

Part 5 :

1. What is tokenization, and how does it differ from lemmatization and stemming?

Answer:

- **Tokenization:**

- **Definition:** Tokenization is the process of splitting a text into smaller units, called tokens. These tokens could be words, sentences, or subwords.
- **Example:** In the sentence "I love machine learning", tokenization splits it into `["I", "love", "machine", "learning"]`.

- **Lemmatization:**

- **Definition:** Lemmatization is the process of reducing a word to its base form (lemma) using a dictionary. It takes into account the **meaning of the word**.
- **Example:** "running" → "run", "better" → "good".

- **Stemming:**

- **Definition:** Stemming removes prefixes or suffixes to reduce a word to its **root form**, but **does not guarantee a meaningful word**.
- **Example:** "running" → "run", "happiness" → "happi".

Differences:

Feature	Tokenization	Lemmatization	Stemming
Purpose	Split text into tokens	Reduce words to their base form	Remove prefixes/suffixes to get root
Result	Words, sentences, or subwords	Meaningful base word	Root word (may not always be meaningful)

Performance	Fast but can break words into unwanted parts	Slower due to linguistic processing	Faster but less accurate
Use Case	Preprocessing step for NLP tasks	For tasks requiring accurate words	For tasks requiring fast processing

2. Explain the concept of Bag of Words (BoW) and its limitations.

Answer:

- **Bag of Words (BoW):**

- **Definition:** BoW is a text representation model where each document is represented as a vector of **word frequencies**.
- **Example:** Given two sentences:
 1. "I love programming."
 2. "Programming is fun."
- The vocabulary will be: `["I", "love", "programming", "is", "fun"]`
- The BoW representation will be:
 - Sentence 1: [1, 1, 1, 0, 0]
 - Sentence 2: [0, 0, 1, 1, 1]

- **Limitations:**

1. **No Context:** BoW does not capture word order or sentence structure.
2. **High Dimensionality:** It creates a large number of features (one for each word in the vocabulary), leading to sparse matrices.
3. **No Semantics:** BoW treats words as independent and does not capture synonyms or word meanings.

3. How does TF-IDF work, and how does it differ from simple word frequency?

Answer:

- **TF-IDF (Term Frequency-Inverse Document Frequency):**

- **Definition:** TF-IDF is a statistic used to measure the importance of a word in a document relative to the entire corpus.
- **Components:**
 - **Term Frequency (TF):** Measures **how frequently a term** appears in a document.
 - **Inverse Document Frequency (IDF):** Measures **how rare or common a term** is across all documents in the corpus.
- **Difference from Word Frequency:**

Feature	TF (Word Frequency)	TF-IDF
Focus	Counts word frequency in a document	Adjusts frequency by importance across all documents
Advantage	Simple and fast	Weighs more important words more accurately
Limitation	Ignores context and document distribution	Can still give too much weight to rare terms

4. What is word embedding, and why is it important in NLP?

Answer:

- **Word Embedding:**
 - **Definition:** Word embedding is the representation of words in a continuous **vector space**, where similar words are placed close to each other in the vector space.
 - **Examples:** **Word2Vec, GloVe, FastText.**
 - **How it works:** Words are mapped to vectors, capturing semantic meaning. For instance, "king" and "queen" might be close in the vector space due to shared contexts.
- **Importance:**

1. **Captures Semantic Meaning:** Unlike BoW, embeddings capture semantic relationships.
 2. **Reduces Dimensionality:** They map words to a dense vector, reducing the feature space size.
 3. **Improves Performance:** Embeddings help improve model accuracy in NLP tasks like classification, translation, and sentiment analysis.
-

5. What are some common applications of NLP in real-world systems?

Answer:

- **Sentiment Analysis:** Analyzing text to determine the sentiment (positive, negative, or neutral).
 - **Chatbots and Virtual Assistants:** Siri, Alexa, and Google Assistant use NLP to process and respond to user queries.
 - **Language Translation:** Tools like Google Translate use NLP to translate text between languages.
 - **Text Summarization:** Extractive and abstractive summarization to generate shorter versions of documents.
 - **Named Entity Recognition (NER):** Identify entities like names, dates, and locations in text for information extraction.
-

6. What is Named Entity Recognition (NER), and where is it applied?

Answer:

- **Named Entity Recognition (NER):**
 - **Definition:** NER is the process of **identifying and classifying entities** in text into predefined categories such as persons, organizations, locations, dates, etc.
 - **Example:** In the sentence "Barack Obama was born in Hawaii in 1961", NER would identify:

- "Barack Obama" as a **Person**
 - "Hawaii" as a **Location**
 - "1961" as a **Date**
 - **Applications:**
 1. **Information Retrieval:** Enhances search engines by allowing them to recognize and prioritize important entities.
 2. **Medical Text Mining:** Extracts information like drug names, diseases, etc., from clinical notes.
 3. **Financial Analysis:** Identifies companies, stock symbols, and financial terms in news articles and reports.
-

7. What are transformers in NLP, and how have they changed the field?

Answer:

- **Transformers:**
 - **Definition:** Transformers are a type of deep learning architecture based on self-attention mechanisms. They allow models to focus on different parts of a sentence when making predictions, regardless of word order.
 - **Example:** BERT, GPT, and T5 are transformer models.
 - **Impact:**
 1. **Handling Long-Range Dependencies:** Transformers can consider the entire context of a sentence at once, which improves performance in tasks like translation and summarization.
 2. **Pretrained Models:** Models like BERT and GPT can be fine-tuned for specific tasks, leading to better generalization and reduced training time.
-

8. How do you handle out-of-vocabulary (OOV) words in NLP models?

Answer:

- **OOV Words:**
 - **Definition:** Words that do not appear in the vocabulary of the model.
- **Handling Techniques:**
 1. **Subword Tokenization:** Break OOV words into smaller subword units (e.g., BPE, SentencePiece).
 2. **Using Special Tokens:** Assign a special token (like `<UNK>`) to represent OOV words.
 3. **Embedding Models:** Use word embeddings like FastText that can generate representations for OOV words based on their subword components.

This is a detailed response to the key questions related to NLP, with points, comparisons, and applications to help in an interview setting. Diagrams and tables can be added where required to visualize certain concepts, especially for TF-IDF, LDA, and transformers.

4o mini

9. What is a language model, and how do you evaluate its performance?

Answer:

- **Language Model:**
 - **Definition:** A language model (LM) is a probabilistic model used to predict the next word in a sequence or to calculate the likelihood of a sentence or phrase. It can be unidirectional (predicts next word) or bidirectional (considers context from both directions, as in BERT).
 - **Types:**
 1. **Statistical Language Models:** Like n-grams.
 2. **Neural Language Models:** Such as LSTM, GRU, Transformer models.

- **Evaluation:**

1. **Perplexity:** Measures how well the model predicts a sample. Lower perplexity means better prediction.
 2. **Accuracy:** Percentage of correct word predictions in a given text.
 3. **BLEU Score:** For machine translation, evaluates how closely a machine-generated translation matches human reference translations.
 4. **Log-Likelihood:** Measures how probable a given sequence of words is, given the model's parameters.
-

10. How does Part-of-Speech (POS) tagging work, and why is it important in NLP?

Answer:

- **POS Tagging:**

- **Definition:** POS tagging is the process of labeling each word in a sentence with its **grammatical category** (e.g., noun, verb, adjective).
- **Example:** In the sentence "The cat sleeps", the tags would be:
 - "The" → Determiner (DET)
 - "cat" → Noun (NN)
 - "sleeps" → Verb (VBZ)

- **How it works:**

- **Rule-based methods:** Use predefined grammar rules.
- **Statistical methods:** Use machine learning models trained on labeled corpora.
- **Neural models:** Use deep learning, such as BiLSTM-CRF models.

- **Importance:**

1. **Syntax Understanding:** Helps in understanding sentence structure.
2. **Named Entity Recognition (NER):** Often used as a preprocessing step for NER tasks.

3. **Improving Sentiment Analysis:** POS tagging can help in identifying the sentiment of specific words.
-

11. How do n-grams help in language modeling?

Answer:

- **N-grams:**
 - **Definition:** N-grams are contiguous sequences of n items from a given sample of text or speech. In language modeling, n-grams represent sequences of n words.
 - **Unigrams:** Single words (e.g., "I", "love", "Python").
 - **Bigrams:** Pairs of consecutive words (e.g., "I love", "love Python").
 - **Trigrams:** Triplets of words (e.g., "I love Python").
 - **Role in Language Modeling:**
 1. **Captures Local Context:** Helps predict the next word in a sequence based on previous words.
 2. **Simplicity:** N-grams are easier to implement compared to more complex models like neural networks.
 3. **Limitations:** Higher-order n-grams (e.g., 5-grams) can lead to sparse data issues, requiring smoothing techniques.
 - **Evaluation:**

N-grams are evaluated based on **perplexity** and **log-likelihood**, as discussed earlier.
-

12. Explain the concept of syntactic parsing in NLP and its applications.

Answer:

- **Syntactic Parsing:**

- **Definition:** Syntactic parsing involves analyzing the structure of a sentence to establish relationships between words based on grammar. It produces a parse tree that represents sentence structure.
 - **Types:**
 1. **Constituency Parsing:** Breaks the sentence into sub-phrases (e.g., noun phrase, verb phrase).
 2. **Dependency Parsing:** Focuses on the grammatical relationships between words.
 - **Applications:**
 1. **Machine Translation:** Understanding the structure of sentences in both languages.
 2. **Question Answering:** Extracting relevant information from text based on its grammatical structure.
 3. **Sentiment Analysis:** Helps identify key sentence components like adjectives or verbs to determine sentiment.
-

13. What is the purpose of stopword removal in NLP?

Answer:

- **Stopword Removal:**
 - **Definition:** Stopwords are common words that typically do not add significant meaning in text analysis (e.g., "the", "is", "in").
 - **Purpose:**
 1. **Reduce Noise:** By removing stopwords, the model focuses on more meaningful words.
 2. **Improve Efficiency:** Reduces the size of the data, making processing faster.
- **Challenges:**
 1. **Context Sensitivity:** Some stopwords may carry meaning in certain contexts.

2. **Language Dependency:** Stopwords vary from language to language.

14. What are the differences between shallow and deep NLP models?

Answer:

Aspect	Shallow NLP Models	Deep NLP Models
Approach	Relies on hand-crafted features (e.g., POS tags, n-grams)	Utilizes deep learning models (e.g., RNN, LSTM, Transformers)
Data Requirement	Works well with limited data	Requires large datasets and computational resources
Performance	Performs well for simple tasks (e.g., text classification)	Performs better for complex tasks (e.g., machine translation)
Model Complexity	Less complex and easier to interpret	Complex, less interpretable, but more powerful
Examples	Logistic regression, Naive Bayes	BERT, GPT, Transformer-based models