

# Explanation of NLP Embedding Methods

## Word Embeddings:

- **Definition:** Word embeddings are a technique in Natural Language Processing (NLP) used to represent words as dense, numerical vectors in a high-dimensional space. These vectors capture the semantic relationships between words, meaning that words with similar meanings are located close to each other in this vector space.
- **Purpose:**
  - Convert words into a format that machine learning models can understand.
  - Capture semantic relationships (e.g., "king" is similar to "queen" in some ways, and "man" is similar to "woman").
  - Reduce the dimensionality of text data compared to one-hot encoding.
- **Why they are important:** Traditional methods of representing words, such as one-hot encoding, result in sparse vectors and do not capture semantic relationships. Word embeddings address these limitations, enabling NLP models to perform better on tasks like sentiment analysis, text classification, and machine translation.

## Key Word Embedding Techniques:

- **Word2Vec:**
  - A predictive model that learns word embeddings by predicting the surrounding words (context) given a target word, or vice versa.
  - **CBOW (Continuous Bag of Words):** Predicts a target word from its surrounding context words.
  - **Skip-gram:** Predicts context words given a target word.
- **GloVe (Global Vectors for Word Representation):**
  - A count-based model that leverages global word co-occurrence statistics from a corpus.
  - It constructs a word co-occurrence matrix, which captures how frequently words appear together.
- **FastText:**
  - An extension of Word2Vec that takes subword information into account.
  - It represents words as bags of character n-grams, allowing it to handle out-of-vocabulary words and capture morphological information.
- **TF-IDF (Term Frequency-Inverse Document Frequency):**
  - A numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.
  - It is often used as a weighting factor in searches of information retrieval, text

mining, and user modeling. The TF-IDF value increases proportionally to the number of times a word appears in the document, but is offset by the frequency of the word in the corpus, which helps to control for the fact that some words are generally more common than others.

### **Pre-trained Word Embeddings:**

- Large pre-trained word embeddings can be downloaded, allowing you to use them in your models.
- **Benefits:**
  - Save significant training time and computational resources.
  - Often capture more general language semantics.

### **Embedding Layer in Neural Networks:**

- An embedding layer is a neural network layer that converts discrete data, such as word indices, into continuous vector representations (word embeddings).
- It is commonly used as the first layer in NLP models.

### **Characteristics of Good Word Embeddings:**

- Capture semantic relationships: Words with similar meanings should have similar vector representations.
- Low dimensionality: Should be dense vectors with a manageable number of dimensions.
- Contextual relationships: Should ideally capture how word meaning varies across different contexts.

### **Limitations of Traditional Word Embeddings (Word2Vec, GloVe):**

- Context-Insensitivity: Traditional word embeddings like Word2Vec and GloVe assign a single vector to each word, regardless of the context in which it appears. This means they struggle to capture the different meanings of polysemous words (words with multiple meanings).