilai awal untuk variabel tersembunyi (topik dari kata-kata dalam dokumen).

Proses perhitungan dilakukan berulang kali untuk memperbaiki estimasi variabel laten. Langkah-langkah ini melibatkan perhitungan proporsi topik dalam dokumen dan proporsi kata dalam topik.

**Metode Variational Inference atau Gibbs Sampling:** Metode ini sering digunakan dalam LDA untuk mendekati distribusi posterior dari variabel tersembunyi. Dalam variational inference, tujuannya adalah untuk mendekati distribusi posterior dengan memilih distribusi yang paling dekat secara matematis. Gibbs sampling, metode lain yang digunakan, melibatkan pengambilan sampel acak dari distribusi probabilitas yang diinginkan.

### **C. Penyesuaian Parameter**

Selama iterasi, nilai-nilai parameter model (seperti  $\alpha$  dan  $\beta$ ) disesuaikan untuk memperbaiki estimasi distribusi topik dan kata-kata di dalam dokumen.

**Evaluasi dan Output:** Setelah iterasi yang cukup banyak, model LDA menghasilkan distribusi topik untuk setiap dokumen dan distribusi kata untuk setiap topik. Hasil ini memberikan representasi yang lebih baik tentang topik-topik yang mendominasi dokumen dan kata-kata yang paling terkait dengan masing-masing topik.

Implementasi matematis LDA melibatkan perhitungan probabilistik yang kompleks, termasuk penggunaan distribusi Dirichlet dan perhitungan untuk memperbaiki estimasi variabel tersembunyi. Langkah-langkah ini memungkinkan model untuk mengekstraksi informasi tersembunyi dari kumpulan dokumen secara efisien.



## D. Persamaan dalam Model LDA

Rumus Umum LDA: LDA dapat direpresentasikan sebagai model generatif probabilitas yang mencakup beberapa variabel laten. Untuk setiap kata dalam dokumen, ada dua variabel laten utama yang penting dalam LDA, yaitu variabel topik  $z$  dan variabel distribusi topik  $\phi$ .

$\theta$ : Distribusi topik dalam dokumen.

$z$ : Topik yang dipilih untuk setiap kata dalam dokumen.

$\beta$ : Distribusi kata dalam topik.

$w$ : Kata yang diamati dalam dokumen.

Model LDA direpresentasikan dengan rumus umum sebagai berikut:

$$P(\theta, z, \beta | w) \\ = P(\theta) \cdot \prod_{d=1}^D P(\theta_d) \cdot \prod_{n=1}^N (Z_{d,n} | \theta_d \cdot P(W w_{d,n} | \beta_{z_{d,n}})$$

Dengan  $D$  adalah jumlah dokumen dalam kumpulan dokumen,  $N$  adalah jumlah kata dalam dokumen, dan  $P(\theta)$  serta  $P(\beta)$  adalah distribusi prior dari variabel  $\theta$  dan  $\beta$ , masing-masing.

## E. Perhitungan Distribusi Posterior

Untuk mengestimasi distribusi posterior dari variabel laten ( $\theta$  dan  $\beta$ ) dalam LDA, diperlukan metode seperti variational inference atau Gibbs sampling. Metode-metode ini digunakan untuk mendekati distribusi posterior dari variabel laten, yang tidak dapat dihitung secara langsung. Rumus umum LDA memberikan kerangka kerja untuk memahami bagaimana dokumen dibangun dari kombinasi topik dan bagaimana kata-kata dalam dokumen berasal dari topik-topik tertentu. Langkah-langkah selanjutnya dalam implementasi LDA melibatkan perhitungan untuk mendekati distribusi posterior dari variabel laten ini.

Tentu, berikut adalah rumus matematis yang mendasari Model Latent Dirichlet Allocation (LDA):

Representasi Dokumen:

D: Jumlah dokumen dalam kumpulan dokumen.

Nd: Jumlah kata dalam dokumen d.

V: Jumlah kata unik dalam kosa kata.

wd,n: Kata ke-n dalam dokumen d.

Parameter Model:

K: Jumlah topik yang diinginkan.

$\alpha$ : Parameter distribusi Dirichlet untuk distribusi topik dalam dokumen.

$\beta$ : Parameter distribusi Dirichlet untuk distribusi kata dalam topik.

Variabel Tersembunyi:

zd,n: Topik yang diatribusikan untuk kata ke-n dalam dokumen d.

Rumus Estimasi LDA:

a. Distribusi topik dalam dokumen:

$$P(\theta_d | \alpha = \text{Dirichlet}(\theta_d | \alpha))$$

$$P(\theta_d) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma \alpha_i} \prod_{i=1}^K \theta_{d,i}^{\alpha_i - 1}$$

b. Distribusi kata dalam topik

$$P(\phi_k | \beta = \text{Dirichlet}(\phi_k | \beta))$$

$$P(\phi_k) = \frac{\Gamma(\sum_{i=1}^V \beta_i)}{\prod_{i=1}^V \Gamma \beta_i} \prod_{i=1}^V \phi_{k,i}^{\beta_i - 1}$$

c. Probabilitas word assignment

$$P(w_{d,n} | \theta_d, \phi_{z_{d,n}}) = \theta_{d,z_{d,n}} \phi_{z_{d,n},w_{d,n}}$$

Rumus-rumus di atas menggambarkan cara LDA memodelkan distribusi topik dalam dokumen, distribusi kata dalam topik, dan probabilitas penugasan kata ke topik dalam dokumen. Metode variational inference atau Gibbs sampling sering digunakan untuk mendekati atau menemukan solusi numerik dari model LDA ini. Tentu, berikut adalah rumus matematis dari Model Latent Dirichlet Allocation (LDA):

Notasi dan Variabel yang Digunakan:

$D$  = Jumlah dokumen dalam kumpulan dokumen

$N$  = Jumlah kata dalam dokumen  $d$

$K$  = Jumlah topik yang diinginkan

$V$  = Jumlah kata dalam kosa kata

$w_{d,n}$  = Kata ke- $n$  dalam dokumen ke- $d$

$z_{d,n}$  = Topik yang ditugaskan kepada kata ke- $n$  dalam dokumen ke- $d$

$\alpha$  = Parameter distribusi Dirichlet untuk distribusi topik dalam dokumen

$\beta$  = Parameter distribusi Dirichlet untuk distribusi kata dalam topik

## F. Rumus-rumus LDA

Representasi Distribusi Topik dalam Dokumen:  $\theta_{d,k} \sim \text{Dirichlet}(\alpha)$

Representasi Distribusi Kata dalam Topik:  $\phi_{k,v} \sim \text{Dirichlet}(\beta)$

Pembentukan Variabel Laten:  $w_{d,n} \sim \text{Multinomial}(\phi_{z_{d,n}})$

Distribusi Posterior untuk  $\theta$  dan  $\phi$ :

$$p(\theta_d | w, \alpha) = \frac{p(\theta_d) \times p(w | \theta_d, \phi) \times p(\phi | \beta)}{p(w | \alpha, \beta)}$$

Estimasi Distribusi Posterior:

$$p(\theta_d | w, \alpha) = \frac{\text{Dirichlet}(\alpha + \sum_{n=1}^N \text{Count}(w_{d,n}))}{\text{Dirichlet}(\alpha) + N}$$

$$p(\phi_k | w, \beta) = \frac{\text{Dirichlet}(\beta + \sum_{d=1}^D \sum_{n=1}^N \text{Count}(w_{d,n}) \times (z_{d,n}=k))}{\text{Dirichlet}(\beta) + \sum_{d=1}^D \sum_{n=1}^N \text{Count}(w_{d,n})}$$

rumus xx

Persamaan xx menunjukkan adalah rumus posterior dari distribusi probabilitas suatu topik ( $\phi_k$ ) dalam model Latent Dirichlet Allocation (LDA). Dalam persamaan ini:

$p(\phi_k | w, \beta)$  adalah probabilitas posterior dari topik  $\phi_k$ , dengan diberikan dokumen-dokumen ( $w$ ) dalam korpus dan parameter hyperparameter  $\beta$ .

Dirichlet  $(\beta + \sum_{(d=1)}^D \sum_{(n=1)}^N \text{Count}(w_{(d,n)}))$   $x(z_{(d,n)}=k)$  adalah distribusi Dirichlet dengan parameter  $\beta$  yang diubah dengan jumlah kata dalam dokumen yang terasosiasi dengan topik  $k$  ( $\phi_k$ ).

Dirichlet( $\beta$ ) adalah distribusi Dirichlet dengan parameter  $\beta$ .  $\sum_{(d=1)}^D \sum_{(n=1)}^N \text{Count}(w_{(d,n)})$   $x(z_{(d,n)}=k)$  mewakili jumlah kemunculan kata dalam dokumen yang dikaitkan dengan topik  $k$ .  $\sum_{(d=1)}^D \sum_{(n=1)}^N \text{Count}(w_{(d,n)})$  mewakili jumlah total kata dalam semua dokumen.

Secara intuitif, rumus tersebut menggambarkan bagaimana probabilitas distribusi topik tertentu dihitung berdasarkan jumlah kemunculan kata dalam dokumen yang terkait dengan topik tersebut, dibandingkan dengan jumlah total kata dalam semua dokumen, dan diperbarui dengan parameter  $\beta$  yang merupakan prior distribusi topik. Dalam LDA, tujuannya adalah untuk mengetahui distribusi topik kata ( $\phi_k$ ) dan distribusi topik dokumen ( $\theta_d$ ), dan rumus ini membantu dalam memperbarui estimasi distribusi topik kata berdasarkan dokumen yang diamati dalam korpus yang diberikan.

## G. Proses Model LDA

Secara singkat dapat dijelaskan sebagai berikut:

Model Latent Dirichlet Allocation (LDA) adalah model generatif yang menggunakan konsep probabilitas untuk menggambarkan hubungan antara variabel tersembunyi (topik) dan variabel pengamatan (kata-kata dalam dokumen). Berikut adalah representasi matematis dari Model LDA:

## H. Variabel Tersembunyi:

### 1. Distribusi Topik Dokumen $\theta$

$$p(\theta_i|\alpha) = \text{Dir}(\theta_i|\alpha)$$

$\theta_i$  adalah distribusi dari topik-topik dalam dokumen  $i$ , diambil dari distribusi Dirichlet dengan parameter  $\alpha$ .

## 2. Distribusi Kata dalam Topik $\beta$

$$p(\beta_k|\eta) = \text{Dir}(\beta_k|\eta)$$

$\beta_k$  adalah distribusi dari kata-kata dalam topik  $k$ , diambil dari distribusi Dirichlet dengan parameter  $\eta$ .

### Proses Generatif:

- Untuk setiap dokumen  $i$ :
  - $\theta_i \sim \text{Dir}(\alpha)$
- Untuk setiap kata ke- $j$  dalam dokumen  $i$ :
  - $z_{ij} \sim \text{Multinomial}(\theta_i)$
  - $w_{ij} \sim \text{Multinomial}(\beta_{z_{ij}})$
- Penjelasan:
- $\alpha$  adalah parameter prior untuk distribusi topik dokumen.
- $\eta$  adalah parameter prior untuk distribusi kata dalam topik.
- $\theta_i$  adalah distribusi dari topik-topik dalam dokumen  $i$ .
- $\beta_k$  adalah distribusi dari kata-kata dalam topik  $k$ .
- $z_{ij}$  adalah variabel tersembunyi yang menunjukkan topik yang diatribusikan ke kata ke- $j$  dalam dokumen  $i$ .
- $w_{ij}$  adalah kata yang diamati ke- $j$  dalam dokumen  $i$ .

Model LDA melakukan proses generatif untuk menghasilkan dokumen-dokumen dengan cara mengambil distribusi topik dari distribusi Dirichlet untuk setiap dokumen, kemudian memilih topik dari distribusi topik dokumen untuk setiap kata dalam dokumen tersebut, dan akhirnya memilih kata dari distribusi kata dalam topik yang terkait dengan topik yang telah dipilih sebelumnya.

## I. Implementasi LDA

```
import zipfile
import pandas as pd
import os
```

`import zipfile`: Ini mengimpor modul `zipfile`, yang memungkinkan Anda untuk bekerja dengan file zip di Python. Dengan menggunakan modul ini, Anda dapat mengekstrak file dari arsip zip, membuat file zip, dan melakukan operasi terkait file zip lainnya. `import pandas as pd`: Ini mengimpor modul `pandas` dengan alias `pd`. `Pandas` adalah pustaka yang sangat populer dalam Python untuk analisis data. Dengan menggunakan `pd` sebagai alias, Anda bisa mengakses fungsi-fungsi dan objek-objek dari pustaka `pandas` dengan menggunakan `pd` sebagai awalan. `import os`: Ini mengimpor modul `os`, yang memberikan fungsionalitas sistem operasi, seperti interaksi dengan sistem file, mengelola variabel lingkungan, dan melakukan operasi terkait sistem operasi lainnya di dalam program Python Anda.

```
from google.colab import drive
drive.mount('/content/drive')
```

Perintah `from google.colab import drive` digunakan dalam lingkungan Google Colab, yang merupakan lingkungan pengembangan berbasis cloud dari Google yang memungkinkan Anda untuk menulis dan mengeksekusi kode Python di browser. Perintah ini mengimpor fungsi `drive` dari modul `google.colab`. Fungsi `drive` ini digunakan untuk melakukan `mount` atau menghubungkan Google Drive ke sesi Colab Anda. Dengan cara ini, Anda bisa mengakses file yang ada di Google Drive dari lingkungan Colab untuk membaca, menulis, atau melakukan operasi lainnya pada file tersebut melalui kode Python.

Pada perintah `drive.mount('/content/drive')`, `drive.mount()` adalah panggilan fungsi yang memicu proses `mount` Google Drive ke sesi Colab. Argument `'/content/drive'` adalah path atau lokasi di mana Google Drive akan di-mount di dalam lingkungan Colab. Setelah menjalankan perintah ini, Colab akan meminta autentikasi dengan akun Google Anda dan memberikan kode untuk autentikasi, yang perlu di-copy-paste untuk mengotorisasi akses ke Google Drive Anda. Setelah otorisasi berhasil, Google Drive akan di-mount ke path yang telah ditentukan, dalam contoh ini ke `'/content/drive'`.

```
#papers = pd.read_csv('drive/My
Drive/dataset/fintechP2P/2023/16ribu.csv') #lokasi file
papers = pd.read_csv('drive/My
Drive/dataset/fintechP2P/2023/dataset/12-feb-2023masterurut-
p2p.csv') #lokasi file
```

Perintah ini adalah contoh penggunaan dari pustaka pandas di Python untuk membaca sebuah file CSV ke dalam variabel papers. Dalam kode yang Anda berikan:

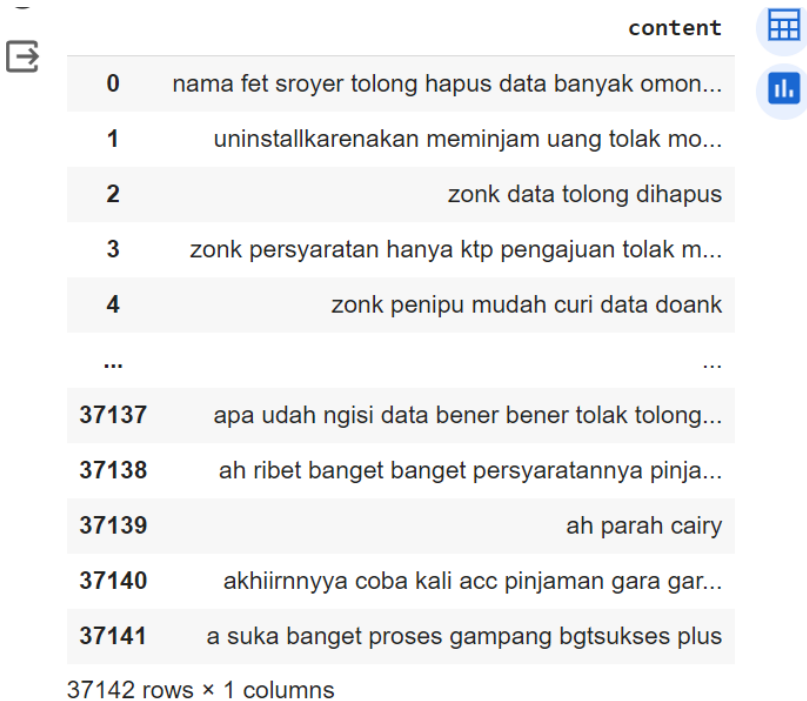
`pd.read_csv()` adalah fungsi dari pustaka pandas yang digunakan untuk membaca file CSV. `'drive/My Drive/dataset/fintechP2P/2023/dataset/12-feb-2023masterurut-p2p.csv'` adalah path atau lokasi dari file CSV yang akan dibaca. Jadi, perintah membaca file CSV yang terletak pada path tersebut dan menyimpannya ke dalam variabel papers. Setelah eksekusi perintah ini, data dari file CSV tersebut akan dimuat ke dalam variabel papers, yang kemudian bisa digunakan untuk analisis lebih lanjut atau manipulasi data menggunakan pustaka pandas.

`papers`

Papers adalah sebuah variabel yang digunakan untuk menyimpan data yang dibaca dari file CSV dengan menggunakan pustaka pandas di Python. Dalam konteks ini, papers mungkin berisi kumpulan data yang terdapat dalam file CSV yang telah dibaca menggunakan fungsi `pd.read_csv()`.

Variabel ini bisa berupa DataFrame, struktur data yang sangat berguna dari pustaka pandas. DataFrame memungkinkan untuk menyimpan data dalam bentuk tabel dengan baris dan kolom, mirip dengan spreadsheet. Setiap kolom dalam DataFrame mungkin merepresentasikan jenis data atau atribut tertentu, sedangkan setiap baris mungkin merepresentasikan entri atau contoh dari data tersebut.

Dengan menggunakan variabel `papers`, Anda bisa melakukan berbagai operasi analisis data, seperti manipulasi data, pengolahan statistik, visualisasi, dan banyak lagi, karena `papers` sekarang berisi data dari file CSV yang telah dimuat menggunakan pustaka `pandas`.



	content
0	nama fet sroyer tolong hapus data banyak omon...
1	uninstallkarenakan meminjam uang tolak mo...
2	zonk data tolong dihapus
3	zonk persyaratan hanya ktp pengajuan tolak m...
4	zonk penipu mudah curi data doank
...	...
37137	apa udah ngisi data bener bener tolak tolong...
37138	ah ribet banget banget persyaratannya pinja...
37139	ah parah cairy
37140	akhiirnnya coba kali acc pinjaman gara gar...
37141	a suka banget proses gampang bgtsukses plus

37142 rows × 1 columns

Variabel `papers` merupakan sebuah `DataFrame` yang berisi data teks atau komentar-komentar terkait dengan informasi tertentu. `DataFrame` ini memiliki satu kolom dengan nama `'content'` yang berisi teks komentar.

Dari potongan data yang Anda tunjukkan, terdapat 37143 baris (0 sampai 37142) dan 1 kolom (`'content'`). Isi dari kolom `'content'` ini tampaknya berupa komentar-komentar atau teks yang berkaitan dengan suatu topik, mungkin terkait dengan pendapat atau ulasan terhadap suatu layanan atau produk.



Contoh beberapa baris dari data yang tersimpan dalam variabel `papers`:

Baris ke-0: "nama fet sroyer tolong hapus data banyak omon..."

Baris ke-1: "uninstallkarenakan meminjam uang tolak mo..."

Baris ke-2: "zonk data tolong dihapus"

Baris ke-3: "zonk persyaratan hanya ktp pengajuan tolak m..."

Baris ke-4: "zonk penipu mudah curi data doank"

Setiap baris berisi komentar atau informasi yang mungkin dapat dianalisis lebih lanjut, misalnya, untuk mengidentifikasi sentimen atau pola-pola tertentu dalam teks tersebut menggunakan teknik pemrosesan bahasa alami atau untuk melakukan analisis sentimen terhadap pendapat-pendapat tersebut.

```
papers['word_count'] = papers['content'].str.split().map(len)
```

Perintah ini menambahkan kolom baru ke dalam DataFrame `papers` dengan nama `'word_count'`.

Mari kita bahas lebih rinci: `papers['content']`: Merujuk pada kolom `'content'` dalam DataFrame `papers`. Ini adalah kolom yang berisi teks atau komentar-komentar.

`str.split()`: Ini adalah metode dari objek Series di pandas yang digunakan untuk membagi setiap teks dalam kolom `'content'` menjadi kata-kata (dengan menggunakan spasi sebagai pemisah). Hasilnya adalah daftar kata-kata untuk setiap teks. `.map(len)`: Setelah kata-kata dipisahkan untuk setiap teks dalam kolom `'content'`, `map(len)` diaplikasikan pada setiap daftar kata-kata. Fungsinya adalah untuk menghitung panjang dari setiap daftar kata-kata, yang pada dasarnya adalah jumlah kata dalam setiap teks. Dengan menggunakan `map(len)`, dihitunglah panjang setiap daftar kata-kata, yang sebenarnya adalah jumlah kata dalam setiap baris.

`papers['word_count']`: Ini menugaskan hasil dari perhitungan jumlah kata ke dalam kolom baru yang bernama 'word\_count' di dalam DataFrame `papers`. Dengan demikian, setelah eksekusi perintah ini, `papers` akan memiliki kolom tambahan yang menampilkan jumlah kata dalam setiap teks yang ada di dalam kolom 'content'.



	content	word_count
0	nama fet sroyer tolong hapus data banyak omon...	9
1	uninstallkarenakan meminjam uang tolak mo...	11
2	zonk data tolong dihapus	4
3	zonk persyaratan hanya ktp pengajuan tolak m...	18
4	zonk penipu mudah curi data doank	6
...	...	...
37137	apa udah ngisi data bener bener tolak tolong...	22
37138	ah ribet banget banget persyaratannya pinja...	12
37139	ah parah cairy	3
37140	akhiirnnya coba kali acc pinjaman gara gar...	18
37141	a suka banget proses gampang bgtsukses plus	7

37142 rows x 2 columns

Tabel ini menunjukkan hasil dari penghitungan jumlah kata dalam setiap teks yang terdapat dalam kolom 'content' dari DataFrame `papers`. Kolom baru yang diberi nama 'word\_count' menampilkan jumlah kata dalam setiap baris teks yang sesuai.

Contohnya, untuk beberapa baris tertentu:

1. Baris pertama ('nama fet sroyer tolong hapus data banyak omon...'): Memiliki 9 kata.
2. Baris kedua ('uninstallkarenakan meminjam uang tolak mo...'): Memiliki 11 kata.
3. Baris ketiga ('zonk data tolong dihapus'): Memiliki 4 kata.
4. Baris keempat ('zonk persyaratan hanya ktp pengajuan tolak m...'): Memiliki 18 kata.

5. Baris kelima ('zonk penipu mudah curi data doank'): Memiliki 6 kata.
6. Baris terakhir ('a suka banget proses gampang bgtsukses plus'): Memiliki 7 kata.

Jadi, kolom 'word\_count' ini memberikan informasi tentang jumlah kata yang ada dalam setiap baris teks yang ada di dalam kolom 'content' DataFrame papers. Informasi ini bisa berguna untuk analisis statistik atau pemahaman lebih lanjut tentang panjang atau kompleksitas teks dalam dataset tersebut.

```
papers = papers[papers['word_count'] > 4]
```

Perintah `papers = papers[papers['word_count'] > 4]` adalah contoh dari penggunaan filter di dalam Python dengan menggunakan pustaka pandas untuk DataFrame papers. Mari kita bahas langkah-langkahnya:

1. `papers['word_count'] > 4`: Ini adalah sebuah kondisi yang diterapkan pada kolom 'word\_count' di dalam DataFrame papers. Kondisi ini mengevaluasi setiap baris dalam kolom 'word\_count' dan menghasilkan nilai True jika nilai dalam baris tersebut lebih besar dari 4, dan False jika tidak.
2. `papers[papers['word_count'] > 4]`: Ini adalah teknik filter DataFrame di dalam pandas. Menggunakan kondisi di atas, perintah ini memilih hanya baris-baris dari DataFrame papers di mana kondisi `papers['word_count'] > 4` bernilai True. Dengan kata lain, hanya baris-baris yang memiliki jumlah kata lebih dari 4 yang akan tetap ada dalam DataFrame yang baru. DataFrame yang dihasilkan akan berisi hanya baris-baris tersebut.

Dengan menggunakan perintah ini, DataFrame papers diubah sedemikian rupa sehingga hanya menyertakan baris-baris di mana jumlah kata dalam teks (diwakili oleh kolom 'word\_count') lebih dari 4. Ini memungkinkan untuk memfilter data berdasarkan kriteria tersebut, membuang baris-baris yang tidak memenuhi syarat tersebut.

```
#menghapus duplikasi data
papers.drop_duplicates(subset="content",keep = False, inplace =
True)
```

Perintah `papers.drop_duplicates(subset="content", keep=False, inplace=True)` digunakan untuk menghapus baris-baris duplikat dari DataFrame `papers` berdasarkan kolom 'content'.

Mari kita bahas detailnya:

1. `papers`: Merujuk pada DataFrame yang sedang dioperasikan.
2. `.drop_duplicates()`: Ini adalah metode dari pandas yang digunakan untuk menghapus baris-baris yang merupakan duplikat dari DataFrame.
3. `subset="content"`: Parameter subset menentukan kolom mana yang akan diperiksa untuk mendeteksi duplikat. Di sini, kita menggunakan kolom 'content', yang berisi teks atau komentar-komentar.
4. `keep=False`: Parameter keep menentukan bagaimana mempertahankan hasil penghapusan. Nilai False berarti semua baris yang memiliki nilai yang sama di kolom yang ditentukan akan dihapus, termasuk baris pertama dan yang kedua (semua duplikat).
5. `inplace=True`: Parameter inplace menentukan apakah perubahan akan diterapkan pada DataFrame itu sendiri atau apakah hasilnya akan disimpan dalam DataFrame baru. Dengan nilai True, perubahan akan diterapkan pada `papers` tanpa membuat DataFrame baru.

Jadi, setelah eksekusi perintah ini, baris-baris yang memiliki nilai yang sama dalam kolom 'content' akan dihapus dari DataFrame `papers`. Hal ini membantu memastikan bahwa setiap baris dalam DataFrame tersebut memiliki nilai yang unik dalam kolom 'content'. Dari proses menghapus data yang duplikasi dan data yang digunakan yang lebih dari 4 hata maka dihasilkan data sebagai berikut:

Tabel yang Anda sertakan menunjukkan DataFrame setelah operasi penghapusan duplikat dan setelah melakukan filter untuk baris-baris di mana jumlah kata (kolom 'word\_count') lebih besar

dari 4. Kolom 'content' berisi teks atau komentar-komentar, sementara kolom 'word\_count' berisi jumlah kata dalam teks tersebut.

Contohnya:

1. Baris pertama ('nama fet sroyer tolong hapus data banyak omon...') memiliki 9 kata.
2. Baris kedua ('uninstallkarenakan meminjam uang tolak mo...') memiliki 11 kata.
3. Baris ketiga ('zonk persyaratan hanya ktp pengajuan tolak m...') memiliki 18 kata.
4. Baris keempat ('zonk penipu mudah curi data doank') memiliki 6 kata.
5. Baris kelima ('zonk kali repot membayar tanggal mei jatuh ...') memiliki 39 kata.
6. Baris terakhir ('a suka banget proses gampang bgtsukses plus') memiliki 7 kata.

Tabel tersebut menampilkan baris-baris unik (tanpa duplikat), di mana setiap baris memiliki jumlah kata lebih besar dari 4, seperti yang telah dijelaskan sebelumnya. Jumlah total baris dalam DataFrame yang ditampilkan setelah operasi filter tersebut adalah 29505.

```
papers.to_csv('./drive/My
Drive/dataset/fintechP2P/2023/20februari-bersih.csv',
index=False)
```

Perintah diatas merupakan sebuah perintah dalam Python menggunakan pustaka pandas untuk menyimpan DataFrame papers ke dalam format file CSV.

1. papers: Merujuk pada DataFrame yang ingin disimpan.
2. .to\_csv(): Ini adalah metode dari pustaka pandas yang digunakan untuk menyimpan DataFrame ke dalam format file CSV.
3. './drive/My Drive/dataset/fintechP2P/2023/20februari-bersih.csv': Ini adalah path atau lokasi file di mana DataFrame papers akan disimpan sebagai file CSV. Dalam

kasus ini, file tersebut akan disimpan di lokasi yang ditentukan dengan nama file '20februari-bersih.csv'.

4. `index=False`: Parameter `index` digunakan untuk menentukan apakah indeks dari DataFrame juga akan disimpan sebagai kolom dalam file CSV. Dengan nilai `False`, indeks tidak akan disertakan dalam file CSV yang dihasilkan.

Jadi, perintah ini akan menyimpan DataFrame `papers` ke dalam file CSV dengan nama '20februari-bersih.csv' di lokasi yang ditentukan. File CSV yang dihasilkan akan berisi data dari DataFrame `papers`, dan indeks DataFrame tidak akan disertakan dalam file CSV tersebut.

## Tahap 2: Data Cleaning

```
papers = papers.sample(10000)
```

Perintah `papers = papers.sample(10000)` adalah perintah yang digunakan pada DataFrame dalam pustaka `pandas` di Python untuk mengambil sampel acak sejumlah 10.000 baris dari DataFrame `papers`.

`papers`: Merujuk pada DataFrame yang sedang dioperasikan.  
`.sample()`: Ini adalah metode dari pustaka `pandas` yang digunakan untuk mengambil sampel acak dari DataFrame. `10000`: Argumen ini menunjukkan jumlah baris yang ingin diambil sebagai sampel dari DataFrame. Dalam hal ini, dipilih untuk mengambil 10.000 baris sebagai sampel acak dari DataFrame `papers`.

Ketika perintah ini dieksekusi, DataFrame `papers` akan berisi 10.000 baris yang diambil secara acak dari data aslinya. Sampel tersebut dapat digunakan untuk analisis yang lebih cepat atau untuk mengurangi ukuran data yang digunakan tanpa kehilangan representasi signifikan dari keseluruhan data.

```
papers['paper_text'] = papers['content'].str.lower()
```

1. Perintah `papers['paper_text'] = papers['content'].str.lower()` digunakan untuk membuat kolom baru dalam DataFrame `papers` dengan nama `'paper_text'`, yang berisi teks dari kolom `'content'` yang telah diubah menjadi huruf kecil (lowercase).
2. `papers['content']`: Merujuk pada kolom `'content'` dalam DataFrame `papers`. Kolom ini berisi teks atau komentar-komentar.
3. `.str.lower()`: Ini adalah metode dari objek Series di pandas yang digunakan untuk mengonversi setiap teks dalam kolom `'content'` menjadi huruf kecil atau lowercase.
4. `papers['paper_text']`: Ini adalah penugasan hasil dari konversi teks menjadi huruf kecil ke dalam kolom baru dengan nama `'paper_text'` di dalam DataFrame `papers`.

Jadi, setelah perintah ini dieksekusi, DataFrame `papers` akan memiliki kolom baru `'paper_text'` yang berisi teks dari kolom `'content'` dengan semua huruf diubah menjadi huruf kecil. Hal ini sering digunakan untuk mempermudah pemrosesan dan analisis teks, karena mengubah teks menjadi lowercase membantu untuk konsistensi dalam pencarian dan pengelompokan teks dalam analisis data.

## 1. Menghapus tanda baca/huruf kecil

Selanjutnya, mari kita lakukan prapemrosesan pada konten kolom `paper_text` agar lebih mudah dianalisis dan hasilnya dapat diandalkan. Untuk melakukannya, kami akan menggunakan ekspresi reguler untuk menghapus tanda baca apa pun, lalu huruf kecil pada teksnya

```
# Load the regular expression library
import re
```

```
# Remove punctuation
papers['paper_text_processed'] =
papers['paper_text'].map(lambda x: re.sub('[,\.!?!]', "", x))
```

```
# Convert the titles to lowercase
papers['paper_text_processed'] =
papers['paper_text_processed'].map(lambda x: x.lower())

# Print out the first rows of papers
papers['paper_text_processed'].head()
```

Perintah ini adalah bagian dari proses pra-pemrosesan teks di dalam DataFrame `papers` menggunakan modul `re` (regular expression) dan pustaka `pandas` di Python. `import re`: Ini adalah perintah untuk memuat modul regular expression (`re`) yang memungkinkan penggunaan ekspresi reguler untuk manipulasi teks.

`papers['paper_text_processed'] = papers['paper_text'].map(lambda x: re.sub('[\.\!?\']', '', x))`: Perintah ini menghapus tanda baca dari teks di dalam kolom `'paper_text'` di DataFrame `papers`. Ini dilakukan dengan menggunakan ekspresi reguler untuk mengganti (substitusi) tanda baca seperti koma, titik, tanda seru, dan tanda tanya dengan string kosong (`''`). Fungsi `lambda` digunakan di sini untuk menerapkan perubahan ini ke setiap baris di kolom `'paper_text'`.

`papers['paper_text_processed'] = papers['paper_text_processed'].map(lambda x: x.lower())`: Setelah menghapus tanda baca, perintah ini mengonversi teks di dalam kolom `'paper_text_processed'` menjadi huruf kecil (`lowercase`). Ini dilakukan menggunakan fungsi `lambda` untuk menerapkan operasi `lowercase` ke setiap baris teks di kolom `'paper_text_processed'`. `papers['paper_text_processed'].head()`: Perintah ini mencetak beberapa baris pertama dari kolom `'paper_text_processed'` dari DataFrame `papers`, menampilkan teks yang telah melalui proses penghapusan tanda baca dan konversi ke huruf kecil.

Dengan demikian, proses ini adalah bagian dari tahap pra-pemrosesan teks yang umum dilakukan sebelum melakukan analisis teks lebih lanjut, seperti pemodelan atau pemrosesan



lanjutan untuk tujuan tertentu seperti analisis sentimen atau pemodelan bahasa alami.

## 2. Tokenize words and further clean-up text

Let's tokenize each sentence into a list of words, removing punctuations and unnecessary characters altogether.

```
import gensim
from gensim.utils import simple_preprocess

def sent_to_words(sentences):
    for sentence in sentences:
        yield(gensim.utils.simple_preprocess(str(sentence),
        deacc=True)) # deacc=True removes punctuations

data = papers.paper_text_processed.values.tolist()
data_words = list(sent_to_words(data))

print(data_words[:1][0][:30])
```

Kode menggunakan pustaka gensim dalam Python, yang umumnya digunakan untuk pemodelan teks dan pemrosesan bahasa alami. `import gensim`: Ini adalah perintah untuk memuat pustaka gensim, yang memiliki alat dan fungsi untuk pemodelan teks dan pemrosesan bahasa alami.

`from gensim.utils import simple_preprocess`: Ini mengimpor fungsi `simple_preprocess` dari `gensim.utils`. Fungsi ini berguna untuk memproses teks secara sederhana, seperti membagi teks menjadi kata-kata kecil (lowercase) dan menghapus aksara.

`def sent_to_words(sentences): ...`: Ini adalah definisi dari sebuah fungsi bernama `sent_to_words`. Fungsi ini menerima daftar kalimat atau teks (`sentences`) dan memprosesnya menjadi kata-kata kecil tanpa aksara (punctuation) menggunakan `simple_preprocess`. Fungsi ini menggunakan generator (`yield`)

untuk menghasilkan kata-kata dari setiap kalimat yang diberikan ke fungsi.

### 3. Pemodelan Frase: Model Bigram dan Trigram

Bigram adalah dua kata yang sering muncul bersamaan dalam dokumen. Trigram adalah 3 kata yang sering muncul. Beberapa contoh dalam contoh kita adalah: 'back\_bumper', 'oil\_leakage', 'maryland\_college\_park' dll. Model Frase Gensim dapat membangun dan mengimplementasikan bigram, trigram, quadgram, dan lainnya. Dua argumen penting pada Frase adalah `min_count` dan ambang batas. Semakin tinggi nilai param ini, semakin sulit kata-kata untuk digabungkan.

`data = papers.paper_text_processed.values.tolist()`: Ini mengambil kolom 'paper\_text\_processed' dari DataFrame `papers` dan mengonversinya ke dalam bentuk daftar (list). Kolom ini berisi teks yang telah diolah sebelumnya.

`data_words = list(sent_to_words(data))`: Fungsi `sent_to_words` yang telah didefinisikan sebelumnya diterapkan ke data (kolom 'paper\_text\_processed') untuk memproses teks menjadi daftar kata-kata kecil tanpa aksara. Hasilnya disimpan dalam variabel `data_words`.

`print(data_words[:1][0][:30])`:

Perintah ini mencetak 30 kata pertama dari hasil pemrosesan teks (`data_words`) untuk satu baris teks pertama yang telah diproses sebelumnya. Jadi, keseluruhan kode ini digunakan untuk mengubah teks yang terdapat dalam kolom 'paper\_text\_processed' dari DataFrame `papers` menjadi daftar kata-kata kecil tanpa aksara (punctuation) menggunakan pustaka `gensim` dan fungsi `simple_preprocess` untuk analisis teks lebih lanjut.

Kode ini menggunakan pustaka `gensim` untuk mengidentifikasi dan membentuk bigram (pasangan dua kata) dan trigram (pasangan tiga kata) dari daftar kata-kata yang telah diproses

sebelumnya dalam variabel `data_words`. Mari kita bahas langkah-langkahnya:

```
bigram = gensim.models.Phrases(data_words, min_count=5,
threshold=100):
```

Di sini, `gensim.models.Phrases` digunakan untuk membentuk bigram dari `data_words`. Parameter `min_count` mengontrol jumlah minimum kemunculan kata dalam teks agar menjadi bigram, sedangkan `threshold` adalah nilai yang menentukan seberapa sering pasangan kata harus muncul agar dianggap sebagai bigram. Semakin tinggi nilai `threshold`, semakin sedikit bigram yang dihasilkan.

```
trigram = gensim.models.Phrases(bigram[data_words],
threshold=100):
```

Langkah ini menggunakan bigram yang telah dibuat sebelumnya sebagai dasar untuk membentuk trigram. Dengan menggunakan `bigram[data_words]`, kita menggunakan bigram yang telah dihitung sebelumnya sebagai acuan untuk menemukan trigram. Ini membantu dalam pembentukan trigram berdasarkan bigram yang telah dibentuk sebelumnya.

`bigram_mod = gensim.models.phrases.Phraser(bigram):` Untuk mempercepat proses pembentukan bigram, `gensim.models.phrases.Phraser` digunakan untuk membuat objek `bigram_mod` dari bigram yang telah dibuat sebelumnya. Objek `bigram_mod` ini dapat digunakan untuk menerapkan bigram ke teks.

`trigram_mod = gensim.models.phrases.Phraser(trigram):` Sama seperti langkah sebelumnya, `trigram` yang telah dihitung sebelumnya dikonversi menjadi objek `trigram_mod` menggunakan `gensim.models.phrases.Phraser`. Ini memungkinkan penggunaan trigram dalam pemrosesan teks.

Dengan menggunakan langkah-langkah ini, bigram dan trigram diidentifikasi dari daftar kata-kata, dan objek `bigram_mod` dan

trigram\_mod dapat digunakan untuk menerapkan bigram dan trigram tersebut pada teks dengan cepat dan efisien. Ini berguna dalam pemodelan teks atau analisis berikutnya yang memerlukan penggunaan bigram dan trigram.

#### 4. Remove Stopwords, Make Bigrams and Lemmatize

The phrase models are ready. Let's define the functions to remove the stopwords, make trigrams and lemmatization and call them sequentially.

```
# NLTK Stop words
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords

#stop_words = stopwords.words('english')
stop_words = stopwords.words('indonesian')
#stop_words.extend(['from', 'subject', 're
stop_words.extend(['yg", "dg", "rt", "dgn", "ny", "d", 'klo',
'sy','saya','kalo', 'amp', 'biar', 'bikin',
'pen', 'u', 'nan', 'loh', 'rt', '&', 'yah'])
```

Kode ini menggunakan pustaka Natural Language Toolkit (NLTK) di Python untuk mengunduh dan menggunakan kata-kata stop (stop words) dalam bahasa Indonesia untuk pemrosesan teks lebih lanjut.

import nltk: Ini mengimpor pustaka NLTK, pustaka yang sering digunakan dalam pemrosesan bahasa alami di Python.

nltk.download('stopwords'): Ini adalah perintah untuk mengunduh dataset kata-kata stop dari NLTK. Dataset ini berisi daftar kata-kata yang umumnya dianggap tidak memiliki makna penting dalam analisis teks karena mereka sangat umum dan sering muncul dalam bahasa tertentu.

from nltk.corpus import stopwords: Setelah dataset stop words diunduh, kita mengimpor modul stopwords dari corpus NLTK. Modul ini berisi daftar kata-kata stop dalam berbagai bahasa.

stop\_words = stopwords.words('indonesian'): Di sini, kita menggunakan daftar kata-kata stop dalam bahasa Indonesia yang telah diunduh dari NLTK. Variabel stop\_words akan berisi daftar kata-kata tersebut, yang akan digunakan untuk menghapus kata-kata ini dari teks dalam proses pra-pemrosesan.

stop\_words.extend([...]): Baris ini digunakan untuk menambahkan kata-kata tambahan ke dalam daftar stop words. Dalam contoh ini, terdapat beberapa kata tambahan yang ditambahkan ke dalam daftar stop words bahasa Indonesia seperti "yg", "dg", "rt", dan lain-lain. Ini bisa dilakukan untuk menyesuaikan daftar kata-kata stop sesuai dengan kebutuhan analisis atau pemrosesan teks yang sedang dilakukan.

Jadi, kode ini membantu untuk memuat daftar kata-kata stop dalam bahasa Indonesia dan menambahkan beberapa kata tambahan ke dalam daftar tersebut agar dapat digunakan dalam proses pra-pemrosesan teks. Hal ini berguna untuk menghilangkan kata-kata yang tidak relevan atau yang biasanya tidak memberikan informasi penting dalam analisis teks.

## 5. Transformasi data: Korpus dan Kamus

Dua masukan utama pada model topik LDA adalah kamus (id2word) dan korpus. Mari kita buat.

```
import gensim.corpora as corpora
```

```
# Create Dictionary
```

```
#id2word = corpora.Dictionary(data_lemmatized)
```

```
id2word = corpora.Dictionary(data_words)
```

```
# Create Corpus
```

```
texts = data_words
```

```
# Term Document Frequency
corpus = [id2word.doc2bow(text) for text in texts]

# View
print(corpus[:1][0][:30])
```

Kode ini menggunakan pustaka Gensim di Python untuk membangun representasi numerik dari teks yang disebut "Bag-of-Words" (BoW). BoW mengubah teks ke dalam representasi vektor di mana setiap kata diwakili sebagai fitur, dan nilai di setiap fitur menunjukkan jumlah kemunculan kata tersebut dalam teks.

`import gensim.corpora as corpora`: Ini mengimpor modul `corpora` dari pustaka Gensim, yang berguna untuk membangun model-model teks.

`id2word = corpora.Dictionary(data_words)`: Di sini, sebuah kamus (`dictionary`) dibuat menggunakan `Dictionary` dari modul `corpora`. `Dictionary` ini memetakan kata-kata dalam `data_words` ke indeks numerik. Setiap kata dalam `data_words` akan diberikan sebuah ID numerik yang unik.

`texts = data_words`: Data yang telah di-preprocess (`data_words`) disimpan dalam variabel `texts`.

`corpus = [id2word.doc2bow(text) for text in texts]`: Langkah ini membangun representasi BoW dari teks yang telah dipreprocess (`texts`). Metode `doc2bow` dari objek `id2word` digunakan untuk mengonversi setiap dokumen (teks) dalam `texts` menjadi representasi BoW. BoW ini terdiri dari tupel (`word_id`, `word_frequency`), yang menunjukkan ID kata dan frekuensi kata dalam teks.

`print(corpus[:1][0][:30])`: Perintah ini mencetak 30 elemen pertama dari representasi BoW dari teks pertama yang telah dihasilkan sebelumnya.

Jadi, keseluruhan kode ini bertujuan untuk membuat representasi BoW dari teks yang telah dipreprocess dan membangun `corpus`

BoW yang siap digunakan untuk model-topik atau analisis lanjutan lainnya menggunakan pustaka Gensim.

## J. Model Dasar

Kami memiliki semua yang diperlukan untuk melatih model LDA dasar. Selain korpus dan kamus, Anda juga perlu menyediakan jumlah topik. Selain itu, alpha dan eta merupakan hyperparameter yang mempengaruhi ketersebaran topik. Menurut dokumen Gensim, keduanya default ke 1.0/num\_topics sebelumnya (kami akan menggunakan default untuk model dasar).

chunksize mengontrol berapa banyak dokumen yang diproses sekaligus dalam algoritma pelatihan. Meningkatkan ukuran potongan akan mempercepat pelatihan, setidaknya selama potongan dokumen tersebut mudah masuk ke dalam memori.

pass mengontrol seberapa sering kita melatih model di seluruh korpus (disetel ke 10). Kata lain untuk pass mungkin adalah "zaman". iterasi agak bersifat teknis, namun pada dasarnya ini mengontrol seberapa sering kita mengulangi perulangan tertentu pada setiap dokumen. Penting untuk menetapkan jumlah "pass" dan "iterasi" yang cukup tinggi.

```
# Build LDA model
```

```
lda_model = gensim.models.LdaMulticore(corpus=corpus,  
                                         id2word=id2word,  
                                         num_topics=10,  
                                         random_state=100,  
                                         chunksize=100,  
                                         passes=10,  
                                         per_word_topics=True)
```

Perintah ini menggunakan pustaka Gensim di Python untuk membangun model LDA (Latent Dirichlet Allocation) yang merupakan metode untuk menemukan topik-topik tersembunyi dalam koleksi dokumen.

```
lda_model = gensim.models.LdaMulticore(corpus=corpus,  
id2word=id2word, num_topics=10, random_state=100,  
chunksize=100, passes=10, per_word_topics=True)
```

`gensim.models.LdaMulticore`: Ini adalah metode untuk membuat model LDA dengan implementasi multicore yang memungkinkan penggunaan beberapa core CPU untuk pelatihan yang lebih cepat.

`corpus=corpus`: Parameter ini adalah representasi BoW dari dokumen yang telah disiapkan sebelumnya dengan menggunakan fungsi `corpora.Dictionary` dan `id2word.doc2bow`.

`id2word=id2word`: Parameter ini adalah kamus (dictionary) yang telah dibuat untuk memetakan kata-kata ke indeks numerik.

`num_topics=10`: Ini adalah jumlah topik yang ingin diidentifikasi dalam model LDA. Dalam contoh ini, model diatur untuk mencari 10 topik tersembunyi dalam koleksi dokumen.

`random_state=100`: Parameter ini mengatur nilai awal untuk pengacakan yang memastikan hasil yang konsisten saat model dilatih ulang.

`chunksize=100`: Ukuran blok untuk pemrosesan paralel dalam model multicore.

`passes=10`: Jumlah iterasi untuk melatih model pada seluruh corpus.

`per_word_topics=True`: Parameter ini mengatur untuk menghasilkan informasi topik untuk setiap kata dalam dokumen, bukan hanya topik utama dari dokumen itu sendiri.

Perintah ini akan membuat sebuah model LDA yang akan mencoba mengidentifikasi 10 topik tersembunyi dalam koleksi dokumen berdasarkan representasi BoW yang telah dibuat sebelumnya. Model ini kemudian dapat digunakan untuk mengeksplorasi dan menganalisis topik dalam dokumen-dokumen tersebut.



Model LDA di atas dibangun dengan 10 topik berbeda dimana setiap topik merupakan kombinasi kata kunci dan setiap kata kunci memberikan kontribusi bobot tertentu pada topik.

Anda dapat melihat kata kunci untuk setiap topik dan bobot (pentingnya) setiap kata kunci menggunakan `lda_model.print_topics()`

```
from pprint import pprint

# Print the Keyword in the 10 topics
pprint(lda_model.print_topics())
doc_lda = lda_model[corpus]
```

Perintah ini menggunakan pustaka `pprint` untuk mencetak topik-topik yang telah ditemukan oleh model LDA (Latent Dirichlet Allocation) yang telah dilatih sebelumnya.

`from pprint import pprint`: Ini mengimpor fungsi `pprint` dari pustaka `pprint`. `pprint` (pretty-print) digunakan untuk mencetak output dengan tata letak yang lebih baik dan lebih mudah dibaca daripada fungsi `print` biasa.

`pprint(lda_model.print_topics())`: Perintah ini mencetak topik-topik yang telah ditemukan oleh model LDA yang telah dilatih sebelumnya. Fungsi `print_topics()` pada objek model LDA mengembalikan daftar topik dengan kata-kata kunci yang paling berkaitan dengan setiap topik.

`doc_lda = lda_model[corpus]`: Ini adalah cara untuk menerapkan model LDA yang telah dilatih pada seluruh dataset (corpus) yang telah digunakan sebelumnya untuk melatih model. Hasilnya disimpan dalam variabel `doc_lda`. Ini akan memberikan distribusi topik untuk setiap dokumen dalam corpus, yaitu, seberapa kuat setiap dokumen terhubung dengan setiap topik.

Dengan menggunakan `pprint(lda_model.print_topics())`, kita mendapatkan tampilan yang terstruktur dan mudah dibaca tentang kata-kata kunci yang paling berkaitan dengan setiap topik yang

ditemukan oleh model LDA. Ini membantu dalam memahami topik-topik yang mungkin ada dalam koleksi dokumen yang telah diproses menggunakan model LDA tersebut.

Hasilnya sebagai berikut:

```
[(0,
  '0.075*"bayar" + 0.048*"udah" + 0.034*"tempo" +
  0.031*"pinjam" + '
  '0.031*"limit" + 0.027*"jatuh" + 0.024*"kecewa" +
  0.023*"telat" + '
  '0.023*"pengajuan" + 0.022*"pembayaran")),
(1,
  '0.134*"data" + 0.061*"tolong" + 0.047*"hapus" +
  0.039*"mohon" + '
  '0.038*"nama" + 0.032*"pengajuan" + 0.031*"pinjaman" +
  0.030*"tolak" + '
  '0.022*"uninstall" + 0.021*"saya")),
(2,
  '0.050*"sistem" + 0.048*"lunas" + 0.031*"ngajuin" +
  0.020*"perbaikan" + '
  '0.019*"kaya" + 0.019*"telepon" + 0.019*"pundi" +
  0.018*"skor" + '
  '0.018*"kirim" + 0.017*"chat")),
(3,
  '0.055*"pinjaman" + 0.035*"semoga" + 0.035*"cepat" +
  0.033*"acc" + '
  '0.033*"kredit" + 0.031*"membantu" + 0.024*"mudah" +
  0.022*"pengajuan" + '
  '0.022*"dana" + 0.022*"proses")),
(4,
  '0.048*"bank" + 0.047*"bunganya" + 0.030*"cicilan" +
  0.029*"tenor" + '
  '0.026*"bunga" + 0.024*"sulit" + 0.023*"kebutuhan" +
  0.022*"tunggu" + '
  '0.022*"selesai" + 0.016*"kredit")),
(5,
  '0.031*"ajukan" + 0.028*"tanggal" + 0.027*"pakai" +
  0.024*"aman" + '
  '0.022*"selesai" + 0.016*"kredit"))]
```

```

'0.019*"suruh" + 0.019*"melunasi" + 0.018*"menit" +
0.018*"sekarang" + '
'0.018*"knp" + 0.016*"upgrade"),
(6,
'0.038*"susah" + 0.028*"tolong" + 0.027*"bukti" +
0.025*"meng" + 0.021*"wa" '
'+ 0.021*"nomor" + 0.021*"masukan" + 0.021*"nomer" +
0.017*"hp" + '
'0.016*"kembalikan"),
(7,
'0.125*"bintang" + 0.057*"acc" + 0.033*"download" +
0.020*"ngajuin" + '
'0.019*"coba" + 0.016*"parah" + 0.013*"di" + 0.013*"dr" +
0.012*"sudah" + '
'0.012*"menunggu"),
(8,
'0.034*"disetujui" + 0.032*"rb" + 0.029*"dah" + 0.027*"trus" +
0.020*"gitu" '
'+ 0.019*"topup" + 0.018*"limitnya" + 0.018*"sehari" +
0.015*"ngisi" + '
'0.014*"eror"),
(9,
'0.053*"masuk" + 0.033*"gagal" + 0.026*"email" +
0.025*"uang" + 0.023*"akun" '
'+ 0.021*"pake" + 0.021*"rekening" + 0.018*"saldo" +
0.017*"cs" + '
'0.017*"maret")]
```

### *Compute Model Perplexity and Coherence Score*

```

from gensim.models import CoherenceModel

# Compute Coherence Score
coherence_model_lda = CoherenceModel(model=lda_model,
texts=data_words, dictionary=id2word, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('Coherence Score: ', coherence_lda)
```

Perintah ini digunakan untuk menghitung skor koherensi (coherence score) dari model Latent Dirichlet Allocation (LDA) yang telah dilatih sebelumnya. Skor koherensi memberikan gambaran tentang seberapa koheren atau terkait topik-topik yang dihasilkan oleh model.

`from gensim.models import CoherenceModel`: Ini mengimpor `CoherenceModel` dari pustaka `Gensim`. `CoherenceModel` digunakan untuk menghitung koherensi dari model topic modeling.

`coherence_model_lda = CoherenceModel(model=lda_model, texts=data_words, dictionary=id2word, coherence='c_v')`: Ini membuat objek `coherence_model_lda` menggunakan `CoherenceModel`. Parameter-parameter yang digunakan adalah: `model=lda_model`: Merujuk pada model LDA yang telah dilatih sebelumnya.

`texts=data_words`: Merupakan teks yang telah diproses sebelumnya, dalam bentuk daftar kata-kata.

`dictionary=id2word`: Kamus yang memetakan kata-kata ke indeks numerik.

`coherence='c_v'`: Jenis koherensi yang digunakan. Dalam hal ini, `'c_v'` adalah metode koherensi yang dikenal sebagai `Coherence 'c_v'`.

`coherence_lda = coherence_model_lda.get_coherence()`: Langkah ini menghitung skor koherensi dengan menggunakan metode `get_coherence()` dari objek `coherence_model_lda`. Skor koherensi akan memberikan gambaran tentang seberapa baik topik-topik yang dihasilkan oleh model LDA.

`print('Coherence Score: ', coherence_lda)`: Perintah ini mencetak skor koherensi yang telah dihitung sebelumnya.

Dengan menggunakan skor koherensi, kita mendapatkan ukuran kualitas dari topik-topik yang dihasilkan oleh model. Semakin tinggi skor koherensi, semakin baik atau lebih terkait topik-topik yang dihasilkan oleh model tersebut.

## K. Penyetelan tuning hyper parameter

Pertama, mari kita bedakan antara hyperparameter model dan parameter model:

Hyperparameter model dapat dianggap sebagai pengaturan untuk algoritma pembelajaran mesin yang disetel oleh data scientist sebelum pelatihan. Contohnya adalah jumlah pohon di hutan acak, atau dalam kasus kami, jumlah topik K. Parameter model dapat dianggap sebagai apa yang dipelajari model selama pelatihan, seperti bobot setiap kata dalam topik tertentu. Sekarang kita memiliki skor koherensi dasar untuk model LDA default, mari kita lakukan serangkaian uji sensitivitas untuk membantu menentukan hyperparameter model berikut:

Jumlah Topik (K), Dirichlet hyperparameter alpha: Kepadatan Topik Dokumen, Hyperparameter beta Dirichlet: Kepadatan Topik Kata

Kami akan melakukan pengujian ini secara berurutan, satu parameter pada satu waktu dengan menjaga parameter lainnya tetap konstan dan menjalankannya pada dua set korpus validasi perbedaan. Kami akan menggunakan C\_v sebagai metrik pilihan kami untuk perbandingan kinerja

```
# supporting function
def compute_coherence_values(corpus, dictionary, k, a, b):
    lda_model = gensim.models.LdaMulticore(corpus=corpus,
                                             id2word=dictionary,
                                             num_topics=k,
                                             random_state=100,
                                             chunksize=100,
```

```

        passes=10,
        alpha=a,
        eta=b)
    coherence_model_lda = CoherenceModel(model=lda_model,
    texts=data_words, dictionary=id2word, coherence='c_v')

    return coherence_model_lda.get_coherence()

```

Fungsi ini digunakan untuk menghitung skor koherensi (coherence score) dari model LDA (Latent Dirichlet Allocation) yang dibangun dengan berbagai nilai alpha dan eta. Fungsi ini membantu dalam mengevaluasi bagaimana nilai-nilai parameter ini mempengaruhi kualitas topik-topik yang dihasilkan oleh model.

Argumen-argumen dalam fungsi `compute_coherence_values`:  
 corpus: Representasi Bag-of-Words (BoW) dari dokumen yang akan digunakan untuk melatih model.  
 dictionary: Kamus yang memetakan kata-kata ke indeks numerik.  
 k: Jumlah topik yang diuji.

1. Parameter alpha yang digunakan dalam model LDA. Ini mengontrol distribusi topik dalam dokumen. Nilai alpha yang lebih tinggi menyebabkan dokumen memiliki distribusi topik yang lebih merata.
2. Parameter eta yang digunakan dalam model LDA. Ini mengontrol distribusi kata dalam topik. Nilai eta yang lebih tinggi menyebabkan topik memiliki distribusi kata yang lebih merata.

Langkah-langkahnya adalah sebagai berikut:

`lda_model = gensim.models.LdaMulticore(...)`: Fungsi ini menggunakan model LDA multicore dari Gensim untuk melatih model LDA dengan parameter yang diberikan (corpus, dictionary, num\_topics=k, alpha=a, eta=b, dst).

`coherence_model_lda = CoherenceModel(...)`: Langkah ini membuat objek `coherence_model_lda` menggunakan `CoherenceModel` dari Gensim dengan menggunakan model LDA yang telah dilatih sebelumnya.

`return coherence_model_lda.get_coherence()`: Fungsi mengembalikan skor koherensi yang dihitung menggunakan metode `get_coherence()` dari objek `coherence_model_lda`. Ini memberikan informasi tentang seberapa koheren topik-topik yang dihasilkan oleh model dengan kombinasi parameter `alpha` dan `eta` tertentu.

Dengan menggunakan fungsi ini, Anda dapat menguji dan membandingkan berbagai nilai `alpha` dan `eta` untuk mengevaluasi bagaimana hal itu mempengaruhi kualitas topik yang dihasilkan oleh model LDA. memanggil fungsinya, dan ulangi pada rentang nilai parameter topik, `alfa`, dan `beta`

```
#bagian 1
import numpy as np
import tqdm

grid = {}
grid['Validation_Set'] = {}

# Topics range
min_topics = 2
max_topics = 11
step_size = 1
topics_range = range(min_topics, max_topics, step_size)

# Alpha parameter
alpha = list(np.arange(0.01, 1, 0.3))
alpha.append('symmetric')
alpha.append('asymmetric')

# Beta parameter
beta = list(np.arange(0.01, 1, 0.3))
```

```

beta.append('symmetric')

# Validation sets
num_of_docs = len(corpus)
corpus_sets = [gensim.utils.ClippedCorpus(corpus,
int(num_of_docs*0.75)),
                corpus]

corpus_title = ['75% Corpus', '100% Corpus']

model_results = {'Validation_Set': [],
                 'Topics': [],
                 'Alpha': [],
                 'Beta': [],
                 'Coherence': []
                }

#Bagian 2
if 1 == 1:
    pbar =
    tqdm.tqdm(total=(len(beta)*len(alpha)*len(topics_range)*len(co
corpus_title)))

    # iterate through validation corpuses
    for i in range(len(corpus_sets)):
        # iterate through number of topics
        for k in topics_range:
            # iterate through alpha values
            for a in alpha:
                # iterare through beta values
                for b in beta:
                    # get the coherence score for the given parameters
                    cv =
compute_coherence_values(corpus=corpus_sets[i],
dictionary=id2word,
                        k=k, a=a, b=b)
                # Save the model results

```



```

        model_results['Validation_Set'].append(corpus_title[
i])

        model_results['Topics'].append(k)
        model_results['Alpha'].append(a)
        model_results['Beta'].append(b)
        model_results['Coherence'].append(cv)

        pbar.update(1)
        pd.DataFrame(model_results).to_csv('./drive/My
Drive/dataset/fintechP2P/2023/13feb_tuning_results.csv',
index=False)
    pbar.close()

```

Bagian 1 Ini adalah kode untuk menyiapkan berbagai parameter yang akan digunakan dalam eksplorasi model LDA (Latent Dirichlet Allocation) dengan menggunakan teknik grid search. Import Numpy dan tqdm: Baris pertama mengimpor modul NumPy untuk operasi numerik dan tqdm, sebuah modul yang membantu membuat bar progres saat iterasi yang lama.

Variabel grid: Ini adalah wadah (dictionary) yang akan digunakan untuk menyimpan hasil eksperimen.

Range Topik: Variabel `min_topics`, `max_topics`, dan `step_size` digunakan untuk menentukan rentang nilai topik yang akan dieksplorasi dalam pencarian model LDA yang optimal.

Parameter Alpha dan Beta: `alpha` dan `beta` adalah daftar yang berisi rentang nilai untuk parameter `alpha` dan `beta` yang akan dieksplorasi dalam pencarian model LDA.

Validation Sets: Variabel `num_of_docs`, `corpus_sets`, dan `corpus_title` digunakan untuk menyiapkan data corpus yang akan digunakan untuk validasi. `num_of_docs` adalah jumlah dokumen dalam corpus, `corpus_sets` adalah persentase corpus yang akan digunakan (75% dan 100%), dan `corpus_title` adalah label yang sesuai dengan corpus yang digunakan.

Model Results: Ini adalah wadah untuk menyimpan hasil evaluasi berbagai parameter LDA, seperti nilai koherensi.

Kode ini adalah langkah pertama dalam proses eksperimen yang melibatkan iterasi berbagai kombinasi parameter LDA untuk menentukan parameter mana yang memberikan hasil terbaik dalam hal koherensi topik pada model. Setelah parameter dan data disiapkan seperti ini, dilakukan iterasi dalam eksperimen grid search untuk mengevaluasi kombinasi parameter yang berbeda terhadap kualitas model yang dihasilkan.

Bagian 2 ini adalah bagian yang melakukan iterasi melalui kombinasi parameter LDA yang telah ditentukan sebelumnya untuk mengevaluasi dan menyimpan skor koherensi untuk setiap kombinasi tersebut.

if 1 == 1:: Ini adalah pernyataan yang selalu benar. Ini menunjukkan bahwa iterasi yang dilakukan di bawahnya akan dieksekusi.

pbar = tqdm.tqdm(total=(len(beta)\*len(alpha)\*len(topics\_range)\*len(corpus\_title))): Membuat progress bar menggunakan tqdm untuk mengukur kemajuan dalam iterasi. Total jumlah iterasi yang diharapkan dihitung berdasarkan jumlah kombinasi yang akan dieksekusi.

Iterasi Loop Bersarang:

1. Loop pertama (for i in range(len(corpus\_sets))): Melakukan iterasi melalui berbagai dataset validasi.
2. Loop kedua (for k in topics\_range): Melakukan iterasi melalui rentang nilai topik.
3. Loop ketiga (for a in alpha): Melakukan iterasi melalui nilai alpha yang telah ditentukan sebelumnya.
4. Loop keempat (for b in beta): Melakukan iterasi melalui nilai beta yang telah ditentukan sebelumnya.

Dalam setiap iterasi kombinasi parameter yang berbeda (corpus, dictionary, k, a, dan b), fungsi `compute_coherence_values` dipanggil untuk menghitung skor koherensi (cv) dari model LDA

yang dibuat berdasarkan parameter-parameter tersebut. Hasil skor koherensi kemudian disimpan dalam model\_results sesuai dengan kombinasi parameter yang digunakan. Setelah semua iterasi selesai, hasil dari model\_results disimpan sebagai file CSV di direktori yang ditentukan menggunakan `pd.DataFrame(model_results).to_csv('./drive/My Drive/dataset/fintechP2P/2023/13feb_tuning_results.csv', index=False)`.

`pbar.close()`: Setelah proses iterasi selesai, progress bar ditutup. Ini adalah langkah yang memakan waktu untuk mengevaluasi dan mencari parameter terbaik yang memberikan skor koherensi yang paling optimal untuk model LDA. Proses ini melibatkan berbagai kombinasi parameter yang mungkin dan mengevaluasi setiap kombinasi tersebut untuk memilih yang terbaik. Hasil perhitungan dari script di atas adalah sebagai berikut:

No	Validation_Set	Topics	Alpha	Beta	Coherence
1	100% Corpus	10	0.01	0.91	0.554553496
2	100% Corpus	3	0.91	0.91	0.554553496
3	100% Corpus	3	0.91	0.31	0.546011003
4	100% Corpus	3	0.91	symmetric	0.546011003
5	100% Corpus	3	asymmetric	0.01	0.544522395
6	100% Corpus	3	0.61	0.61	0.541360905
7	100% Corpus	3	0.91	0.01	0.540641211
8	100% Corpus	3	0.61	symmetric	0.536954678
9	100% Corpus	3	0.61	0.91	0.535848054
10	100% Corpus	3	0.61	0.31	0.534917947
11	100% Corpus	8	symmetric	0.91	0.532990744
12	75% Corpus	3	0.61	0.31	0.531145051
13	75% Corpus	3	0.91	0.61	0.531057059
14	75% Corpus	3	0.91	0.91	0.531057059
15	100% Corpus	3	0.61	0.01	0.527612144
16	75% Corpus	3	0.61	0.01	0.526229031
17	75% Corpus	3	0.61	symmetric	0.523419221
18	75% Corpus	3	0.91	symmetric	0.52139985

Dari hasil di atas dapat dilihat bahwa nilai coherence yang lebih dari 0.50 terlihat dan terdapat pada 3 dan 8 topic. Yang perlu diperhatikan adalah apakah dengan 8 topic masih mendapatkan makna yang berbeda pada setiap topik, jika maknanya berbeda jumlah topik dapat di set menjadi 8 (nomor urut 11). Namun jika banyak topik yang sama, maka dapat dipilih 10 topic dengan nilai coherence tertinggi terlihat pada nomor urut 1.

Berdasarkan evaluasi eksternal (Kode akan ditambahkan dari analisis berbasis Excel), mari kita latih model akhir dengan parameter yang menghasilkan skor koherensi tertinggi.

```
num_topics = 10
```

```
lda_model = gensim.models.LdaMulticore(corpus=corpus,
                                         id2word=id2word,
                                         num_topics=num_topics,
                                         random_state=100,
                                         chunksize=100,
                                         passes=10,
                                         # alpha='asymmetric',
                                         alpha=0.01,
                                         #eta=0.9
                                         eta=0.91)
```

Perintah di atas adalah bagian dari proses pembentukan model LDA (Latent Dirichlet Allocation) menggunakan pustaka Gensim dalam Python. Ini digunakan untuk membuat model topik dari data teks yang telah diubah menjadi representasi Bag-of-Words (BoW).

`num_topics = 10`: Parameter ini menentukan jumlah topik yang ingin ditemukan dalam korpus teks. Dalam kasus ini, ditetapkan sebagai 10, artinya model akan berusaha menemukan 10 topik yang berbeda dalam teks yang diberikan.

`gensim.models.LdaMulticore`: Ini adalah fungsi untuk membuat model LDA menggunakan teknik multicore dari pustaka Gensim.  
`corpus=corpus`: Representasi BoW dari dokumen yang akan digunakan untuk melatih model.

`id2word=id2word`: Kamus yang memetakan kata-kata ke indeks numerik, digunakan untuk memahami representasi numerik kata-kata dalam model.

`random_state=100`: Seed untuk inisialisasi bilangan acak. Ini memastikan bahwa hasil dari model yang sama akan konsisten ketika dijalankan kembali dengan seed yang sama.

`chunksize=100`: Jumlah dokumen yang akan digunakan dalam satu batch selama proses pelatihan model. Ini mempengaruhi efisiensi dan penggunaan memori saat pelatihan.

`passes=10`: Jumlah iterasi yang dilakukan oleh model saat melatih dataset. Setiap iterasi melibatkan pembaruan bobot pada model.

`alpha=0.01`: Parameter alpha mengontrol seberapa banyak topik yang ada dalam satu dokumen. Nilai yang lebih rendah seperti 0.01 mengindikasikan bahwa dokumen hanya akan memiliki sedikit topik yang dominan.

`eta=0.91`: Parameter eta mengontrol seberapa banyak kata yang terkait dengan satu topik tertentu. Nilai yang lebih tinggi seperti 0.91 menunjukkan bahwa setiap topik akan memiliki banyak kata yang kuat terkait dengannya. Kombinasi dari parameter-parameter ini membentuk model LDA yang akan menemukan 10 topik dalam data teks, dengan distribusi yang dikendalikan oleh nilai alpha dan eta yang telah ditetapkan.

```
from pprint import pprint

# Print the Keyword in the 10 topics
pprint(lda_model.print_topics())
doc_lda = lda_model[corpus]
```

Perintah ini menggunakan pustaka pprint untuk mencetak topik-topik yang telah ditemukan oleh model LDA (Latent Dirichlet Allocation) yang telah dilatih sebelumnya.

Langkah-langkahnya adalah sebagai berikut:

`from pprint import pprint`: Ini mengimpor fungsi pprint dari pustaka pprint. pprint (pretty-print) digunakan untuk mencetak output dengan tata letak yang lebih baik dan lebih mudah dibaca daripada fungsi print biasa.

`pprint(lda_model.print_topics())`: Perintah ini mencetak topik-topik yang telah ditemukan oleh model LDA yang telah dilatih sebelumnya. Fungsi `print_topics()` pada objek model LDA mengembalikan daftar topik dengan kata-kata kunci yang paling berkaitan dengan setiap topik.

`doc_lda = lda_model[corpus]`: Ini adalah cara untuk menerapkan model LDA yang telah dilatih pada seluruh dataset (corpus) yang telah digunakan sebelumnya untuk melatih model. Hasilnya disimpan dalam variabel `doc_lda`. Ini akan memberikan distribusi topik untuk setiap dokumen dalam corpus, yaitu, seberapa kuat setiap dokumen terhubung dengan setiap topik.

## DAFTAR PUSTAKA

- Abdulumumin, I., & Galadanci, B. S. (2019). hauWE: Hausa Words Embedding for Natural Language Processing. *2019 2nd International Conference of the IEEE Nigeria Computer Chapter (NigeriaComputConf)*, 1–6. <https://doi.org/10.1109/NigeriaComputConf45974.2019.8949674>
- Agarwal, M. (2019). An Overview of Natural Language Processing. *International Journal for Research in Applied Science and Engineering Technology*, 7(5), 2811–2813. <https://doi.org/10.22214/ijraset.2019.5462>
- Atkinson-Abutridy, J. (2022). Text analytics: An introduction to the science and applications of unstructured information analysis. In *Text Analytics: An Introduction to the Science and Applications of Unstructured Information Analysis*. <https://doi.org/10.1201/9781003280996>
- Basha, M. J., Vijayakumar, S., Jayashankari, J., Alawadi, A. H., & Durdona, P. (2023). Advancements in Natural Language Processing for Text Understanding. *E3S Web of Conferences*, 399. <https://doi.org/10.1051/e3sconf/202339904031>
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). *Latent Dirichlet Allocation*. 3, 993–1022.
- Gavhane, J., Prasad, R., & Kumar, R. (2022). Graph embeddings for natural language processing. *Graph Learning and Network Science for Natural Language Processing*, 7148, 77–95. <https://doi.org/10.1201/9781003272649-4>
- Gross, A., & Murthy, D. (2014). Modeling virtual organizations with Latent Dirichlet Allocation: A case for natural language processing. *Neural Networks*, 58, 38–49. <https://doi.org/10.1016/j.neunet.2014.05.008>
- Hannon, B., Kumar, Y., Gayle, D., Li, J. J., & Morreale, P. (2024). Robust Testing of AI Language Model Resiliency with Novel Adversarial Prompts. *Electronics (Switzerland)*,

- 13(5). <https://doi.org/10.3390/electronics13050842>
- Heimerl, F., & Gleicher, M. (2018). Interactive Analysis of Word Vector Embeddings. *Computer Graphics Forum*, 37(3), 253–265. <https://doi.org/10.1111/cgf.13417>
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Science*, 349(6245), 261–266. <https://doi.org/10.1126/science.aaa8685>
- Joshi, I., Koringa, P., & Mitra, S. (2019). Word Embeddings in Low Resource Gujarati Language. *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, 110–115. <https://doi.org/10.1109/ICDARW.2019.40090>
- Kjell, O., Giorgi, S., & Schwartz, H. A. (2023). The Text-Package: An R-Package for Analyzing and Visualizing Human Language Using Natural Language Processing and Transformers. *Psychological Methods*, 28(6), 1478–1498. <https://doi.org/10.1037/met0000542>
- Li, J., Sun, A., Han, J., & Li, C. (2022). A Survey on Deep Learning for Named Entity Recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1), 50–70. <https://doi.org/10.1109/TKDE.2020.2981314>
- Meng, Y., Li, X., Sun, Z., & Li, J. (2019). *Query-Based Named Entity Recognition*. <http://arxiv.org/abs/1908.09138>
- Mishra, P. (2019). Natural Language Processing Using PyTorch. In *PyTorch Recipes* (pp. 165–178). Apress. [https://doi.org/10.1007/978-1-4842-4258-2\\_7](https://doi.org/10.1007/978-1-4842-4258-2_7)
- Sun, P., Yang, X., Zhao, X., & Wang, Z. (2018). An Overview of Named Entity Recognition. *2018 International Conference on Asian Language Processing (IALP)*, 273–278. <https://doi.org/10.1109/IALP.2018.8629225>
- Tunstall, L., 22TUNA, Von Werra, L., & Wolf, T. (2022). *NLP with transformers building language applications*.
- Vajjala, S., Majumder, B., Gupta, A., & Surana, H. (2012). Statistical Natural Language Processing. In *SpringerReference*. [https://doi.org/10.1007/springerreference\\_205170](https://doi.org/10.1007/springerreference_205170)



- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems, 2017-Decem*(Nips), 5999–6009.
- Wang, D., Thint, M., & Al-Rubaie, A. (2012). Semi-Supervised Latent Dirichlet Allocation and Its Application for Document Classification. *2012 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*, 306–310. <https://doi.org/10.1109/WI-IAT.2012.211>
- Yu, J., Ji, B., Li, S., Ma, J., Liu, H., & Xu, H. (2022). S-NER: A Concise and Efficient Span-Based Model for Named Entity Recognition. *Sensors*, 22(8). <https://doi.org/10.3390/s22082852>
- Ziyadi, M., Sun, Y., Goswami, A., Huang, J., & Chen, W. (2020). *Example-Based Named Entity Recognition*. 1–15. <http://arxiv.org/abs/2008.10570>

Pengantar NLP dan Topik Model LDA adalah panduan komprehensif yang menyelami dasar-dasar pemrosesan bahasa alami (NLP) dan pemodelan topik menggunakan Latent Dirichlet Allocation (LDA), dua bidang krusial yang semakin penting di era data digital. Buku ini dirancang untuk memberikan pembaca pemahaman menyeluruh tentang konsep, algoritma, dan aplikasi praktis NLP, sekaligus mengeksplorasi pendekatan pemodelan topik untuk mengidentifikasi pola-pola tersembunyi dalam teks. Pada bagian awal, buku ini memperkenalkan fondasi NLP —meliputi definisi, sejarah, teknik dasar, dan relevansinya di berbagai sektor seperti kesehatan, industri, dan pendidikan. Setelah itu, pembaca diajak memahami algoritma populer yang digunakan dalam NLP, seperti Naive Bayes, SVM, dan jaringan saraf buatan, serta tantangan yang dihadapi dalam penerapan teknologi ini.

Buku ini kemudian beralih ke pembahasan Latent Dirichlet Allocation, metode probabilistik yang kuat untuk menemukan topik dalam kumpulan teks besar. Pembaca diperkenalkan pada prinsip kerja LDA, termasuk penyesuaian parameter dan langkah-langkah praktis dalam membangun model topik. Dengan contoh-contoh aplikatif, buku ini memudahkan pembaca dalam memahami bagaimana NLP dan LDA bekerja di dunia nyata. Pengantar NLP dan Topik Model LDA adalah sumber berharga bagi mereka yang ingin memulai perjalanan di dunia NLP dan pemodelan topik, serta bagi mereka yang ingin memperdalam pemahaman mereka tentang penerapan teknologi ini di berbagai bidang profesional dan penelitian.

Diterbitkan oleh



ISBN 978-623-10-4853-0

