



Natural Language Processing

Introduction

- ▶ Natural Language Processing (NLP) is a subfield of artificial intelligence (AI) that focuses on the interaction between computers and human language. The goal of NLP is to enable machines to read, understand, and generate human language in a way that is both meaningful and useful.
- ▶ NLP combines linguistics, computer science, and machine learning techniques to solve problems involving language. It is widely used in applications such as virtual assistants (e.g., Siri, Alexa), translation services (e.g., Google Translate), sentiment analysis in social media, and much more.
- ▶ The challenges in NLP include ambiguity, slang, context understanding, and the variety of languages and dialects. With advancements in deep learning, NLP has made significant progress in recent years, particularly with models like BERT and GPT, which understand and generate text with impressive accuracy.

Introduction

- ▶ NLP encompasses a variety of tasks, including:
 1. **Text Classification:** Assigning categories to text, such as sentiment analysis (positive/negative) or spam detection.
 2. **Named Entity Recognition (NER):** Identifying and classifying entities in text, such as names of people, organizations, dates, etc.
 3. **Machine Translation:** Automatically translating text from one language to another.
 4. **Text Generation:** Creating new, coherent text based on given input (e.g., chatbots, content creation).
 5. **Speech Recognition:** Converting spoken language into text.
 6. **Part-of-Speech Tagging:** Identifying the grammatical components (nouns, verbs, etc.) in a sentence.
 7. **Sentiment Analysis:** Determining the sentiment (positive, negative, neutral) expressed in a piece of text.

Text Processing

- ▶ Tokenization
- ▶ Stop word removal
- ▶ Lemmatization and Stemming



Cosine Similarity

- ▶ **Cosine Similarity** is a metric used to measure how similar two text documents (or vectors) are, based on the cosine of the angle between them.
- ▶ It's often used in text analysis and natural language processing (NLP) to determine the similarity between two documents, regardless of their size. The cosine similarity ranges from 0 (completely dissimilar) to 1 (completely similar).

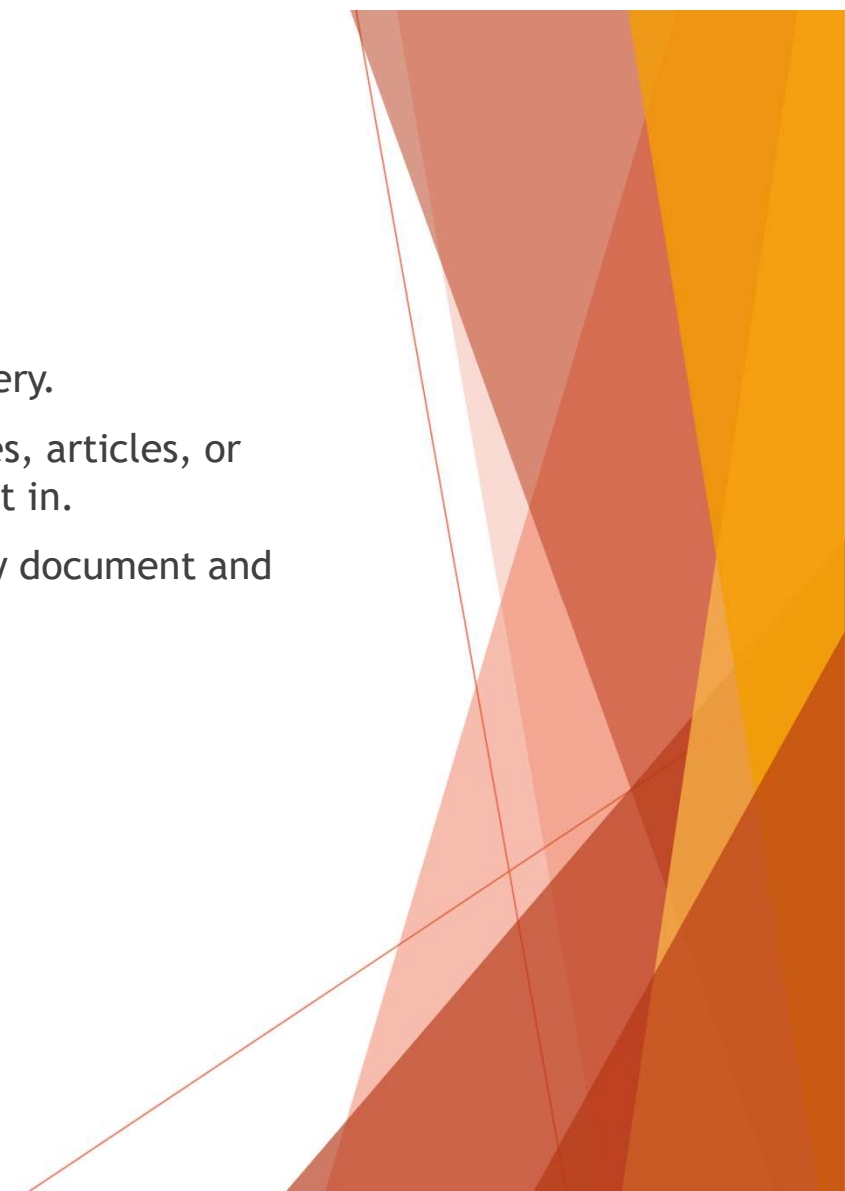
$$\text{Cosine Similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

Where:

- $A \cdot B$ is the dot product of vectors A and B .
- $\|A\|$ and $\|B\|$ are the magnitudes (norms) of vectors A and B .

Cosine Similarity: Applications

- ▶ Search Engines: Finding relevant documents based on a query.
- ▶ Recommendation Systems: Suggesting items (such as movies, articles, or products) that are similar to what a user has shown interest in.
- ▶ Text Classification: Measuring the similarity between a new document and existing labeled categories or clusters.



Bag of Words Model

- ▶ The **Bag of Words (BoW)** model is a simple and commonly used approach for text representation in natural language processing (NLP).
- ▶ It transforms a collection of text (a corpus) into a matrix of token counts, disregarding grammar and word order but keeping track of the frequency of words in the documents.
- ▶ **Key Concept:**
 - **Bag of Words** ignores the syntax and word order and treats the text as a collection of words (or tokens).
 - Each document is represented by a vector where each dimension corresponds to a unique word from the corpus (vocabulary), and the value in each dimension is the frequency (or count) of the word in the document.

BoW: Applications

- ▶ Text Classification: Converting text into a structured format for machine learning algorithms to classify (e.g., spam detection).
- ▶ Sentiment Analysis: Analyzing the sentiment of text by counting the frequency of positive or negative words.
- ▶ Document Similarity: Comparing the frequency of common words between two documents to assess their similarity (e.g., news articles, product reviews).

BoW: Advantages, Disadvantages

► Advantages:

- Simple and effective.
- Works well when the order of words does not significantly impact the meaning.

► Disadvantages:

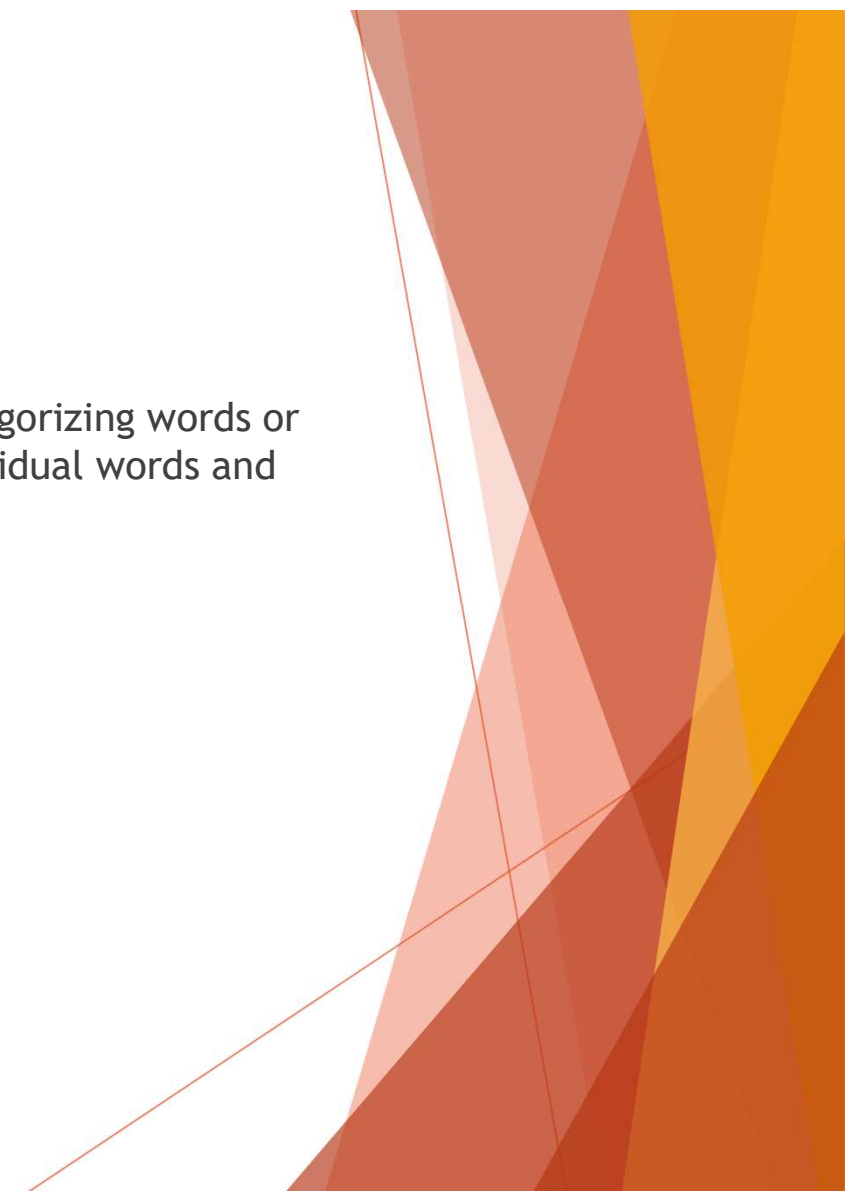
- Ignores word order, context, and syntax.
- Creates high-dimensional sparse vectors, which can lead to large memory usage in the case of a large corpus.

Linguistic Analysis

- ▶ In Natural Language Processing (NLP), various levels of linguistic analysis are performed to understand text. The main types of analysis are:
 - ▶ Lexical
 - ▶ Semantic
 - ▶ Syntactic
 - ▶ Pragmatic

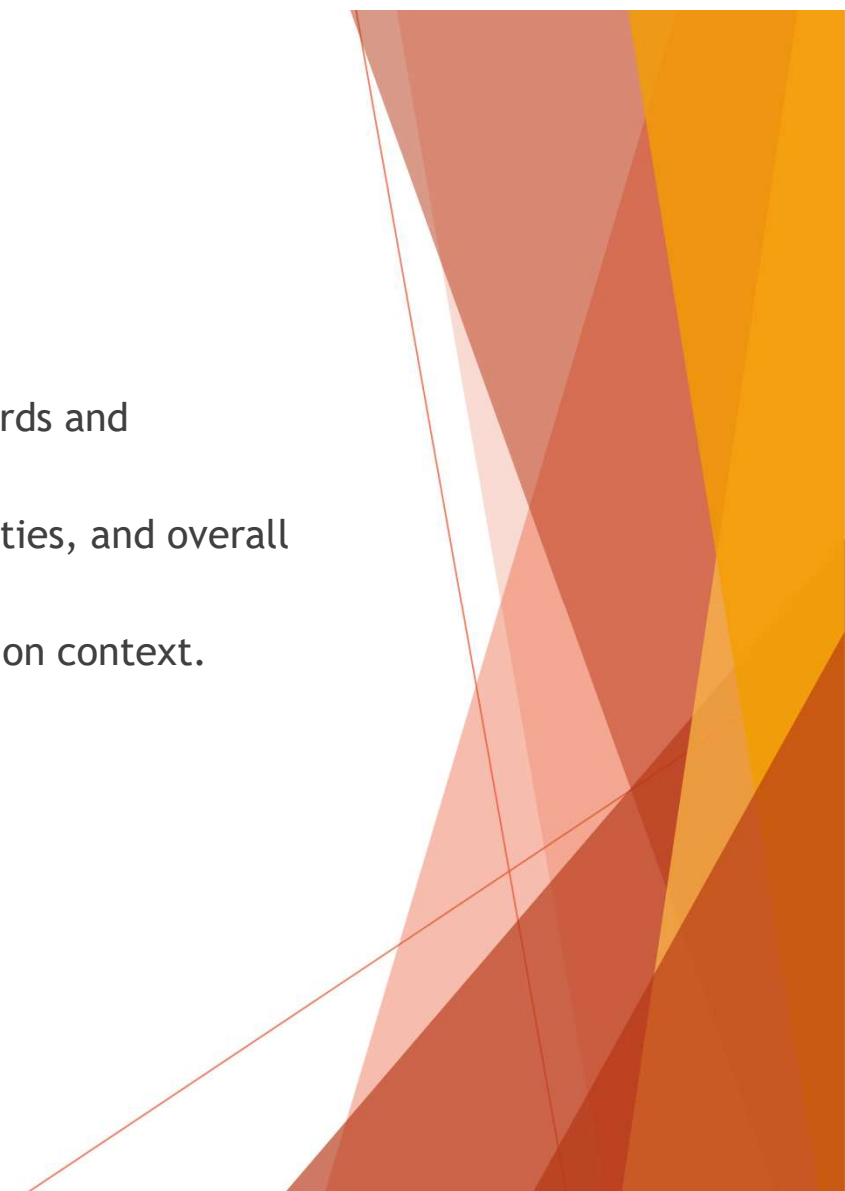
Lexical Analysis

- ▶ Lexical analysis refers to the process of identifying and categorizing words or tokens in text, including breaking down sentences into individual words and their corresponding parts of speech (POS).
- ▶ It is the first step in text processing.



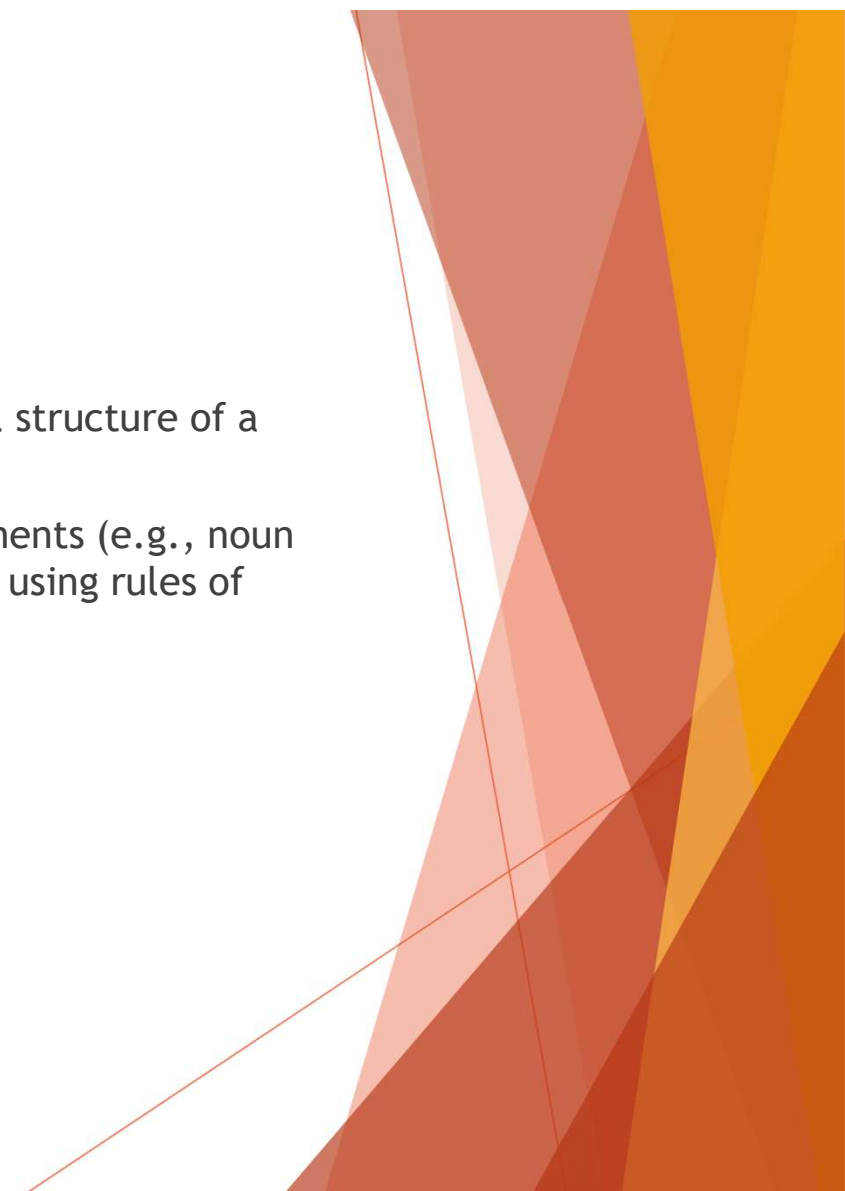
Semantic Analysis

- ▶ Semantic analysis involves understanding the meaning of words and sentences.
- ▶ It includes identifying the relationships between words, entities, and overall meaning.
- ▶ This analysis aims to extract the meaning from a text based on context.



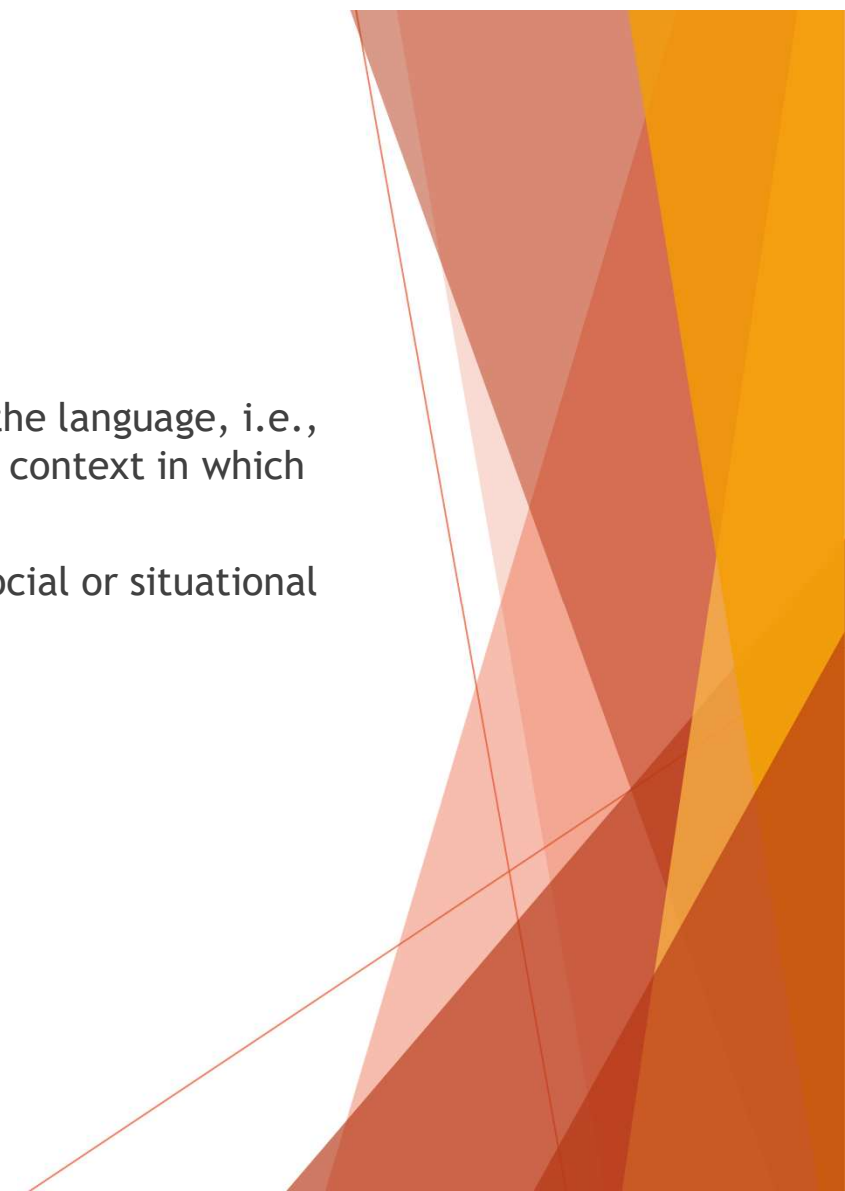
Syntactic Analysis

- ▶ Syntactic analysis focuses on understanding the grammatical structure of a sentence.
- ▶ It involves parsing the sentence into its grammatical components (e.g., noun phrases, verb phrases) and analyzing the sentence structure using rules of syntax.



Pragmatic Analysis

- ▶ Pragmatic analysis deals with understanding the context of the language, i.e., interpreting the intended meaning based on the situation or context in which the text was produced.
- ▶ It involves recognizing implied meaning, sarcasm, and the social or situational context.



Vectorization

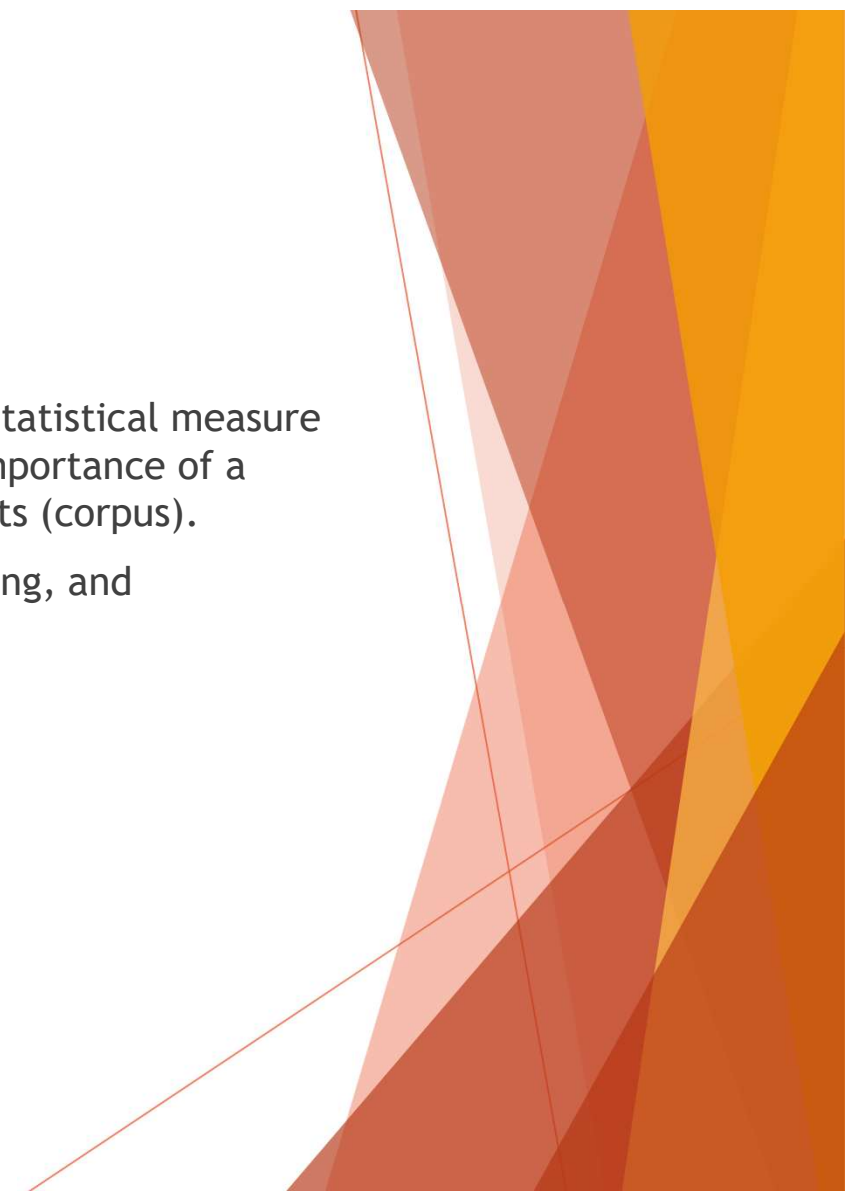
- ▶ **Definition:** Vectorization is the process of converting text data (like words, phrases, or documents) into numerical format by representing each token as a vector of numbers. It's a broad term that includes any method of transforming text into a vector space.
- ▶ **Techniques:** Some common vectorization techniques are:
 - ▶ **One-Hot Encoding:** Represents each word as a binary vector with a single "1" and the rest "0"s, indicating its presence in a large vocabulary.
 - ▶ **Count Vectorization:** Counts the occurrences of each word in a text and represents text as frequency vectors.
 - ▶ **TF-IDF (Term Frequency-Inverse Document Frequency):** A weighting technique that considers the frequency of words while reducing the weight of common terms across documents.

Embedding

- ▶ **Definition:** Embedding is a specific type of vectorization that captures semantic relationships between words by representing them in a dense, low-dimensional vector space. Embeddings are learned representations that map similar words close together in this vector space based on their meanings or contexts.
- ▶ **Techniques:**
 - ▶ **Word2Vec:** Learns word embeddings by predicting word contexts (Skip-gram) or predicting words given contexts (CBOW).
 - ▶ **GloVe (Global Vectors):** Captures global statistical information in word co-occurrence matrices to learn embeddings.
 - ▶ **BERT and GPT-based Embeddings:** Use transformer models to generate contextual embeddings that change depending on surrounding words.

TF-IDF (Term Frequency-Inverse Document Frequency)

- ▶ **TF-IDF** (Term Frequency-Inverse Document Frequency) is a statistical measure used in natural language processing (NLP) to evaluate the importance of a word within a document relative to a collection of documents (corpus).
- ▶ It is used for tasks like text classification, document clustering, and information retrieval.



Components of TF-IDF

- ▶ **Term Frequency** measures how often a word appears in a document. The intuition is that the more a word appears in a document, the more important it is for that document.
- ▶ **Inverse Document Frequency** measures how common or rare a word is across all documents in the corpus. The intuition is that words that appear in many documents are less useful for distinguishing between documents, while words that appear in fewer documents are more valuable.

$$\text{TF}(t, d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

$$\text{IDF}(t) = \log \left(\frac{\text{Total number of documents}}{\text{Number of documents containing term } t} \right)$$

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

Why use TF-IDF?

- ▶ Relevance Weighting: TF-IDF helps identify words that are important to a document while reducing the weight of common words (e.g., "the," "and," "is") that appear in many documents.
- ▶ Feature for Machine Learning: In text classification or clustering, words with high TF-IDF values are typically the most informative features for distinguishing between documents.
 - ❖ The TF-IDF Matrix contains rows (documents) and columns (words in the vocabulary).
 - ❖ Words with higher TF-IDF scores (e.g., "love", "programming") have a larger impact in distinguishing documents.
 - ❖ Words that are common across documents (e.g., "python") will have lower TF-IDF scores because their IDF is lower, indicating they are not very informative.

Simple TF-IDF Calculation

► Consider a corpus of 3 documents:

1. Document 1: "I love programming in Python."
2. Document 2: "Python programming is fun."
3. Document 3: "I love coding in Python."

Step 1: Calculate TF for "Python" in Document 1

- Document 1: "I love programming in Python."
- The word "Python" appears **1 time** out of 5 total words.

$$\text{TF}(\text{Python}, \text{Document 1}) = \frac{1}{5} = 0.2$$

Simple TF-IDF Calculation

Step 2: Calculate IDF for "Python"

- The word "Python" appears in all 3 documents.
- The total number of documents is 3.
- The number of documents containing "Python" is 3.

$$\text{IDF}(\text{Python}) = \log\left(\frac{3}{3}\right) = \log(1) = 0$$

Step 3: Calculate TF-IDF for "Python" in Document 1

$$\text{TF-IDF}(\text{Python}, \text{Document 1}) = 0.2 \times 0 = 0$$

Since "Python" appears in all documents, its IDF value is 0, indicating that it is not a very distinctive word across the corpus. Therefore, its TF-IDF value is 0, meaning it doesn't help much in distinguishing Document 1 from others.

Applications of TF-IDF

- ▶ **Information Retrieval:** Ranking documents based on the relevance of their terms.
- ▶ **Text Classification:** Using the TF-IDF values of words as features to classify documents.
- ▶ **Clustering:** Grouping documents with similar content based on their TF-IDF scores.
- ▶ **Search Engines:** Determining the relevance of documents or web pages for a search query.

Word2Vec

- ▶ **Word2Vec** is a popular technique used to convert words into vector representations (embeddings) in a continuous vector space. It is based on neural networks and leverages the context in which words appear to learn dense representations, where similar words are placed closer together in the vector space.
- ▶ Word2Vec uses two primary models:
 1. **CBOW (Continuous Bag of Words)**: Predicts a target word based on context words surrounding it.
 2. **Skip-gram**: Predicts context words based on a target word.

Applications of Word2vec

- ▶ **Semantic Similarity:** Finding words that are similar in meaning or context (e.g., "king" and "queen").
- ▶ **Text Classification:** Converting words into vectors and using them as features for classification algorithms.
- ▶ **Recommendation Systems:** Recommending similar items (like words, products, etc.) based on vector similarity.
- ▶ **Named Entity Recognition (NER):** Using embeddings to recognize entities in text.

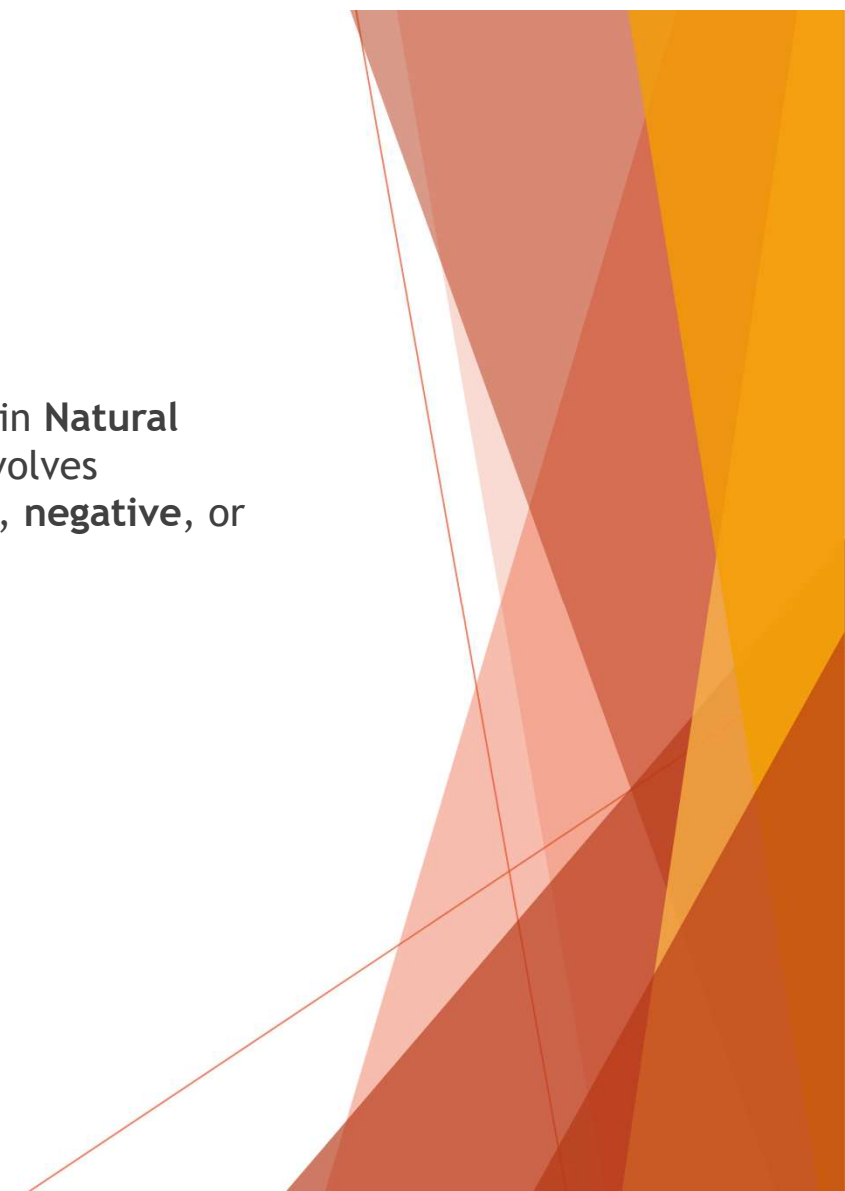
NLP Tasks

- ▶ Text Classification
- ▶ Named Entity Recognition (NER)
- ▶ Part-of-Speech (POS) Tagging
- ▶ Sentiment Analysis



Text Classification

- ▶ One of the most common applications of **text classification** in **Natural Language Processing (NLP)** is **sentiment analysis**, which involves categorizing text into different sentiments, such as **positive**, **negative**, or **neutral**.



Parts of Speech Tagging

- ▶ **Part of Speech (POS) Tagging** is a fundamental task in **Natural Language Processing (NLP)** that involves assigning a specific grammatical category (part of speech) to each word in a sentence. These categories include **nouns, verbs, adjectives, adverbs, pronouns, prepositions, conjunctions, and interjections**, among others.
- ▶ **Why is POS Tagging Important?**
 - It helps understand the syntactic structure of sentences.
 - It enables more advanced NLP tasks like named entity recognition (NER), machine translation, and information extraction.
 - It improves the performance of other NLP applications by providing syntactic information about words.

Parts of Speech Tagging

► Common POS Tags:

- **NN:** Noun, singular
- **NNS:** Noun, plural
- **VB:** Verb, base form
- **VBD:** Verb, past tense
- **JJ:** Adjective
- **RB:** Adverb
- **PRP:** Pronoun
- **IN:** Preposition

Application of Part of Speech Tagging

- ▶ Named Entity Recognition (NER): Identifying entities like names of people, organizations, and locations in text requires understanding parts of speech.
- ▶ Machine Translation: Understanding sentence structure is essential when translating text from one language to another.
- ▶ Information Extraction: Extracting meaningful information (e.g., dates, places) from unstructured text.
- ▶ Question Answering Systems: Identifying the type of words (nouns, verbs) in a question helps in providing accurate answers.

N-grams and Language Models

- ▶ An N-Gram is a contiguous sequence of n items (words, characters, etc.) from a given text or speech. In the context of Natural Language Processing (NLP), N-grams typically refer to sequences of words.
 - ▶ 1-Gram (Unigram): A single word. Example: "dog"
 - ▶ 2-Gram (Bigram): A sequence of two words. Example: "the dog"
 - ▶ 3-Gram (Trigram): A sequence of three words. Example: "the big dog"
 - ▶ n -Gram: A sequence of n words

Contextual Understanding: N-Grams capture contextual information by considering the relationships between neighboring words. For instance, "New York" (bigram) has different meaning than "New" and "York" individually (unigrams).

Application of N-Grams

- ▶ Text Generation: Generate new text by predicting the next word based on previous $n-1$ words.
- ▶ Speech Recognition: Use N-grams to predict the most likely sequence of words given a sequence of spoken words.
- ▶ Machine Translation: Translate sentences by leveraging the probability of word sequences.
- ▶ Spell Check: Identify the most likely correction for misspelled words based on surrounding words.

Text Generation using N-Grams

- ▶ Text generation using N-Grams involves creating a sequence of words based on the probabilities of word sequences learned from a training corpus. We can generate new text by predicting the next word based on the previous $n-1$ words.
- ▶ **Steps to Implement Text Generation Using N-Grams:**
 1. **Create a corpus:** Use a small example sentence or text data.
 2. **Generate N-Grams:** Build bigrams (or trigrams, etc.) from the corpus.
 3. **Text Generation:** Start with an initial word, then predict the next word using the bigrams.

NER

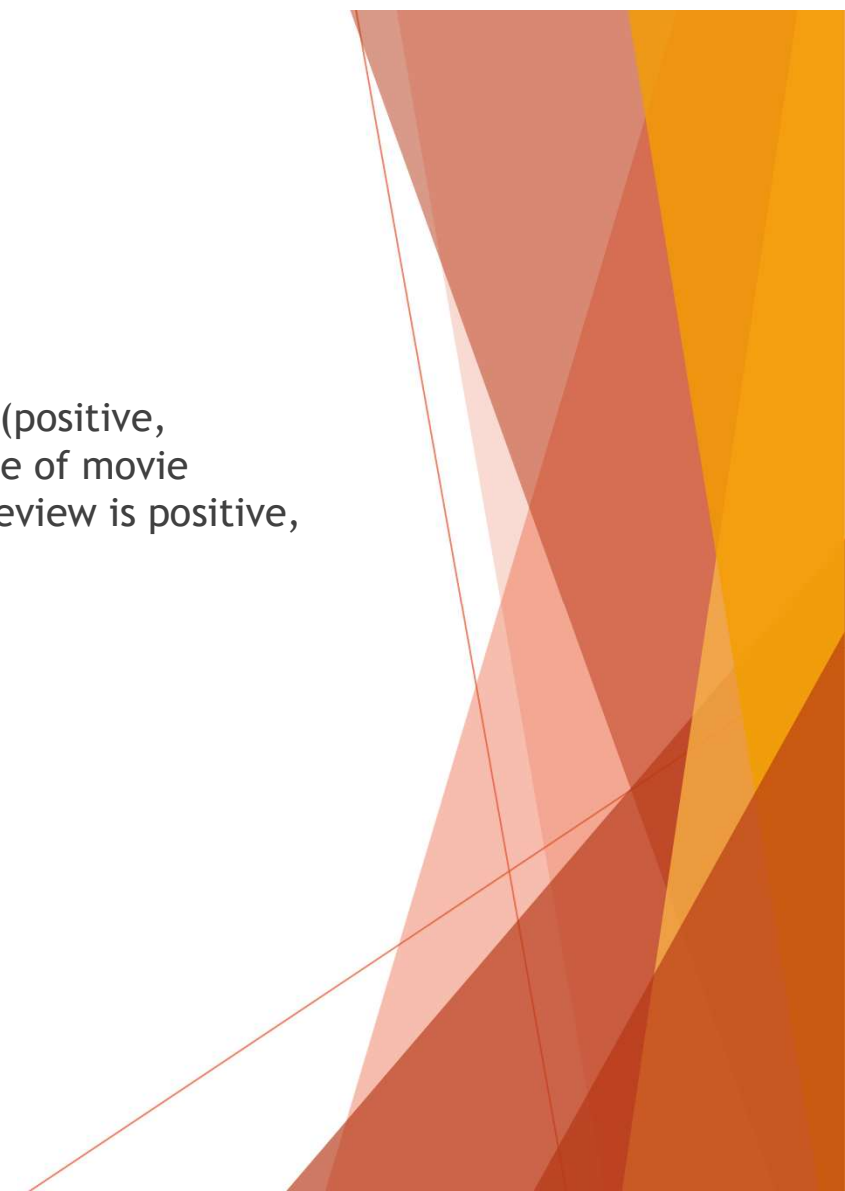
- ▶ **Named Entity Recognition (NER)** is a subtask of **Information Extraction (IE)** in Natural Language Processing (NLP). It involves identifying and classifying named entities in text into predefined categories, such as the names of persons, organizations, locations, dates, monetary values, etc.
- ▶ Let's say we have a sentence that contains various named entities, such as people's names, organizations, and locations. NER helps us extract these entities to better understand the structure of the sentence.

NER Example

- ▶ "Apple Inc. is planning to open a new store in Berlin on 15th January 2025. Tim Cook, the CEO of Apple, announced it in a press conference."
- ▶ In this example:
 - "Apple Inc." is an organization.
 - "Berlin" is a location.
 - "15th January 2025" is a date.
 - "Tim Cook" is a person.
 - "Apple" is another organization (same as before).

Sentiment Analysis

- ▶ Sentiment analysis is the task of determining the sentiment (positive, negative, or neutral) expressed in a piece of text. In the case of movie reviews, sentiment analysis can help determine whether a review is positive, negative, or neutral.



Applications of Sentiment Analysis

- ▶ **Customer Reviews:** Sentiment analysis is widely used in analyzing customer reviews on e-commerce platforms to gauge the overall satisfaction with products or services.
- ▶ **Social Media Monitoring:** Brands use sentiment analysis to track how people feel about their products or services on social media platforms.
- ▶ **Market Research:** It is used to understand public opinion about different brands, products, or political candidates.
- ▶ **Political Sentiment Analysis:** Helps in analyzing public opinion from speeches, debates, or social media to gauge the favorability of politicians or political decisions.

Evaluation Metrics for NLP

- ❖ **Accuracy** - Proportion of correct predictions (classification).
- ❖ **Precision** - Proportion of true positives among predicted positives (classification).
- ❖ **Recall** - Proportion of true positives among actual positives (classification).
- ❖ **F1-Score** - Harmonic mean of precision and recall (classification).
- ❖ **BLEU** - Measures the quality of text generation (e.g., machine translation).
- ❖ **ROUGE** - Measures overlap between n-grams in generated vs. reference text (summarization).
- ❖ **Perplexity** - Evaluates language models by measuring how well they predict a sample.
- ❖ **Confusion Matrix** - Visualizes performance with counts of TP, TN, FP, FN (classification).
- ❖ **AUC-ROC** - Measures classifier performance with various thresholds (binary classification).
- ❖ **Mean Squared Error (MSE)** - Measures the average squared difference (regression).