
Graduate Certificate in Adopting AI for Infection Prevention and Control

Natural Language Processing and Text Mining

Bag of Words: A simple yet effective way to analyze text data in Natural Language Processing (NLP) and Text Mining. It involves treating each document as a collection of words, without considering the grammar or word order. The frequency of each word in the document is then used as a feature for analysis.

Corpus: A large collection of texts that are used as a sample for analysis in NLP and Text Mining. A corpus can be made up of texts from a specific domain, such as medical articles or social media posts, and can contain millions of documents.

Data Cleaning: The process of preparing text data for analysis by removing irrelevant information, such as punctuation, numbers, and stop words (common words like "and" and "the"). Data cleaning is an important step in NLP and Text Mining as it can improve the accuracy of the analysis.

Feature Extraction: The process of converting raw text data into numerical features that can be used for analysis in NLP and Text Mining. This can include techniques such as Bag of Words, TF-IDF (Term Frequency-Inverse Document Frequency), and word embeddings.

Inverse Document Frequency (IDF): A numerical value that reflects the importance of a word in a corpus. Words that appear in many documents have a low IDF, while words that appear in few documents have a high IDF. IDF is often used in conjunction with Bag of Words and TF-IDF to weight the features in NLP and Text Mining.

Latent Dirichlet Allocation (LDA): A type of topic modeling used in NLP and Text Mining to automatically identify topics in a corpus. LDA assumes that each document is a mixture of topics, and each topic is a probability distribution over words.

Natural Language Processing (NLP): A field of computer science that deals with the interaction between computers and human language. NLP techniques are used to analyze, understand, and generate human language in a form that computers can understand.

Named Entity Recognition (NER): A type of NLP task that involves identifying and classifying named entities, such as people, organizations, and locations, in text data. NER is often used in Text Mining to extract structured information from unstructured text.

Part-of-Speech (POS) Tagging: A type of NLP task that involves identifying the part of speech, such as nouns, verbs, and adjectives, of each word in a sentence. POS tagging is often used as an intermediate step in NLP and Text Mining to improve the accuracy of other NLP tasks.

Sentiment Analysis: A type of NLP task that involves automatically determining the sentiment or attitude expressed in text data. Sentiment analysis is often used in Text Mining to analyze customer feedback, social media posts, and other forms of text data to gain insights into public opinion.

Stemming: A technique used in NLP and Text Mining to reduce words to their base or root form. For example, the words "running," "runs," and "ran" can be reduced to the root word "run" through stemming.

Stop Words: Common words like "and" and "the" that are often removed from text data during the data cleaning process in NLP and Text Mining. Stop words are usually removed because they do not carry much meaning and can add noise to the analysis.

Support Vector Machine (SVM): A type of machine learning algorithm that is often used in NLP and Text Mining for classification tasks. SVMs work by finding the optimal boundary or hyperplane that separates data points into different classes.

Term Frequency (TF): A numerical value that reflects the frequency of a word in a document. TF is often used in conjunction with IDF to create TF-IDF features in NLP and Text Mining.

Text Mining: The process of extracting valuable insights and knowledge from text data through the use of NLP techniques and machine learning algorithms. Text Mining is often used in fields such as healthcare, finance, and marketing to gain insights into customer behavior, public opinion, and other important trends.

Tokenization: The process of breaking text data into smaller units, such as words or sentences, for analysis in NLP and Text Mining. Tokenization is often the first step in NLP and Text Mining, as it allows for the analysis of individual words and phrases.

Topic Modeling: A type of NLP task that involves automatically identifying topics in a corpus. Topic modeling algorithms, such as LDA, assume that each document is a mixture of topics and use statistical methods to identify the topics and their distributions in the corpus.

Transfer Learning: A technique used in NLP and Text Mining to improve the performance of machine learning models by using pre-trained models on large datasets. Transfer learning is often used to improve the performance of deep learning models, such as word embeddings and transformers.

Transformers: A type of deep learning model used in NLP and Text Mining for tasks such as machine translation, sentiment analysis, and question answering. Transformers use self-attention mechanisms to weigh the importance of each word in a sentence, allowing for more accurate analysis.

Unstructured Data: Data that does not have a pre-defined format or structure, such as text data. Unstructured data is often difficult to analyze using traditional methods, and NLP and Text Mining techniques are often used to extract insights from unstructured data.

Word Embeddings: A type of feature extraction technique used in NLP and Text Mining to represent words

as numerical vectors. Word embeddings capture semantic relationships between words, such as similarity and relatedness, and are often used in deep learning models.

Word Sense Disambiguation (WSD): A type of NLP task that involves determining the meaning or sense of a word in a specific context. WSD is often used in Text Mining to improve the accuracy of other NLP tasks, such as sentiment analysis and information extraction.

In conclusion, NLP and Text Mining are powerful techniques for analyzing and extracting insights from text data. From Bag of Words to Word Sense Disambiguation, the glossary terms covered in this response provide a solid foundation for understanding the key concepts and techniques used in NLP and Text Mining. Whether you're a student, researcher, or practitioner, mastering these concepts will help you unlock the full potential of text data and gain valuable insights into the world around us.