# NLP - ISE 2

## Semantic Analysis

Semantic analysis is the process of extracting meaning from text and natural language. It involves understanding the relationships between words and concepts, and how they are used to convey meaning in different contexts.

Semantic analysis is a complex task because it requires a deep understanding of language, including its structure, syntax, and semantics. However, it is also an essential task for many natural language processing (NLP) applications, such as machine translation, question answering, and text summarization.

There are two main types of semantic analysis: supervised and unsupervised.

**Supervised semantic analysis** uses labeled data to train a model to predict the meaning of text. The labeled data typically consists of text samples that have been annotated with semantic labels, such as the sentiment of a text or the topic of a text.

Once the model is trained, it can be used to predict the meaning of new text samples. For example, a supervised semantic analysis model could be used to predict the sentiment of a product review or to identify the topic of a news article.

Examples - Machine Translation, Question Answering and Text Summarization

**Unsupervised semantic analysis** does not use labeled data. Instead, it relies on statistical methods to learn patterns in the data and extract meaning from text.

Unsupervised semantic analysis is often used to discover new relationships between words and concepts. For example, an unsupervised semantic analysis model could be used to identify clusters of words that are related to each other, or to identify the hidden meaning of a text.

Examples - Topic Modeling, Sentiment Analysis and word sense disambiguation


**Definite and indefinite noun phrases**

Definite and indefinite noun phrases are two types of noun phrases that differ in how they refer to the noun they modify.

**Definite noun phrases** refer to a specific noun that is known to the speaker and/or listener. They are typically marked by the presence of the definite article "the". For example:

- The cat is sitting on the mat.

- I saw the president yesterday.

- The book I am reading is very interesting.

**Indefinite noun phrases** refer to a general noun or a noun that is not known to the speaker and/or listener. They are typically marked by the presence of an indefinite article ("a" or "an") or by a numeral. For example:

- A cat is sitting on the mat.

- I saw a president yesterday.

- I am reading a book.

Definite and indefinite noun phrases can be used in a variety of ways. For example, they can be used to introduce new information, to refer to previously mentioned information, or to make generalizations.

## Word Sense Disambiguation

**Word sense disambiguation (WSD)** is the process of identifying which sense of a word is meant in a sentence or other segment of context. Word sense disambiguation, in natural language processing (NLP), may be defined as the ability to determine which meaning of word is activated by the use of word in a particular context. For example, the word "bank" has multiple senses, such as the side of a river, a financial institution, or a steep slope. In the sentence "I went to the bank to deposit a check," the sense of "bank" is clearly the financial institution.

WSD is a challenging task because many words have multiple senses, and the correct sense can often only be determined by considering the context in which the word is used. For example, in the sentence "I saw a bank of fish swimming in the ocean," the sense of "bank" is clearly the side of a river, not a financial institution or a steep slope.

WSD is an important task in natural language processing (NLP) because it can be used to improve the performance of a variety of NLP tasks, such as machine translation, text summarization, and question answering. For example, a machine translation system

that can correctly disambiguate the word "bank" will be able to produce more accurate translations.

**Methods of WSD:**

- **Dictionary-based and knowledge-based methods:** These methods rely on dictionaries, thesauri, and other lexical knowledge bases to identify the correct sense of a word. For example, the Lesk algorithm is a dictionary-based method that uses the definitions of words and their synonyms to disambiguate word senses.

- **Supervised methods:** These methods use machine learning to train a model on a dataset of text that has been manually annotated with the correct senses of words. The trained model can then be used to disambiguate the senses of words in new text.

- **Semi-supervised methods:** These methods use a combination of labeled and unlabeled data to train a WSD model. Semi-supervised methods are useful for WSD because it can be difficult and time-consuming to manually annotate large amounts of text with the correct senses of words.

- **Unsupervised methods:** These methods do not use any labeled data to train a WSD model. Instead, they rely on statistical methods to identify patterns in the data that can be used to disambiguate word senses.

**Applications of WSD -**

**Machine translation (MT)** is the process of translating text from one language to another using a computer program. WSD is important for MT because it can help to ensure that the correct translation is selected for each word or phrase. For example, if the word "bank" is used in a sentence, the WSD system can help the MT system to determine whether the word should be translated as "financial institution" or "riverbank."

**Information retrieval (IR)** is the process of finding relevant information in a collection of documents. WSD is important for IR because it can help to improve the accuracy of search results. For example, if a user searches for the word "bank," the WSD system can help the IR system to return results that are relevant to both financial institutions and rivers.

**Text mining and information extraction (IE)** are the processes of extracting knowledge from text. WSD is important for text mining and IE because it can help to identify the correct meanings of words and phrases. For example, if a text mining

system is trying to identify the different types of products that a company sells, the WSD system can help to identify the different senses of the word "product" (e.g., physical product, software product, service).

**Lexicography** is the study of words and their meanings. WSD is important for lexicography because it can help lexicographers to write accurate definitions for words. For example, when writing a definition for the word "bank," the lexicographer can use the WSD system to identify the most common senses of the word and to write definitions for those senses.

**Difficulties in WSD**

- **Context dependency:** The meaning of a word can depend on the context in which it is used. For example, the word "sharp" can mean "having a cutting edge" or "keen or intelligent." In the sentence "The knife was sharp," the first sense of "sharp" is intended. However, in the sentence "The student was a sharp thinker," the second sense of "sharp" is intended.

- **Lack of labeled data:** Training data for WSD systems is often difficult and time-consuming to create. This is because each word must be manually annotated with the correct sense in each context in which it is used.

- **Data sparsity:** Some words are used very infrequently in certain contexts. This can make it difficult for WSD systems to accurately disambiguate the senses of these words.

**Dictionary-based methods**

Dictionary-based methods of word sense disambiguation (WSD) are the oldest and simplest type of WSD method. They rely on dictionaries to provide information about the different senses of words.

The most common dictionary-based WSD algorithm is the Lesk algorithm. The Lesk algorithm works by comparing the definitions of words in a context window to find the pair of definitions with the greatest overlap. The sense of the target word is then chosen to be the sense of the context word that overlaps the most with the target word's definition.

For example, consider the following sentence:

> I went to the bank to deposit a check.

The Lesk algorithm would first look up the definitions of the words "bank" and "deposit" in a dictionary. It would then compare the two definitions and find that they overlap significantly. For example, both definitions mention that a bank is a financial institution and that a deposit is a sum of money that is placed in a bank account. The Lesk algorithm would then conclude that the sense of the word "bank" in the sentence is "financial institution."

Dictionary-based WSD methods are relatively simple to implement and they do not require any training data. However, they can be inaccurate if the dictionary is incomplete or inaccurate. Additionally, dictionary-based WSD methods are often unable to disambiguate words that have multiple senses that are very similar in meaning.

# Discourse

## Natural Discourse Processing

Natural discourse processing (NDP) is a subfield of natural language processing (NLP) that deals with the understanding and generation of discourse. Discourse is a sequence of connected sentences that is used to convey a meaning. NDP systems aim to analyze the structure and meaning of discourse in order to perform tasks such as:

- **Coreference resolution:** identifying which expressions in a discourse refer to the same entity. For example, in the sentence "John went to the bank. He deposited a check," the expressions "John" and "he" refer to the same entity.

- **Discourse segmentation:** dividing a discourse into smaller units, such as paragraphs or sentences.

- **Discourse relation classification:** Identifying the relationships between different units of a discourse, such as causal, temporal, and adversative relationships.

- **Discourse summarization:** generating a concise summary of a discourse.

- **Discourse generation:** Generating new discourse, such as stories or news articles.

NDP is a challenging task because discourse is often complex and ambiguous. For example, a discourse may contain multiple coreferences, and the relationships between different units of a discourse may be implicit. Additionally, discourse is often influenced by factors such as the speaker's or writer's background and the context in which the discourse is produced.

**Challenges**

**Anaphoric Ambiguity:** Anaphoric ambiguity is a type of ambiguity that occurs when a pronoun or other anaphoric expression can refer to multiple different entities in a text. This can make it difficult for readers to understand the meaning of the text.

For example, consider the following sentence:

> John went to the bank. He deposited a check.

The pronoun "he" in the second sentence could refer to either John or to someone else. This is because there are two different entities that have been mentioned in the text so far: John and the bank teller.

Anaphoric ambiguity can be caused by a number of different factors, including:

- **Distance:** The more distance there is between the anaphoric expression and its antecedent, the more likely it is that the expression will be ambiguous.

- **Salience:** The more salient the antecedent is, the more likely it is that the anaphoric expression will refer to it.

- **Syntactic complexity:** The more syntactically complex the sentence is, the more likely it is that the anaphoric expression will be ambiguous.

**Word sense discreteness:** This refers to the fact that different senses of a word are distinct entities, with their own unique meanings. For example, the word "bank" has two distinct senses: a financial institution and the side of a river. These two senses are discrete in the sense that they are not related to each other in any way.

**Coherence:** This refers to the fact that a discourse should be semantically coherent, i.e., the different units of the discourse should be related to each other in a meaningful way. For example, the following discourse is coherent:

> John went to the bank. He deposited a check. Then, he went to the store to buy some groceries.

This discourse is coherent because the different units of the discourse are all related to each other. John's trip to the bank and his purchase of groceries are both part of a larger story about his day.

**Discrete segmentation:** This refers to the fact that a discourse can be divided into smaller units, such as paragraphs or sentences, that are semantically complete. For example, the following discourse can be divided into two sentences:

> John went to the bank. He deposited a check.

Each of these sentences is semantically complete, i.e., it conveys a complete thought.

**Reference Resolution**

Reference resolution is the task of identifying which entities are referred to by which expressions in a text or discourse. This is a challenging task because it requires the system to understand the meaning of the text, the context in which it is used, and the relationships between the different entities in the text.

For example, in the sentence "John went to the bank. He deposited a check," the pronoun "he" refers to the entity "John". This can be inferred from the context of the sentence, which indicates that John is the only entity that has been mentioned so far.

Reference resolution is important for many natural language processing (NLP) tasks, such as machine translation, question answering, and text summarization. For example, a machine translation system needs to be able to correctly resolve the referents of pronouns in order to generate a correct translation.

There are two main types of reference resolution:

- **Coreference resolution** is the task of identifying which expressions in a text refer to the same entity. For example, in the sentence "John went to the bank. He deposited a check," the expressions "John" and "he" refer to the same entity.

- **Anaphora resolution** is the task of identifying which expressions in a text refer to previously mentioned entities. For example, in the sentence "John went to the bank. He deposited a check," the pronoun "he" refers to the previously mentioned entity "John".

Reference resolution is a complex task, and there is no single algorithm that can solve it perfectly. However, there are a number of different approaches that can be used, and the best approach will vary depending on the specific NLP task being performed.

**More on coreference resolution**

Coreference resolution is the task of identifying which expressions in a text refer to the same entity. This is a challenging task because it requires the system to understand the meaning of the text, the context in which it is used, and the relationships between the different entities in the text.

For example, in the sentence "John went to the bank. He deposited a check," the expressions "John" and "he" refer to the same entity. This can be inferred from the context of the sentence, which indicates that John is the only entity that has been mentioned so far.

Coreference resolution is important for many natural language processing (NLP) tasks, such as machine translation, question answering, and text summarization. For example, a machine translation system needs to be able to correctly resolve the referents of pronouns in order to generate a correct translation.

**Text Coherence**

Text coherence in natural language processing (NLP) refers to the ability of a system to understand the relationships between the different sentences and paragraphs in a text. It is a challenging task because it requires the system to understand the meaning of the text, the context in which it is used, and the relationships between the different entities in the text.

There are a number of different approaches to text coherence in NLP. One common approach is to use a rule-based system. Rule-based systems use a set of hand-crafted rules to identify the relationships between different sentences and paragraphs. For

example, a rule-based system might have a rule that says that two sentences are coherent if they share the same topic.

Another common approach to text coherence in NLP is to use a machine learning-based system. Machine learning-based systems are trained on a dataset of text that has been manually annotated with coherence labels. The system learns to identify the relationships between different sentences and paragraphs by looking for patterns in the data.

Text coherence is an important task for many NLP applications, such as machine translation, question answering, and text summarization.

# Lexical Translation

Lexical translation is the task of translating individual words or phrases, either on their own (e.g., search-engine queries or meta-data tags) or as part of a knowledge-based Machine Translation (MT) system. In contrast with statistical MT, lexical translation does not require aligned corpora as input. Because large aligned corpora are non-existent for many language pairs, and are very expensive to generate, lexical translation is possible for a much broader set of languages than statistical MT.

While it does not solve the full machine-translation problem, lexical translation is valuable for a number of practical tasks including the translation of search queries, meta-tags, and individual words or phrases. For example, Google and other companies have fielded WordTranslator tools that allow the reader of a Web page to view the translation of particular word, which is helpful if you are, say, a Japanese speaker reading an English text and you come across an unfamiliar word.

Lexical translation is a challenging task because it requires the system to have a deep understanding of the meaning of the words and phrases being translated, as well as the context in which they are being used. For example, the English word "bank" can be translated to the Spanish word "banco" in the context of a financial institution, but it can also be translated to the Spanish word "orilla" in the context of a river.

There are a number of different approaches to lexical translation. One common approach is to use a dictionary-based system. Dictionary-based systems use a bilingual dictionary to translate words and phrases from one language to another. However,

dictionary-based systems can be inaccurate, especially for words and phrases that have multiple meanings.

Another common approach to lexical translation is to use a rule-based system. Rule-based systems use a set of hand-crafted rules to translate words and phrases from one language to another. Rule-based systems can be more accurate than dictionary-based systems, but they are also more complex to develop and maintain.

Machine learning-based systems are also being used for lexical translation. Machine learning-based systems are trained on a dataset of text that has been manually annotated with translations. The system learns to translate words and phrases by identifying patterns in the data. Machine learning-based systems have been shown to achieve good results on lexical translation tasks, but they require a large amount of training data.

Lexical translation is an important task for many NLP applications, such as machine translation, question answering, and text summarization. As NLP research continues to advance, lexical translation systems are likely to become more accurate and efficient.

## Machine Translation

Machine translation (MT) is the process of using a computer to translate text from one language to another. MT systems are trained on large amounts of parallel text, which is text that has been translated into two or more languages. The system learns to identify patterns in the data and to use those patterns to translate new text samples.

There are two main types of MT systems: **statistical MT (SMT) and neural MT (NMT)**. SMT systems use statistical methods to translate text, while NMT systems use neural networks. NMT systems have been shown to achieve better results than SMT systems on most MT tasks.

MT systems are still under development, but they have made significant progress in recent years. MT systems can now translate text with high accuracy for many languages and tasks. As MT research continues to advance, MT systems are likely to become even more accurate and efficient.

Here are some examples of how machine translation is used in the real world:

- Google Translate is a popular MT system that can be used to translate text into over 100 languages. It is used by millions of people around the world to translate websites, documents, and other types of text.

- Netflix uses MT to subtitle and dub its content into multiple languages. This allows people from all over the world to enjoy Netflix content in their own language.

- The United Nations uses MT to translate documents and speeches into multiple languages. This helps to ensure that everyone can understand and participate in the UN's work.

**Statistical MT (SMT)**

Statistical machine translation (SMT) is a type of machine translation (MT) that uses statistical models to translate text from one language to another. SMT systems are trained on large amounts of parallel data, which is text that has been translated into two or more languages. The system learns to identify patterns in the data and to use those patterns to translate new text samples.

SMT systems typically use two main models: a translation model and a language model. The translation model predicts the most likely translation of a word or phrase in the source language, given the context of the sentence. The language model predicts the most likely sequence of words in the target language, given the translation of the source language sentence.

SMT systems have been shown to achieve good results on a variety of MT tasks, including translation of news articles, scientific papers, and government documents. However, SMT systems can be sensitive to the quality of the training data, and they may not perform well on tasks that require a deep understanding of the source or target language, such as translation of literary texts or poetry.

**Hybrid MT**

Hybrid machine translation (hybrid MT) is a type of machine translation that combines statistical machine translation (SMT) and neural machine translation (NMT) to improve the quality of the translation.

**Statistical machine translation:** SMT systems are trained on large amounts of parallel data, which is text that has been translated into two or more languages. The system learns to identify patterns in the data and to use those patterns to translate new text samples.

**Neural machine translation:** NMT systems use neural networks to translate text. Neural networks are inspired by the human brain, and they can learn to perform complex tasks without being explicitly programmed.

Hybrid MT systems combine the best features of SMT and NMT systems. SMT systems are good at translating general-purpose text, while NMT systems are good at translating text that is specific to a particular domain, such as medical or legal text.

Hybrid MT systems typically work by first using an SMT system to generate a preliminary translation of the source text. The NMT system is then used to improve the quality of the translation by correcting errors and making the translation more fluent.

Hybrid MT systems have been shown to achieve better results than SMT or NMT systems alone on a variety of MT tasks. However, hybrid MT systems are more complex to train and deploy than SMT or NMT systems.

Here are some examples of how hybrid MT is used in the real world:

- Google Translate uses a hybrid MT system to translate text into over 100 languages.

- Microsoft Translator uses a hybrid MT system to translate text into over 70 languages.

- Amazon Translate uses a hybrid MT system to translate text into over 24 languages.


**Types of Machine Translation**

Machine translation (MT) is the process of using a computer to translate text from one language to another. There are four main types of MT systems:

**1. Word-based MT**

Word-based MT systems translate text on a word-by-word basis. They use a dictionary to translate each word in the source text into the corresponding word in the target text. Word-based MT systems are simple to implement, but they often produce inaccurate translations, especially for complex texts.

**Example:**

Source text: **I love cats.**

Target text: **J'aime les chats.**

The word-based MT system simply translates each word in the source text into the corresponding word in the target text.

**2. Phrase-based MT**

Phrase-based MT systems translate text on a phrase-by-phrase basis. They use a database of bilingual phrases to translate phrases in the source text into the corresponding phrases in the target text. Phrase-based MT systems produce more accurate translations than word-based MT systems, especially for complex texts.

**Example:**

Source text: **I love cats.**

Target text: **J'aime les chats.**

The phrase-based MT system translates the phrase "I love cats" into the corresponding phrase in French, "J'aime les chats".

**3. Syntax-based MT**

Syntax-based MT systems take into account the syntax of the source text when translating. They parse the source text into a syntactic tree and then generate a target text with the same syntactic structure. Syntax-based MT systems produce more accurate translations than word-based and phrase-based MT systems, but they are more complex to implement.

**Example:**

Source text: **The cat sat on the mat.**

Target text: **Le chat était assis sur le tapis.**

The syntax-based MT system parses the source text into a syntactic tree and then generates a target text with the same syntactic structure.

**4. Rule-based MT**

Rule-based MT systems use a set of hand-crafted rules to translate text. The rules are based on the linguistic knowledge of the source and target languages. Rule-based MT systems can produce very accurate translations, but they are time-consuming and expensive to develop.

**Example:**

Source text: **I love cats.**

Target text: **J'aime les chats.**

The rule-based MT system uses a set of rules to translate the source text into French. For example, the system might have a rule that says that the English word "I" is translated to the French word "je" when it is the subject of a sentence.

**EM Algorithm**

The Expectation-Maximization (EM) algorithm is an iterative algorithm for finding the maximum likelihood estimates of the parameters of a statistical model. It is commonly used in machine learning and statistics for parameter estimation in latent variable models.

The EM algorithm works by iterating between two steps:

1. **Expectation (E-step):** Estimate the latent variables in the model, given the observed data and the current estimates of the model parameters.

2. **Maximization (M-step):** Maximize the expected log-likelihood of the observed data, given the current estimates of the latent variables.

The EM algorithm is guaranteed to converge to a local maximum of the log-likelihood, but it is important to note that it may not find the global maximum.

The EM algorithm is used in a variety of machine learning applications, including:

- **Clustering:** The EM algorithm can be used to cluster data into multiple groups, where the group memberships are latent variables.

- **Mixture models:** The EM algorithm can be used to fit mixture models to data, where the mixture components are latent variables.

- **Hidden Markov models:** The EM algorithm can be used to train hidden Markov models (HMMs), where the hidden states are latent variables.

- **Natural language processing:** The EM algorithm can be used for a variety of natural language processing tasks, such as part-of-speech tagging and word sense disambiguation.

Here is a simple example of how the EM algorithm can be used for clustering:

Suppose we have a dataset of points in 2D, and we want to cluster the points into two groups. We can use the following EM algorithm:

1. **E-step:** For each point, estimate the probability that the point belongs to each group.

2. **M-step:** Update the mean and covariance matrix of each group to maximize the expected log-likelihood of the data.

We can repeat steps 1 and 2 until the algorithm converges, or until we have reached a maximum number of iterations.

Once the algorithm has converged, we can assign each point to the group with the highest probability.