

## 1. Explain One-Hot Encoding.

**Answer:** one-hot encoding is often used to represent words or tokens in a text corpus. Each unique word in the vocabulary is assigned a unique integer index. Then, each word is represented as a binary vector, where the index corresponding to the word's position in the vocabulary is set to 1, and all other indices are set to 0.

For example, consider a vocabulary containing the words: ["cat", "dog", "bird"]. Using one-hot encoding, each word would be represented as follows:

- "cat" = [1, 0, 0]
- "dog" = [0, 1, 0]
- "bird" = [0, 0, 1]

So, if you have a sentence like "The cat chased the bird", you would represent it using one-hot encoding as a sequence of binary vectors:

- "The" = [1, 0, 0] (assuming it's in the vocabulary)
- "cat" = [1, 0, 0]
- "chased" = [0, 0, 0] (not in the vocabulary)
- "the" = [0, 0, 0] (not in the vocabulary)
- "bird" = [0, 0, 1]

One-hot encoding in NLP has some limitations, such as the resulting vectors being very large (especially for large vocabularies) and not capturing any semantic information about the words. However, it's a simple and effective way to represent text data for certain tasks, especially when combined with techniques like bag-of-words or bag-of-ngrams.

## 2. Explain Bag of Words.

**Answer:** The Bag of Words (BoW) model is a simple yet powerful technique used in natural language processing (NLP) for text analysis and feature extraction. It represents text data as a collection of words, disregarding grammar and word order, and focusing solely on word occurrence.

Here's how the Bag of Words model works:

1. **Tokenization:** First, the text data is tokenized, which involves splitting it into individual words or tokens. Punctuation and other non-alphanumeric characters are typically removed or treated as separate tokens.
2. **Vocabulary Creation:** Next, a fixed vocabulary is constructed from the unique words in the entire corpus of text data. Each unique word becomes a feature in the vocabulary.
3. **Vectorization:** For each document or text sample in the corpus, a vector representation is created based on the frequency of words in the vocabulary. Each element in the vector corresponds to a word in the vocabulary, and its value represents the frequency of that word in the document. Alternatively, it can represent whether the word is present (binary representation) in the document (1 for presence, 0 for absence).
4. **Normalization (optional):** Optionally, the vectors can be normalized to ensure that longer documents don't have an advantage due to having more words.

5. **Analysis:** Once the documents are represented as vectors, various machine learning algorithms can be applied for tasks such as classification, clustering, or sentiment analysis.

For example, consider the following two sentences:

- Sentence 1: "The cat sat on the mat."
- Sentence 2: "The dog played in the yard."

After tokenization and vocabulary creation, the vocabulary might contain: ["the", "cat", "sat", "on", "mat", "dog", "played", "in", "yard"].

Then, the Bag of Words representation for each sentence could be:

- Sentence 1: [1, 1, 1, 1, 1, 0, 0, 0, 0]
- Sentence 2: [1, 0, 0, 0, 0, 1, 1, 1, 1]

The Bag of Words model is widely used due to its simplicity and effectiveness, especially in tasks like text classification, where the focus is on keyword presence rather than context or word order. However, it does not capture semantics, word relationships, or word order, which can limit its applicability in certain NLP tasks.

### 3. Explain Bag of N-Grams.

**Answer:** The Bag of N-Grams model is an extension of the Bag of Words (BoW) model in natural language processing (NLP) that considers sequences of words, called n-grams, instead of just individual words. An n-gram is a contiguous sequence of n items (which are typically words in the context of NLP).

Here's how the Bag of N-Grams model works:

1. **Tokenization and N-Gram Generation:** Similar to the BoW model, the text data is tokenized into individual words. Then, instead of considering only single words, the model generates all possible sequences of n contiguous words from the text.
2. **Vocabulary Creation:** A fixed vocabulary is created from the unique n-grams in the entire corpus of text data. Each unique n-gram becomes a feature in the vocabulary.
3. **Vectorization:** For each document or text sample in the corpus, a vector representation is created based on the frequency of n-grams in the vocabulary. Each element in the vector corresponds to an n-gram in the vocabulary, and its value represents the frequency of that n-gram in the document. Alternatively, it can represent whether the n-gram is present (binary representation) in the document (1 for presence, 0 for absence).
4. **Normalization (optional):** Optionally, the vectors can be normalized to ensure that longer documents don't have an advantage due to having more words.
5. **Analysis:** Once the documents are represented as vectors, various machine learning algorithms can be applied for tasks such as classification, clustering, or sentiment analysis.

For example, consider the following sentence:

- Sentence: "The cat sat on the mat."

If we use bi-grams (2-grams) for this sentence, the generated bi-grams would be:

- ["The cat", "cat sat", "sat on", "on the", "the mat"]

Then, the Bag of N-Grams representation for this sentence could be:

- [1, 1, 1, 1, 1]

The Bag of N-Grams model captures more contextual information compared to the Bag of Words model because it considers sequences of words. However, it suffers from the curse of dimensionality, especially with larger values of  $n$ , as the number of possible  $n$ -grams increases exponentially with the length of the  $n$ -gram. Additionally, like the BoW model, it does not capture word order beyond the  $n$ -gram size chosen.

#### 4. Explain TF-IDF.

**Answer:** In natural language processing (NLP), TF-IDF is a widely used technique for representing text data numerically, which helps in various tasks such as document classification, clustering, and information retrieval. Here's how TF-IDF is applied in the context of NLP:

1. **Document-Term Matrix:** TF-IDF is used to convert a collection of documents into a numerical matrix representation called the document-term matrix. Each row in this matrix corresponds to a document, and each column corresponds to a unique term (usually words) in the entire corpus.
2. **Term Frequency (TF):** For each document, the term frequency (TF) is calculated, which measures how often each term appears in the document. TF is calculated as the number of times a term appears in a document divided by the total number of terms in that document. It normalizes the frequency of terms within each document.
3. **Inverse Document Frequency (IDF):** For each term, the inverse document frequency (IDF) is calculated, which measures how important the term is across the entire corpus. IDF is calculated as the logarithm of the total number of documents divided by the number of documents containing the term, with a smoothing factor to prevent division by zero.
4. **TF-IDF Calculation:** The TF-IDF score for each term-document pair is calculated by multiplying the term frequency (TF) by the inverse document frequency (IDF). This score reflects the importance of each term in each document relative to the entire corpus.
5. **Normalization (optional):** Optionally, the TF-IDF vectors can be normalized to ensure that longer documents don't have an advantage due to having more terms.
6. **Analysis:** Once the documents are represented as TF-IDF vectors, various machine learning algorithms can be applied for tasks such as text classification, clustering, sentiment analysis, and information retrieval. TF-IDF helps in identifying the most relevant terms in each document and capturing the semantic meaning of the text.

Overall, TF-IDF is a powerful tool in NLP for transforming raw text data into numerical features that can be used by machine learning algorithms, enabling efficient and effective text analysis and processing.

#### 5. What is OOV problem?

**Answer:** The OOV (Out-Of-Vocabulary) problem refers to the situation in natural language processing (NLP) where a word or token encountered during testing or inference is not present in the vocabulary or training data of a model.

Here's why the OOV problem can occur:

1. **Limited Vocabulary:** NLP models are typically trained on a fixed vocabulary, which consists of the words seen during training. If a word appears during testing that is not present in the vocabulary, it becomes an out-of-vocabulary word.
2. **Data Sparsity:** In many NLP tasks, the vocabulary can be vast, and it's impractical to include every possible word during training. This leads to data sparsity, where rare or infrequent words may not be seen often enough during training to be included in the vocabulary.
3. **Word Variation:** Languages exhibit variations due to morphology, spelling errors, or creative language use. These variations can lead to unseen words or variations of known words that are not in the vocabulary.

The OOV problem can negatively impact the performance of NLP models, especially in tasks like machine translation, text generation, and sentiment analysis, where understanding the entire vocabulary is crucial for accurate predictions.

There are several strategies to address the OOV problem:

1. **Fallback Mechanisms:** Use fallback mechanisms such as replacing OOV words with a special token or with a generic placeholder like "<UNK>" (unknown token).
2. **Subword Tokenization:** Use subword tokenization techniques such as Byte Pair Encoding (BPE) or WordPiece, which can handle unseen words by decomposing them into smaller subword units that are present in the vocabulary.
3. **Character-Level Models:** Train models at the character level instead of the word level, allowing them to handle unseen words by generalizing based on character patterns.
4. **Dynamic Vocabulary:** Dynamically expand the vocabulary during training or inference to include new words encountered during testing. This approach requires the model to be able to adapt and learn from new vocabulary items.

Addressing the OOV problem is essential for building robust and effective NLP models that can handle real-world text data where unseen words are common.

## 6. What are word embeddings?

**Answer:** Word embeddings are numerical representations of words in a continuous vector space, where words with similar meanings are represented by vectors that are close to each other in the space. These representations capture semantic relationships between words and are widely used in natural language processing (NLP) tasks.

Here's how word embeddings work:

1. **Dense Vector Representations:** Unlike traditional one-hot encoding or sparse representations, word embeddings represent words as dense vectors of real numbers, typically with hundreds of dimensions. Each dimension in the vector corresponds to a feature or aspect of the word's meaning.
2. **Learned from Data:** Word embeddings are learned from large text corpora using unsupervised learning techniques, such as neural networks. During training, the model learns to predict words in context based on their surrounding words or sentences. As a result, the vectors are learned in a way that similar words have similar vector representations.
3. **Semantic Similarity:** Because word embeddings are learned based on the context in which words appear, words with similar meanings or contexts tend to have similar vector representations. For example, in a well-trained word embedding model, the vectors for "king" and "queen" might be close together because they often appear in similar contexts.
4. **Vector Arithmetic:** Word embeddings exhibit interesting properties that allow for operations such as vector arithmetic. For example, the vector representation of "king" - "man" + "woman" might be close to the vector representation of "queen". This property allows for analogical reasoning and semantic calculations within the vector space.

Word embeddings have revolutionized the field of NLP and have become a fundamental building block for many NLP tasks, including sentiment analysis, machine translation, document classification, and named entity recognition. They provide a dense and continuous representation of words that captures rich semantic information, enabling more effective and accurate modeling of language. Popular word embedding algorithms include Word2Vec, GloVe, and fastText.

## 7. Explain Continuous bag of words (CBOW).

**Answer:** Continuous Bag of Words (CBOW) is a type of word embedding model used in natural language processing (NLP) and neural network-based language models. CBOW aims to predict a target word based on its surrounding context words within a fixed window size.

Here's how the Continuous Bag of Words (CBOW) model works:

1. **Input Representation:** In CBOW, the input to the model is a context window of surrounding words centered around a target word. The context window size is a hyperparameter that determines the number of words to consider on each side of the target word.
2. **Word Embeddings:** Each word in the context window is represented by a word embedding, which is a dense vector representation learned during the training process. These word embeddings are typically initialized randomly or using pre-trained embeddings like Word2Vec or GloVe.
3. **Aggregation:** The word embeddings of the context words are aggregated or averaged to create a single vector representation, which serves as the input to the neural network.
4. **Neural Network Architecture:** The aggregated vector representation is fed into a neural network with one or more hidden layers. The network is trained to predict the target word based on this aggregated representation. The output layer typically consists of softmax units, with each unit corresponding to a word in the vocabulary. The model is trained using a classification objective, where the target word is treated as a label, and the model is trained to minimize the cross-entropy loss between the predicted probabilities and the true label.

5. **Training:** During training, the parameters of the neural network, including the word embeddings, are updated using backpropagation and gradient descent optimization algorithms such as stochastic gradient descent (SGD) or Adam.
6. **Word Embedding Learning:** As the model is trained on a large corpus of text data, the word embeddings are adjusted to capture semantic relationships between words. Words that appear in similar contexts will have similar vector representations, allowing the model to capture semantic similarity and perform tasks such as word analogy and semantic similarity.

CBOW is known for its simplicity and efficiency, especially for smaller datasets and frequent words. However, it may not capture complex semantic relationships as effectively as other models like Skip-gram, which predicts context words given a target word. Nonetheless, CBOW remains a popular choice for word embedding generation, especially when computational resources are limited.

## 8. Explain SkipGram.

**Answer:** Skip-gram is another type of word embedding model used in natural language processing (NLP) and neural network-based language models, particularly for learning distributed representations of words. Unlike Continuous Bag of Words (CBOW), which predicts a target word based on its context, Skip-gram predicts the context words (surrounding words) given a target word.

Here's how the Skip-gram model works:

1. **Input Representation:** The input to the Skip-gram model is a target word from the text corpus. Unlike CBOW, which uses a window of surrounding words, Skip-gram focuses on a single target word at a time.
2. **Word Embeddings:** Each word in the vocabulary is represented by a dense vector, called a word embedding. These word embeddings are learned during the training process and capture the semantic meaning of words based on their context in the corpus.
3. **Context Prediction:** Given a target word, the Skip-gram model is trained to predict the surrounding context words within a fixed window size. The context words are sampled from the same sentence or text window containing the target word.
4. **Neural Network Architecture:** The Skip-gram model typically uses a shallow neural network with a single hidden layer. The input layer consists of the one-hot encoded representation of the target word, while the output layer consists of softmax units corresponding to the vocabulary words. The model is trained to minimize the cross-entropy loss between the predicted probabilities of context words and the true context words.
5. **Training:** During training, the parameters of the neural network, including the word embeddings, are updated using backpropagation and optimization algorithms such as stochastic gradient descent (SGD) or Adam. The goal is to adjust the word embeddings to accurately predict the context words for each target word in the corpus.
6. **Word Embedding Learning:** As the Skip-gram model is trained on a large corpus of text data, the word embeddings are learned to capture semantic relationships between words. Words that appear in similar contexts will have similar vector representations, allowing the model to capture semantic similarity and perform tasks such as word analogy and semantic similarity.

Skip-gram is known for its ability to capture complex semantic relationships and perform well on larger datasets. However, it may require more computational resources and training time compared to CBOW, especially for infrequent words. Nonetheless, Skip-gram remains a popular choice for word embedding generation, particularly when capturing fine-grained semantic information is important.

## 9. Explain GloVe Embeddings.

**Answer:** GloVe (Global Vectors for Word Representation) is a word embedding model used in natural language processing (NLP) tasks to learn vector representations of words. Unlike models like Word2Vec, which focus on predicting context words given a target word (Skip-gram) or predicting a target word given context words (CBOW), GloVe aims to learn word embeddings by capturing global word co-occurrence statistics from a corpus.

Here's how the GloVe model works:

1. **Word Co-occurrence Matrix:** GloVe starts by constructing a co-occurrence matrix  $XX$ , where  $X_{ij}$  represents the number of times word  $i$  co-occurs with word  $j$  in the corpus. The co-occurrence counts can be computed using a fixed window size around each word in the corpus.
2. **Objective Function:** GloVe seeks to learn word embeddings such that the dot product of two word vectors captures the log of their co-occurrence probability. The objective function of GloVe is defined as follows:

$$J = \sum_{i,j=1}^V f(X_{ij})(\mathbf{w}_i^T \mathbf{w}_j + b_i + b_j - \log(X_{ij}))^2 \quad J = \sum_{i,j=1}^V f(X_{ij})(\mathbf{w}_i^T \mathbf{w}_j + b_i + b_j - \log(X_{ij}))^2$$

Where:

- $V$  is the size of the vocabulary.
  - $\mathbf{w}_i$  and  $\mathbf{w}_j$  are the word vectors for word  $i$  and  $j$ , respectively.
  - $b_i$  and  $b_j$  are the bias terms.
  - $f(X_{ij})$  is a weighting function that adjusts the influence of co-occurrence counts. It's typically defined as  $f(x) = \min(1, (x/x_{\max})^\alpha)$ , where  $x_{\max}$  is a threshold and  $\alpha$  is a hyperparameter.
3. **Optimization:** The objective function  $J$  is optimized using gradient descent or other optimization algorithms to learn the word vectors  $\mathbf{w}_i$  and  $\mathbf{w}_j$ . During training, the model adjusts the word vectors to minimize the difference between the dot product of word vectors and the log of their co-occurrence counts.
  4. **Word Embeddings:** Once trained, the word embeddings learned by GloVe capture semantic relationships between words based on their co-occurrence statistics in the corpus. Words with similar meanings or contexts will have similar vector representations.

GloVe embeddings are known for their ability to capture global word co-occurrence statistics and produce high-quality word representations that excel in tasks like word analogy, word similarity, and semantic reasoning. They are widely used in various NLP applications, including sentiment analysis, machine translation, and document classification. GloVe embeddings are often pre-trained on large text corpora and then fine-tuned or used as fixed features in downstream NLP tasks.