

Exploring Natural Language Processing

Ranjit Kolkar and Grok

March 7, 2025

Chapter 1

Introduction to Natural Language Processing

1.1 Origins and Challenges of NLP

Introduction

Natural Language Processing (NLP) is an interdisciplinary field that combines computer science, linguistics, and artificial intelligence to enable computers to understand, interpret, and generate human languages. The origins of NLP can be traced back to the 1950s when researchers first began exploring the possibility of machine translation.

Historical Context:

- **1950s - 1960s:** Early machine translation projects like the Georgetown-IBM experiment in 1954.
- **1970s:** Development of rule-based systems for language parsing.
- **1980s - 1990s:** Shift towards statistical methods with the advent of more powerful computers.
- **2000s onwards:** Rise of machine learning and deep learning techniques, significantly improving NLP capabilities.

Challenges in NLP:

- **Ambiguity:** Words can have multiple meanings or functions (e.g., "bank" can refer to a financial institution or the side of a river).
- **Syntax and Structure:** The complexity of human languages, including different sentence structures and rules.
- **Semantics:** Understanding the meaning behind words and phrases, where context plays a crucial role.

- **Pragmatics:** Interpreting language based on context, including cultural nuances and implied meanings.
- **Variability:** Languages evolve, slang emerges, and dialects vary, making uniform processing difficult.
- **Resource Scarcity:** Many languages lack the digital resources needed for effective NLP systems.

1.2 Language Modeling

Language modeling is a fundamental concept in Natural Language Processing (NLP), aiming to predict the probability of a sequence of words. It serves as the backbone for numerous NLP applications, including speech recognition, machine translation, text generation, and more.

Language models (LMs) are designed to understand and generate human-like text by estimating the likelihood of a word based on its preceding context. For example, given a sequence of words w_1, w_2, \dots, w_{n-1} , a language model predicts the probability $P(w_n | w_1, w_2, \dots, w_{n-1})$.

Key goals of language modeling include:

- Capturing syntactic and semantic relationships between words.
- Handling the ambiguity and variability of natural language.
- Providing robust mechanisms for text prediction and generation.

Modern language models are categorized into different approaches, including grammar-based, statistical, and neural network-based methods, each with its strengths and limitations.

Applications of Language Modeling

- **Predictive Text Input:** Autocompletion systems in search engines and mobile devices.
- **Speech Recognition:** Decoding spoken words into text.
- **Machine Translation:** Generating accurate translations by modeling cross-linguistic probabilities.

1.3 Grammar-based Language Modeling

Grammar-based language modeling relies on formal linguistic rules to construct language models. It uses predefined grammar structures to define the syntactic and sometimes semantic possibilities of language sequences. This approach is highly structured, deterministic, and interpretable.

Core Concepts

- **Context-Free Grammar (CFG):** A type of formal grammar where each production rule maps a single non-terminal symbol to a sequence of non-terminals and terminals. CFG is widely used in parsing and syntax analysis.
- **Probabilistic Context-Free Grammar (PCFG):** An extension of CFG that assigns probabilities to production rules, enabling the ranking of different parses for ambiguous inputs.

Advantages

- Explicit encoding of linguistic rules.
- Transparency and interpretability in model predictions.
- Effective for domain-specific or constrained language tasks.

Limitations

- Inability to handle variability and ambiguity inherent in natural language effectively.
- Lack of adaptability to unseen data without significant manual effort.
- Scalability issues with complex or large-scale language tasks.

Grammar-based approaches laid the foundation for language modeling but are often complemented or replaced by statistical and neural methods in modern NLP applications.

Grammar-based language models are foundational in NLP for understanding and generating language through structured rules. These models focus on:

Context-Free Grammars (CFG): - CFGs define a set of rules describing how to generate valid sentences in a language. Each rule specifies how non-terminal symbols can be replaced by a sequence of terminals or other non-terminals.

- **Example Rule:** $S \rightarrow NPVP$ where S is the start symbol, NP is a noun phrase, and VP is a verb phrase.

Use in Parsing:

- Parsing with CFGs involves determining if a string of words can be derived from the grammar's start symbol, which is useful for syntax analysis.

- **Parse Trees:** These models can generate parse trees, visually representing the grammatical structure of sentences.

Limitations:

- Grammar-based models struggle with natural language's ambiguity and variability since real language often deviates from strict grammatical rules.

- They are less flexible in adapting to new or informal language constructs unless the grammar is meticulously maintained or expanded.

Applications:

- Used in traditional NLP systems for tasks like sentence structure analysis or in education for teaching language syntax.

1.4 Statistical Language Models

Statistical language models (SLMs) predict word sequences based on statistical data, rather than strict rules:

N-gram Models: - An n-gram model considers a sequence of n items from a text or speech corpus, where n is the number of items considered together. Common types include: - **Unigrams** (context-free), - **Bigrams** (context of one previous word), - **Trigrams** (context of two previous words).

- **Example:** A bigram model might estimate the probability of "the cat" following "sat" based on how often "sat the cat" appears in a corpus.

Markov Assumption: - SLMs often use the Markov assumption, which states that the probability of a word depends only on the previous $n - 1$ words, simplifying computation.

Smoothing Techniques: - Since some word sequences might never occur in training data, smoothing techniques like Add-one (Laplace) smoothing or Kneser-Ney smoothing are used to assign non-zero probabilities to unseen n-grams.

Code Example for Bigram Model:

```
from collections import defaultdict

def bigram_model(text):
    words = text.split()
    bigrams = defaultdict(int)
    unigrams = defaultdict(int)

    for i in range(len(words) - 1):
        bigrams[(words[i], words[i+1])] += 1
        unigrams[words[i]] += 1

    bigram_probs = {}
    for (w1, w2), count in bigrams.items():
        # Probability estimation using Maximum Likelihood Estimation
        bigram_probs[(w1, w2)] = count / unigrams[w1]

    return bigram_probs

# Example usage:
text = "the cat sat on the mat the dog sat"
model = bigram_model(text)
print(model)
```

Applications: - Speech recognition, where predicting the next word aids in transcription accuracy. - Machine translation, for choosing the most probable word sequence in the target language.

Statistical models excel in handling the variability and ambiguity of natural language by learning from data, but they require substantial text corpora for training and can sometimes overfit to the training data if not generalized properly.

1.5 Regular Expressions

Regular expressions, often abbreviated as regex or regexp, are a powerful tool for text processing in NLP, offering a method to match, search, and manipulate text based on specific patterns:

Pattern Matching: - Regular expressions allow you to define patterns for matching strings within text. These patterns can be simple or complex, capturing various aspects like character types, repetitions, and positions.

Syntax and Common Patterns: - **Literals:** Match exact characters or strings. - Example: 'cat' matches the word "cat". - **Metacharacters:** Characters like '.', '*', '+', '?', '^', '\$', '[', ']', '(', ')', '{', '}', '|', and '\'. - '.' matches any single character except newline. - '*' matches zero or more occurrences of the previous character or group.

Table 1.1: Common Literals in Regular Expressions

Literal	Meaning
a	Matches the character 'a'
abc	Matches the string "abc"
123	Matches the number 123
.	Matches any single character except newline (dot)
[abc]	Matches any of the characters 'a', 'b', or 'c'
[^abc]	Matches any character except 'a', 'b', or 'c'
\d	Matches any digit (equivalent to [0-9])
\D	Matches any non-digit character (equivalent to [^0-9])
\w	Matches any word character (equivalent to [a-zA-Z0-9_])
\W	Matches any non-word character (equivalent to [^a-zA-Z0-9_])
\s	Matches any whitespace character (spaces, tabs, line breaks)
\S	Matches any non-whitespace character

- **Character Classes:** Define sets of characters: - '[abc]' matches any of the characters a, b, or c. - '\w' matches any word character (equivalent to '[a-zA-Z0-9_]').

- **Anchors:** Specify positions within the string: - '^' matches the start of the string. - '\$' matches the end of the string.

- **Groups and Alternations:** - '(cat|dog)' matches either "cat" or "dog". - '(3)' matches exactly three digits and captures them as a group.

Table 1.2: Metacharacters in Regular Expressions

Metacharacter	Meaning
.	Matches any single character except newline
*	Matches zero or more occurrences of the previous character or group
+	Matches one or more occurrences of the previous character or group
?	Matches zero or one occurrence of the previous character or group (optional)
[]	Defines a character set; matches any character within the brackets
[^]	Defines a negated character set; matches any character not in the brackets
\	Escapes a metacharacter; treats the following character as a literal
	Alternation; matches either the expression before or after the pipe
()	Groups expressions; used for precedence and capturing
{ }	Specifies the exact number of occurrences of the previous character or group
	Matches the start of a string or line
\$	Matches the end of a string or line

Use in NLP: - Tokenization: Regular expressions can be used to split text into tokens based on patterns for words, numbers, or punctuation. - **Data Cleaning:** Remove or replace unwanted characters or patterns. - **Information Extraction:** Extract specific information like email addresses or dates.

Python Example for Regex in NLP:

```
import re

# Example of tokenizing text
text = "The quick brown fox jumps over the lazy dog."
tokens = re.findall(r'\b\w+\b', text)
print("Tokens:", tokens)

# Example of finding all email addresses in a text
email_text = "Contact us at support@example.com or info@company.org"
emails =
    re.findall(r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b',
    email_text)
print("Emails found:", emails)

# Example of replacing text
cleaned_text = re.sub(r'[\W\s]', '', text) # remove all punctuation
```



```
print("Cleaned text:", cleaned_text)
```

Advantages: - Highly flexible and concise for complex pattern matching tasks. - Can be used across many programming languages with slight variations in syntax.

Challenges: - Writing complex regex can be error-prone and hard to maintain. - Performance can degrade with very complex or large-scale patterns.

Regular expressions are indispensable in NLP for their ability to process and manipulate text at a granular level. They are particularly useful when dealing with structured data within unstructured text or when needing to perform operations that are too intricate for simpler string methods.

1.6 Finite-State Automata

Finite-State Automata (FSA), also known as Finite Automata (FA), are abstract machines used in NLP for pattern recognition, text parsing, and language processing:

Types of Finite-State Automata:

- **Deterministic Finite Automata (DFA):**
Each state has exactly one transition for each input symbol.
There's no ambiguity in transitions, making DFAs efficient for pattern matching.
- **Non-Deterministic Finite Automata (NFA):**
Allows multiple transitions for a single input or even no transitions (epsilon transitions).
More expressive but less efficient than DFAs for matching as they require more complex processing.

Components of FSA:

- **States:** Represent conditions or stages in processing.
- **Transitions:** Define how to move between states based on input symbols.
- **Start State:** The initial state of the automaton.
- **Accept States:** States where the automaton successfully recognizes a pattern or word.

Application in NLP:

- **Morphological Analysis:** FSAs can be used to model word structures, recognizing prefixes, suffixes, and stems.
- **Tokenization:** Breaking down text into tokens can be modeled with FSAs to handle different rules for word boundaries.

- **Pattern Matching:** Regular expressions are essentially descriptions of FSAs for pattern matching in strings.

Example of FSA for Word Recognition:

Let's consider an FSA that recognizes the words "cat" or "dog":

- **States:** {q0, q1, q2, q3, q4} where q0 is the start state, q2 and q4 are accept states for "cat" and "dog" respectively.
- **Transitions:**
 - q0 → q1 on 'c'
 - q1 → q2 on 'a'
 - q2 → q3 on 't'
 - q0 → q3 on 'd'
 - q3 → q4 on 'o'
 - q4 → q4 on 'g'

This can be represented visually or in a pseudo-code-like format:

```
q0 -c-> q1 -a-> q2 -t-> q3 (accept state for "cat")
q0 -d-> q3 -o-> q4 -g-> q4 (accept state for "dog")
```

Advantages of Using FSAs in NLP:

- Simplicity and efficiency for certain language tasks.
- Ability to model regular languages, which includes many syntactic structures in natural language.

Limitations:

- FSAs cannot handle context-free or context-sensitive languages directly, limiting their capability in dealing with more complex linguistic structures.

Finite-State Automata provide a foundational approach in computational linguistics for tasks that involve recognizing or generating patterns in text, making them an essential concept in NLP.

1.7 English Morphology

Morphology in linguistics deals with the internal structure of words, including how words are formed from morphemes, which are the smallest meaningful units of language. In English, morphology involves:

Morphemes: - **Free Morphemes:** Can stand alone as words (e.g., "book", "run"). - **Bound Morphemes:** Must be attached to other morphemes (e.g., "-s", "-ed").

Inflection vs. Derivation: - Inflection alters a word for grammatical function without changing its core meaning or category (e.g., "cats" from "cat", changing number). - Derivation creates new words by adding morphemes, often changing the word's meaning or part of speech (e.g., "happiness" from "happy", noun from adjective).

Morphological Processes: - Affixation: Adding prefixes or suffixes (e.g., "un-" in "undo", "-ly" in "quickly"). - Compounding: Combining two or more words (e.g., "toothbrush"). - Conversion: Changing word class without affixation (e.g., "run" from verb to noun).

Morphological Analysis in NLP: - Stemming: Reducing words to their root form using heuristic rules (e.g., "running" to "run"). - Lemmatization: More sophisticated, considers context to reduce words to their dictionary form (e.g., "was" to "be").

Example of Stemming in Python with NLTK:

```
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize

stemmer = PorterStemmer()
text = "The quick brown foxes jump over the lazy dogs"
tokens = word_tokenize(text)
stemmed = [stemmer.stem(word) for word in tokens]
print("Original:", tokens)
print("Stemmed:", stemmed)
```

Understanding English morphology aids in several NLP tasks like machine translation, information retrieval, and text normalization.

1.8 Transducers for Lexicon and Rules

Transducers in NLP are computational models that can transform one sequence of symbols into another, useful for both analysis and generation of language:

Types of Transducers: - Finite State Transducers (FST): Like FSAs but map input symbols to output symbols. They are crucial for: - Morphological Analysis: Turning surface forms into their roots and affixes or vice versa. - Phonological Rules: Mapping between phonetic and orthographic representations.

Components: - States: Similar to FSAs. - Transitions: Now labeled with both input and output symbols (e.g., $a : b$ means "read 'a' and write 'b'"). - Initial and Final States: For starting and ending the transduction process.

Application in NLP: - Lexical Transduction: Converting between different forms of words (e.g., singular to plural). - Rule Application: Applying phonological or morphological rules (e.g., turning "run" into "running").

Example of a Simple FST: Consider a transducer for singular to plural noun conversion:

- States: $\{q_0, q_1, q_2\}$ where q_0 is start, q_2 is final. - Transitions: - $q_0 \rightarrow q_1$ on $\epsilon:\epsilon$ (for words ending in 'y', e.g., "city") - $q_1 \rightarrow q_2$ on $y:\text{ies}$ - $q_0 \rightarrow q_2$ on $\epsilon:s$ (for general case, e.g., "cat")

1.9 Tokenization

Tokenization is the process of breaking down text into tokens, which are the basic units used in further NLP processing. These tokens can be words, punctuation, numbers, or even symbols:

Types of Tokenization: - Word Tokenization: Splitting text into words or phrases. - Sentence Tokenization: Segmenting text into sentences. - Subword Tokenization: Breaking words into smaller units, useful for languages with complex morphology or for handling out-of-vocabulary words.

Challenges: - Ambiguity: Words like "can't" might be split into "can" and "t" or kept as one token. - Punctuation: Deciding whether punctuation should be separated or kept with words. - Multi-word Expressions: Recognizing phrases like "New York" as a single token.

Techniques: - Rule-Based: Using regular expressions or predefined rules for tokenization. - Statistical: Learning from data, like using machine learning models to determine token boundaries.

Python Example for Word Tokenization with NLTK:

```
import nltk
from nltk.tokenize import word_tokenize, sent_tokenize

nltk.download('punkt') # Download necessary data for tokenization

text = "Hello, Mr. Smith! How's your day? The weather in New York is
        amazing today."
words = word_tokenize(text)
sentences = sent_tokenize(text)

print("Words:", words)
print("Sentences:", sentences)
```

Applications in NLP: - Pre-processing for text analysis: Essential for any NLP task where text needs to be analyzed at word or sentence level. - Machine Translation: Tokens serve as the basic units for translation models.

1.10 Detecting and Correcting Spelling Errors

Detecting and correcting spelling errors is crucial in NLP to improve text quality and comprehension:

Methods for Detection: - Dictionary Lookup: Comparing words in the text against a known dictionary to find non-matching words. - Statistical Methods: Probability models can flag words that are uncommon in given contexts.

- Contextual Analysis: Using surrounding words to determine if a word fits contextually.

Correction Techniques: - Minimum Edit Distance: Algorithms like Levenshtein distance to find the closest correct word. - N-gram Models: Suggest corrections based on frequently occurring n-grams in a corpus. - Machine Learning: Training models on examples of errors and corrections to predict fixes.

Minimum Edit Distance (Levenshtein Distance): This measures how dissimilar two strings are by counting the minimum number of single-character edits needed to change one word into another:

- Operations: Insertion, deletion, and substitution.

Python Example for Spelling Correction using Edit Distance:

```
def levenshtein(s1, s2):
    if len(s1) < len(s2):
        return levenshtein(s2, s1)

    if len(s2) == 0:
        return len(s1)

    previous_row = range(len(s2) + 1)
    for i, c1 in enumerate(s1):
        current_row = [i + 1]
        for j, c2 in enumerate(s2):
            insertions = previous_row[j + 1] + 1
            deletions = current_row[j] + 1
            substitutions = previous_row[j] + (c1 != c2)
            current_row.append(min(insertions, deletions, substitutions))
        previous_row = current_row

    return previous_row[-1]

# Example of usage
correct_word = "correct"
misspellings = ["corect", "corret", "corrept"]
for word in misspellings:
    dist = levenshtein(word, correct_word)
    print(f"Distance between {word} and {correct_word}: {dist}")
```

Applications: - Text Editing Tools: Spell checkers in word processors or text editors. - Search Engines: Enhancing search accuracy by suggesting corrections. - Speech Recognition: Correcting misheard or mispronounced words.

These processes involve both algorithmic complexity and linguistic knowledge to achieve high accuracy in spelling correction, often requiring a blend of rule-based and machine-learning approaches.

1.11 Edit Distance in Natural Language Processing

Edit Distance, also known as Levenshtein Distance, is a measure of the minimum number of operations required to convert one string into another. These operations include:

- **Insertion:** Adding a character.
- **Deletion:** Removing a character.
- **Substitution:** Replacing one character with another.

Edit distance is widely used in NLP tasks such as spell correction, DNA sequence alignment, and approximate string matching.

1.11.1 Example: Calculating Edit Distance

Consider the two strings:

- String 1: **kitten**
- String 2: **sitting**

To calculate the edit distance, we use a dynamic programming approach. The table below illustrates the step-by-step computation.

		<i>s</i>	<i>i</i>	<i>t</i>	<i>t</i>	<i>i</i>	<i>n</i>	<i>g</i>
	0	1	2	3	4	5	6	7
<i>k</i>	1	1	2	3	4	5	6	7
<i>i</i>	2	2	1	2	3	4	5	6
<i>t</i>	3	3	2	1	2	3	4	5
<i>t</i>	4	4	3	2	1	2	3	4
<i>e</i>	5	5	4	3	2	2	3	4
<i>n</i>	6	6	5	4	3	3	2	3

Explanation:

- Each cell (i, j) represents the minimum number of operations to transform the first i characters of **kitten** into the first j characters of **sitting**.
- Start from the top-left corner and proceed to the bottom-right corner.

Result: The edit distance between **kitten** and **sitting** is **3**, indicating that three operations (substitute 'k' with 's', substitute 'e' with 'i', and insert 'g') are required.

Algorithm 1 Edit Distance Calculation

```

1: Input: Two strings  $A$  of length  $m$  and  $B$  of length  $n$ 
2: Output: Minimum edit distance between  $A$  and  $B$ 
3: Create a matrix  $dp[m+1][n+1]$ 
4: for  $i = 0$  to  $m$  do
5:      $dp[i][0] = i$  ▷ Cost of deleting all characters from  $A$ 
6: end for
7: for  $j = 0$  to  $n$  do
8:      $dp[0][j] = j$  ▷ Cost of inserting all characters of  $B$ 
9: end for
10: for  $i = 1$  to  $m$  do
11:     for  $j = 1$  to  $n$  do
12:         if  $A[i-1] == B[j-1]$  then
13:              $dp[i][j] = dp[i-1][j-1]$  ▷ No cost if characters are the same
14:         else
15:              $dp[i][j] = 1 + \min(dp[i-1][j], dp[i][j-1], dp[i-1][j-1])$  ▷
            Choose the best among deletion, insertion, or substitution
16:         end if
17:     end for
18: end for
19: return  $dp[m][n]$ 

```

1.11.2 Algorithm: Dynamic Programming for Edit Distance

The following algorithm describes the dynamic programming approach to calculate edit distance:

This algorithm has a time complexity of $O(m \times n)$ and a space complexity of $O(m \times n)$.

Chapter 2

UNIT 2

2.1 Introduction

Language modeling plays a crucial role in various Natural Language Processing (NLP) tasks, including speech recognition, machine translation, and text generation. One of the foundational approaches to language modeling is the N-gram model, which approximates the probability distribution of word sequences based on their frequency in a given corpus. However, these models face challenges such as data sparsity and poor generalization, leading to the need for evaluation techniques, smoothing methods, and probabilistic tagging approaches.

In this unit, we delve into the Unsmoothed N-gram models, their limitations, evaluation methods, and the challenges associated with probabilistic language modeling.

2.2 Unsmoothed N-grams

An N-gram model estimates the probability of a word occurring based on the previous $N - 1$ words. Formally, an N-gram model expresses the probability of a word sequence $W = w_1, w_2, \dots, w_n$ as:

$$P(W) = \prod_{i=1}^n P(w_i | w_{i-(N-1)}, \dots, w_{i-1})$$

where N determines the size of the context window.

2.2.1 Types of N-grams

- Unigram Model (N=1): Assumes each word occurs independently:

$$P(W) = \prod_{i=1}^n P(w_i)$$

Example:

- "The cat sat on the mat" is modeled as:

$$P(\text{The})P(\text{cat})P(\text{sat})P(\text{on})P(\text{the})P(\text{mat})$$

- Bigram Model (N=2): Uses adjacent word probabilities:

$$P(W) = \prod_{i=2}^n P(w_i | w_{i-1})$$

Example:

- "The cat sat" is modeled as:

$$P(\text{cat}|\text{The})P(\text{sat}|\text{cat})$$

- Trigram Model (N=3): Uses three-word sequences:

$$P(W) = \prod_{i=3}^n P(w_i | w_{i-2}, w_{i-1})$$

2.2.2 Limitations of Unsmoothed N-grams

Unsmoothed N-gram models suffer from several drawbacks:

- **Data Sparsity:** Many valid word sequences may not appear in the training corpus, leading to zero probabilities.
- **Lack of Generalization:** The model does not handle unseen words effectively and is heavily dependent on training data.
- **Computational Complexity:** Higher-order N-grams require large amounts of memory and training data.
- **Context Length Limitation:** Longer dependencies are ignored in low-order N-grams (e.g., unigrams lose all context information).

2.3 Evaluating N-grams

Evaluating an N-gram model is crucial for understanding its effectiveness in language modeling. The primary metrics used for evaluation include perplexity, out-of-vocabulary (OOV) rate, entropy, and precision-recall measures.

2.3.1 Perplexity (PP)

Perplexity is a widely used measure of language model quality. It quantifies how well a model predicts a sequence of words. Mathematically, it is defined as:

$$PP(W) = P(w_1, w_2, \dots, w_n)^{-\frac{1}{n}}$$

For an N-gram model:

$$PP(W) = \exp \left(-\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_{i-(N-1)}, \dots, w_{i-1}) \right)$$

A lower perplexity indicates a better model as it suggests less uncertainty in predictions.

2.3.2 Out-of-Vocabulary (OOV) Rate

OOV rate measures the percentage of words in a test set that were not observed during training:

$$\text{OOV Rate} = \frac{\text{OOV Words}}{\text{Total Words}}$$

A high OOV rate indicates poor model coverage and the need for better smoothing techniques.

2.3.3 Comparison of Different N-gram Models

Model	Memory Usage	Accuracy	Perplexity
Unigram	Low	Low	High
Bigram	Medium	Medium	Medium
Trigram	High	High	Low

Table 2.1: Comparison of Different N-gram Models

2.3.4 Additional Evaluation Metrics

- Entropy: Measures the unpredictability of a model's predictions, where lower entropy suggests better prediction stability.
- BLEU Score: Used in machine translation to compare model-generated text with human-written reference translations.
- Precision and Recall: Applied in NLP tasks to measure how accurately the predicted sequences match the reference corpus.

2.4 Interpolation and Backoff

Interpolation and **backoff** are two essential techniques used to improve N-gram models by addressing data sparsity and ensuring better probability estimation for unseen word sequences.

2.4.1 Interpolation

Interpolation combines multiple N-gram models (e.g., unigram, bigram, and trigram) by assigning them weights to create a more robust language model. Instead of relying solely on a single N-gram model, interpolation considers all possible models with a weighted sum.

Mathematical Formulation:

$$P(w_i|w_{i-1}) = \lambda_1 P_{unigram}(w_i) + \lambda_2 P_{bigram}(w_i|w_{i-1}) + \lambda_3 P_{trigram}(w_i|w_{i-2}, w_{i-1})$$

where:

- $\lambda_1, \lambda_2, \lambda_3$ are interpolation weights such that $\lambda_1 + \lambda_2 + \lambda_3 = 1$.
- The probabilities from different N-gram models are combined to improve performance.

Example: Given a trigram model trained on a corpus, if the phrase "The cat sat" appears rarely, the model can assign weight to lower-order N-grams, such as bigram and unigram probabilities, to estimate the likelihood of the next word.

2.4.2 Backoff

Backoff is an alternative approach that falls back to lower-order N-gram models when the higher-order model encounters unseen word sequences.

Mathematical Formulation:

$$P(w_i|w_{i-2}, w_{i-1}) = \begin{cases} P_{trigram}(w_i|w_{i-2}, w_{i-1}), & \text{if trigram exists} \\ \alpha P_{bigram}(w_i|w_{i-1}), & \text{otherwise} \\ \beta P_{unigram}(w_i), & \text{otherwise} \end{cases}$$

where α and β are discounting factors that adjust probability estimates.

Example: If the trigram "The cat sings" is unseen in training data, the model backtracks to the bigram "cat sings." If this is also unseen, it falls back to the unigram probability of "sings."

2.5 Word Classes

Word classes, also known as **Part-of-Speech (PoS) categories**, group words based on their syntactic functions, such as nouns, verbs, adjectives, and adverbs. These classes help in linguistic analysis, language modeling, and syntactic parsing.

2.5.1 Major Word Classes

- **Nouns (N):** Represent entities such as people, places, or things (e.g., "dog," "India," "happiness").
- **Verbs (V):** Indicate actions or states (e.g., "run," "is," "write").
- **Adjectives (Adj):** Describe nouns (e.g., "beautiful," "red," "large").
- **Adverbs (Adv):** Modify verbs, adjectives, or other adverbs (e.g., "quickly," "very," "carefully").
- **Prepositions (P):** Show relationships between nouns (e.g., "in," "on," "under").

2.5.2 Closed vs Open Word Classes

Type	Description	Examples
Open Class	Can easily accept new words	Nouns, Verbs, Adjectives
Closed Class	Rarely accepts new words	Pronouns, Prepositions, Conjunctions

Table 2.2: Comparison of Open and Closed Word Classes

2.5.3 Applications of Word Classes

- PoS Tagging: Assigns syntactic labels to words in a sentence.
- Named Entity Recognition (NER): Identifies specific entities such as names and places.
- Grammar Checking: Helps in syntactic parsing for sentence correction.

2.6 Part-of-Speech Tagging

Part-of-Speech (PoS) tagging is the process of assigning parts of speech such as noun, verb, adjective, etc., to each word in a sentence based on its definition and context. It plays a crucial role in Natural Language Processing (NLP) applications like text analysis, machine translation, and speech recognition.

2.6.1 Example

extbfSentence: *The quick brown fox jumps over the lazy dog.* extbfPoS Tags: DT JJ JJ NN VBZ IN DT JJ NN

2.7 Rule-Based Tagging

Rule-based PoS tagging relies on manually crafted linguistic rules to assign tags to words. These rules consider syntactic and morphological features to determine the appropriate tag.

2.7.1 Algorithm

1. Assign the most common PoS tag to each word based on a lexicon.
2. Apply hand-crafted rules to resolve ambiguities using contextual clues.

2.7.2 Example

extbfSentence: *He can fish.* extbfInitial Tags: PRP MD NN (Incorrect: “fish” as a noun) extbfRule Applied: If “can” is used as a modal verb, the next word is likely a verb. extbfFinal Tags: PRP MD VB (Corrected: “fish” as a verb)

2.8 Stochastic Tagging

Stochastic tagging uses probabilistic models to determine the PoS tag based on statistical properties of the text. It is data-driven and relies on training from large corpora.

2.8.1 Types of Stochastic Tagging

- **Unigram Tagger:** Assigns the most probable PoS tag to a word based on its frequency in the corpus.
- **Bigram/Trigram Tagger:** Uses previous words’ tags to predict the current word’s tag.
- **Hidden Markov Model (HMM):** Computes the probability of a tag sequence using transition and emission probabilities.

2.8.2 Example

extbfSentence: *She will run.* extbfTagging Process: Using bigram probabilities, if “will” is identified as a modal verb, “run” is more likely to be tagged as a verb rather than a noun. extbfFinal Tags: PRP MD VB

2.9 Transformation-Based Tagging

Transformation-Based Tagging, also known as **Brill Tagging**, is a hybrid approach that combines rule-based and machine learning methods. Unlike purely statistical methods, Brill tagging learns a set of transformation rules from a tagged corpus, iteratively refining the tagging results.

2.9.1 Algorithm

1. Assign initial PoS tags using a simple method such as a unigram tagger.
2. Apply transformation rules iteratively to correct errors in tagging.
3. The rules are learned from a corpus and are ranked based on effectiveness.

2.9.2 Example

Sentence: *I saw a bear.*

Initial Tags: PRP VBD DT NN (incorrect: “saw” as a verb)

Transformation Rule: If a word follows “a” and is known as a noun in the corpus, tag it as NN instead of VBD.

Final Tags: PRP VBD DT NN (Corrected: “saw” remains a verb but other words are correctly tagged)

2.10 Issues in PoS Tagging

Despite advances in tagging techniques, several challenges persist in Part-of-Speech tagging:

2.10.1 Ambiguity

Many words in English can belong to multiple PoS categories based on context.

- Example: “Book” (NN) as a noun in “This is a book.” vs. “Book” (VB) as a verb in “Please book a ticket.”

2.10.2 Out-of-Vocabulary (OOV) Words

Unknown words that are not present in the training corpus pose a challenge to statistical and rule-based taggers.

- Example: A newly coined word like “cryptozoology” may be difficult to classify correctly.

2.10.3 Domain-Specific Language

Taggers trained on general corpora may not perform well on domain-specific texts such as medical or legal documents.

- Example: In a medical corpus, “lead” might be a verb, while in a chemistry text, it is a noun referring to the metal.

2.10.4 Multiword Expressions

Fixed phrases and idioms often require special handling as their meaning cannot be inferred from individual words.

- Example: “Kick the bucket” should not be tagged as literal words but as an idiomatic expression.

Chapter 3

Syntactic Analysis

Syntactic analysis, or parsing, is the process of analyzing a sentence based on its syntax. It plays a crucial role in computational linguistics, natural language processing, and compiler design. This chapter discusses the foundational aspects of syntactic analysis, including Context-Free Grammars (CFGs), grammar rules for English, treebanks, and normal forms for grammar.

3.1 Context-Free Grammars

A *Context-Free Grammar* (CFG) is a formal system that defines the syntactic structure of a language. It consists of:

- A set of non-terminal symbols (N)
- A set of terminal symbols (Σ)
- A set of production rules (P) of the form $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in (N \cup \Sigma)^*$
- A start symbol ($S \in N$)

3.1.1 Example of a Simple CFG

A simple CFG for a subset of English could be:

$$\begin{aligned} S &\rightarrow NP VP \\ NP &\rightarrow Det N \\ VP &\rightarrow V NP \\ Det &\rightarrow \{\text{the, a}\} \\ N &\rightarrow \{\text{dog, cat}\} \\ V &\rightarrow \{\text{chased, saw}\} \end{aligned}$$

A derivation using this grammar for "the dog chased a cat" follows:

$$\begin{aligned}
 S &\Rightarrow NP VP \\
 &\Rightarrow Det N VP \\
 &\Rightarrow \text{the dog } VP \\
 &\Rightarrow \text{the dog } V NP \\
 &\Rightarrow \text{the dog chased } Det N \\
 &\Rightarrow \text{the dog chased a cat}
 \end{aligned}$$

3.1.2 Applications of CFGs

CFGs are widely used in various domains:

- **Programming Languages:** Used in the design of compilers and interpreters for defining syntax and parsing code.
- **Natural Language Processing (NLP):** Used for syntactic parsing in machine translation, grammar checking, and speech recognition.
- **Artificial Intelligence:** Helps in understanding and processing structured data.
- **Query Languages:** CFGs define structured queries in databases (e.g., SQL parsing).

3.2 Grammar Rules for English

English grammar follows a hierarchical phrase structure, which can be effectively modeled using a CFG. The primary syntactic components include:

3.2.1 Basic Phrase Structure Rules

English sentences generally follow a subject-verb-object (SVO) order. Some key syntactic structures include:

- **Sentence (S):** Consists of a noun phrase (NP) followed by a verb phrase (VP).
- **Noun Phrase (NP):** Typically consists of a determiner (Det) followed by an optional adjective (Adj) and a noun (N).
- **Verb Phrase (VP):** Contains a verb (V) and optionally a noun phrase (NP) or a prepositional phrase (PP).
- **Prepositional Phrase (PP):** Consists of a preposition (P) followed by a noun phrase (NP).

3.2.2 Example Grammar Rules for English

A simplified CFG representing these structures can be defined as:

$$\begin{aligned}
 S &\rightarrow NP VP \\
 NP &\rightarrow Det (Adj) N \\
 VP &\rightarrow V (NP) (PP) \\
 PP &\rightarrow P NP \\
 Det &\rightarrow \{\text{the, a, an}\} \\
 Adj &\rightarrow \{\text{big, small, red, fast}\} \\
 N &\rightarrow \{\text{dog, cat, car, tree}\} \\
 V &\rightarrow \{\text{chased, saw, ate, climbed}\} \\
 P &\rightarrow \{\text{on, in, under, over}\}
 \end{aligned}$$

3.2.3 Example Sentence Parsing

Consider the sentence: *"The big dog chased a cat on the tree."*

Using the grammar rules, its derivation is:

$$\begin{aligned}
 S &\Rightarrow NP VP \\
 &\Rightarrow Det Adj N VP \\
 &\Rightarrow \text{The big dog } VP \\
 &\Rightarrow \text{The big dog } V NP PP \\
 &\Rightarrow \text{The big dog chased } Det N PP \\
 &\Rightarrow \text{The big dog chased a cat } P NP \\
 &\Rightarrow \text{The big dog chased a cat on } Det N \\
 &\Rightarrow \text{The big dog chased a cat on the tree}
 \end{aligned}$$

This breakdown follows the hierarchical phrase structure, making it easier to understand syntactic relationships.

3.2.4 Challenges in English Grammar Parsing

Parsing natural language grammar comes with challenges:

- **Ambiguity:** Sentences may have multiple valid parse trees (e.g., "I saw the man with a telescope.").
- **Long-Distance Dependencies:** Elements far apart in a sentence may influence each other.
- **Irregular Structures:** English has idiomatic expressions and exceptions that do not fit strict CFG rules.

Despite these challenges, CFG-based syntactic parsing remains a foundational approach in computational linguistics and NLP applications.

3.3 Treebanks

A *Treebank* is a structured linguistic resource where sentences are annotated with their syntactic structures in the form of parse trees. Treebanks serve as a crucial resource in natural language processing (NLP) and computational linguistics for training and evaluating syntactic parsers.

3.3.1 Purpose of Treebanks

Treebanks provide manually or semi-automatically annotated corpora that help in:

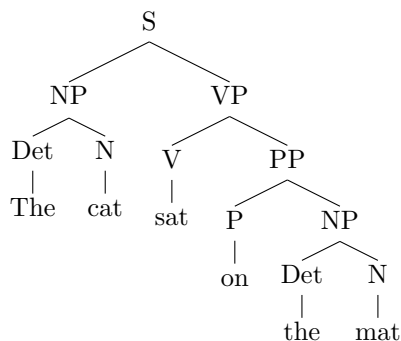
- Developing and training probabilistic parsers.
- Evaluating parsing algorithms against a standard dataset.
- Analyzing syntactic patterns and structures in large text corpora.
- Supporting linguistic research by providing structured representations of syntax.

3.3.2 Example of a Treebank Annotation

Consider the sentence:

The cat sat on the mat.

A treebank might annotate this sentence using a hierarchical tree structure as follows:



This structure represents:

- **S** (Sentence) consisting of a noun phrase (NP) and a verb phrase (VP).
- **NP** (Noun Phrase) consisting of a determiner (Det) and a noun (N).
- **VP** (Verb Phrase) consisting of a verb (V) and a prepositional phrase (PP).
- **PP** (Prepositional Phrase) consisting of a preposition (P) and a noun phrase (NP).

3.3.3 Penn Treebank and Other Resources

One of the most widely used treebanks is the **Penn Treebank**, which provides a large corpus of English text annotated with phrase-structure trees and part-of-speech tags. Other notable treebanks include:

- **Universal Dependencies (UD) Treebanks** – A collection of syntactically annotated treebanks for multiple languages.
- **TIGER Treebank** – A treebank for German, annotated with syntactic and semantic structures.
- **Chinese Treebank** – A structured corpus for Mandarin Chinese syntax.

3.3.4 Applications of Treebanks

Treebanks play a significant role in various NLP applications:

- **Training Statistical Parsers:** Supervised machine learning models use treebanks to learn syntactic parsing.
- **Grammar Induction:** Helps in constructing and refining grammatical rules.
- **Semantic Analysis:** Used in higher-level NLP tasks such as question answering and machine translation.
- **Speech Recognition:** Helps improve syntactic parsing in spoken language systems.

3.3.5 Challenges in Treebank Construction

Building treebanks involves challenges such as:

- **Manual Annotation Effort:** Creating treebanks requires extensive linguistic expertise and time.
- **Inconsistencies in Annotation:** Variations in human annotation lead to inconsistencies in the dataset.
- **Scalability:** Annotating large corpora is difficult and expensive.
- **Cross-Language Adaptability:** Different languages exhibit distinct syntactic structures, making uniform annotation challenging.

Treebanks are indispensable tools in computational linguistics, aiding in syntactic parsing, machine learning, and linguistic analysis. They provide a structured means of representing sentence syntax, allowing researchers to develop and refine NLP models. Despite the challenges, treebanks continue to be a cornerstone in advancing the field of syntactic analysis.

3.4 Normal Forms for Grammar

In the study of formal grammars, particularly in Natural Language Processing (NLP), *normal forms* provide a structured and standardized way of representing grammars. They simplify parsing, improve computational efficiency, and facilitate language processing tasks such as syntax analysis and machine translation.

3.4.1 Chomsky Normal Form (CNF)

A context-free grammar (CFG) is in **Chomsky Normal Form (CNF)** if all production rules conform to one of the following forms:

- $A \rightarrow BC$, where A, B, C are non-terminal symbols (except that B and C cannot be the start symbol).
- $A \rightarrow a$, where a is a terminal symbol.

Additionally, no production should derive the empty string (ϵ) unless explicitly allowed.

Example:

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Advantages of CNF

- Ensures every derivation has a predictable structure, aiding efficient parsing.
- Used in dynamic programming algorithms like the **CYK (Cocke–Younger–Kasami)** algorithm for parsing.
- Simplifies transformations between different grammars.

3.4.2 Greibach Normal Form (GNF)

A context-free grammar is in **Greibach Normal Form (GNF)** if all production rules follow the structure:

$$A \rightarrow a\alpha$$

where:

- A is a non-terminal symbol.
- a is a terminal symbol.
- α is a sequence of non-terminals (which may be empty).

Example:

$$S \rightarrow aA$$

$$A \rightarrow bB$$

$$B \rightarrow c$$

Advantages of GNF

- Facilitates recursive descent parsing, as each rule starts with a terminal.
- Guarantees that every derivation expands in a leftmost manner, making it useful for top-down parsing techniques.
- Helps in the conversion of grammars into pushdown automata.

Normal forms play a crucial role in NLP and computational linguistics by providing structured representations of grammars. CNF is essential for efficient parsing, GNF simplifies recursive descent parsing, and Kuroda Normal Form extends to context-sensitive grammars. Understanding these normal forms is key to designing robust language models and grammar parsers.

3.5 Dependency Grammar

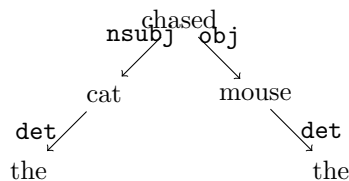
Dependency Grammar (DG) is a linguistic framework that represents syntactic structure based on the dependency relations between words in a sentence. Unlike phrase structure grammar, which uses constituents and hierarchical trees, dependency grammar emphasizes direct relations between words.

3.5.1 Basic Concepts of Dependency Grammar

- **Head-Dependent Relationship:** Each word (except the root) depends on another word, known as its *head*.
- **Dependency Tree:** The sentence structure is represented as a directed tree where edges indicate dependency relations.
- **Labeled Dependencies:** Relationships between words are annotated with labels such as **nsubj** (nominal subject), **obj** (object), **det** (determiner), etc.

3.5.2 Example of Dependency Structure

Consider the sentence: *"The cat chased the mouse."* A dependency tree for this sentence would look like:



3.5.3 Applications of Dependency Grammar

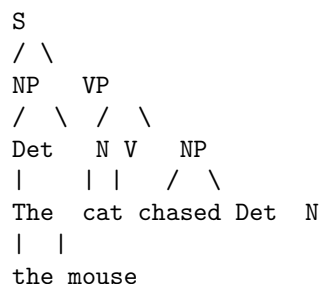
- **Syntactic Parsing:** Used in dependency-based parsers for analyzing sentence structures.
- **Machine Translation:** Helps in understanding sentence meaning and improving translation accuracy.
- **Information Extraction:** Extracts meaningful relations between entities in a sentence.
- **Sentiment Analysis:** Identifies sentiment-carrying words and their dependencies.

3.6 Syntactic Parsing

Syntactic Parsing is the process of analyzing the grammatical structure of a sentence to determine its syntactic composition. It is crucial in natural language understanding (NLU) and machine translation.

3.6.1 Types of Syntactic Parsing

- **Constituency Parsing (Phrase Structure Parsing):**
 - Represents a sentence using hierarchical tree structures.
 - Uses Context-Free Grammars (CFGs) to define valid sentence structures.
 - Example parsing tree for "The cat chased the mouse":



- **Dependency Parsing:**

- Focuses on direct relationships between words.
- Uses Dependency Grammar as its formalism.
- More efficient for tasks like machine translation and sentiment analysis.

3.6.2 Parsing Algorithms

- **Chart Parsing:**
 - Uses dynamic programming (e.g., CYK Algorithm for CFGs).
 - Efficient for structured grammar parsing.
- **Transition-Based Parsing:**
 - Predicts dependencies incrementally using machine learning models.
 - Fast and widely used in modern NLP applications.
- **Graph-Based Parsing:**
 - Treats dependency parsing as a graph optimization problem.
 - Uses algorithms like Maximum Spanning Tree (MST) parsing.

3.6.3 Applications of Syntactic Parsing

- **Machine Translation:** Helps align words grammatically for accurate translation.
- **Question Answering:** Identifies syntactic structures for extracting relevant answers.
- **Speech Recognition:** Improves sentence understanding by ensuring grammatical correctness.
- **Text Summarization:** Enhances semantic comprehension for better summarization results.

Dependency Grammar and Syntactic Parsing play vital roles in NLP by analyzing and representing sentence structures. Constituency parsing helps in phrase-level analysis, while dependency parsing focuses on direct word relations, making both approaches essential for linguistic and AI applications.

3.7 Ambiguity in NLP

Ambiguity is a fundamental challenge in Natural Language Processing (NLP), where a given sentence or phrase can have multiple interpretations. It can occur at different linguistic levels, leading to difficulties in parsing, translation, and semantic analysis.

3.7.1 Types of Ambiguity

- **Lexical Ambiguity:** Occurs when a word has multiple meanings. Example: The word "bank" can mean a financial institution or the side of a river.
- **Syntactic Ambiguity:** Arises when a sentence can have multiple valid syntactic structures. Example: "*I saw the man with a telescope.*" (Did I see a man using a telescope, or did I use a telescope to see the man?)
- **Semantic Ambiguity:** Happens when a sentence has multiple meanings despite having a clear syntax. Example: "*Visiting relatives can be annoying.*" (Are relatives visiting, or is visiting them annoying?)
- **Pragmatic Ambiguity:** Depends on context or background knowledge. Example: "*Can you pass the salt?*" (Asking about ability vs. making a request)
- **Attachment Ambiguity:** Occurs when a modifier (e.g., a prepositional phrase) can be attached to different parts of the sentence. Example: "*She saw the boy with the binoculars.*" (Does "with the binoculars" modify "the boy" or "saw"?)

3.7.2 Handling Ambiguity

- **Probabilistic Parsing:** Assigns probabilities to different parses and selects the most likely one.
- **Semantic Role Labeling (SRL):** Identifies the roles of words in a sentence.
- **Context-based Disambiguation:** Uses contextual clues to resolve meanings.
- **Machine Learning and Deep Learning Models:** Use neural networks to predict the correct meaning based on large datasets.

3.8 Dynamic Programming Parsing

Dynamic Programming (DP) parsing is an efficient approach to syntactic analysis, reducing redundant computations when parsing sentences with recursive structures.

3.8.1 Motivation for DP in Parsing

Parsing algorithms often involve recursive processes that may repeatedly analyze the same substructures. Dynamic programming helps by:

- Avoiding redundant computations.

- Storing intermediate results for future use.
- Improving efficiency compared to naive recursive approaches.

3.8.2 CYK Algorithm (Cocke-Younger-Kasami)

One of the most well-known DP-based parsing algorithms is the CYK algorithm, which applies to Context-Free Grammars (CFGs) in Chomsky Normal Form (CNF).

1. Convert the given CFG into CNF.
2. Construct a DP table where each cell stores possible non-terminal symbols generating the corresponding substring.
3. Fill the table bottom-up, combining smaller substrings to form larger ones.

Example: Parsing the sentence "the cat sleeps" with a simplified CNF grammar.

```

S
/  \
NP   VP
/  \  |
Det  N  V
|    |  |
the  cat sleeps

```

3.8.3 Earley Parsing

Unlike CYK, which requires CNF, the Earley parser works with any CFG. It follows three steps:

- **Prediction:** Expands non-terminals based on grammar rules.
- **Scanning:** Matches input tokens.
- **Completion:** Combines results to form valid parses.

3.8.4 Applications of DP Parsing

- Speech Recognition: Efficiently analyzes sentence structures.
- Machine Translation: Parses complex linguistic patterns.
- Compilers: Syntax analysis in programming languages.

3.9 Shallow Parsing (Chunking)

Shallow Parsing, also known as chunking, is a technique in NLP where text is processed to identify intermediate syntactic structures without full parsing.

3.9.1 Difference Between Shallow and Full Parsing

- **Full Parsing:** Generates complete parse trees with hierarchical structures.
- **Shallow Parsing:** Extracts phrases (e.g., noun phrases, verb phrases) without deeper grammatical analysis.

3.9.2 Techniques for Shallow Parsing

- **Regular Expressions:** Rule-based approaches for extracting syntactic chunks.
- **Machine Learning Models:** Uses classifiers (e.g., CRF, BiLSTM) to recognize chunk boundaries.
- **Part-of-Speech (POS) Tagging:** Helps in identifying noun and verb phrases.

3.9.3 Example of Chunking

Consider the sentence: *"The quick brown fox jumps over the lazy dog."*

Chunking Output:

[Noun Phrase The quick brown fox] [Verb Phrase jumps]
[Prepositional Phrase over] [Noun Phrase the lazy dog]

3.9.4 Applications of Shallow Parsing

- **Named Entity Recognition (NER):** Identifying names of people, places, organizations.
- **Information Extraction:** Extracting structured data from unstructured text.
- **Question Answering:** Finding key phrases relevant to queries.
- **Text Summarization:** Extracting essential information.

Ambiguity is a major challenge in NLP, affecting parsing and meaning interpretation. Dynamic programming-based parsers, such as CYK and Earley parsing, improve efficiency by avoiding redundant computations. Shallow parsing, or chunking, provides a fast and effective way to analyze sentence structures without full syntactic parsing, making it valuable for many NLP applications.

3.10 Probabilistic Context-Free Grammar (PCFG)

A **Probabilistic Context-Free Grammar (PCFG)** extends a traditional Context-Free Grammar (CFG) by assigning probabilities to production rules. This probabilistic approach helps in selecting the most likely parse when multiple valid parses exist.

3.10.1 Definition of PCFG

A PCFG is defined as a tuple $G = (N, \Sigma, R, S, P)$, where:

- N is a set of non-terminal symbols.
- Σ is a set of terminal symbols.
- R is a set of production rules of the form $A \rightarrow \alpha$.
- S is the start symbol.
- P is a probability function $P(A \rightarrow \alpha)$, assigning a probability to each production rule such that:

$$\sum_{\alpha} P(A \rightarrow \alpha) = 1, \quad \forall A \in N.$$

3.10.2 Example of PCFG

Consider the following PCFG:

$$S \rightarrow NPVP \quad (1.0)$$

$$NP \rightarrow Det N \quad (0.6) \quad | \quad Pronoun \quad (0.4)$$

$$VP \rightarrow V NP \quad (0.7) \quad | \quad V \quad (0.3)$$

$$Det \rightarrow \text{"the"} \quad (1.0), \quad N \rightarrow \text{"cat"} \quad (1.0)$$

$$V \rightarrow \text{"chases"} \quad (1.0), \quad Pronoun \rightarrow \text{"it"} \quad (1.0)$$

3.10.3 Parsing with PCFG

Given an input sentence, a parser computes the most probable parse tree using dynamic programming approaches like:

- **Viterbi Parsing:** Finds the most probable derivation using the probabilities of rules.
- **Inside-Outside Algorithm:** Used for training PCFGs based on observed data.

3.10.4 Applications of PCFG

- **Statistical Parsing:** Determines the most likely syntactic structure.
- **Speech Recognition:** Resolves syntactic ambiguities.
- **Machine Translation:** Helps in probabilistic phrase alignment.

3.11 Probabilistic CYK Parsing

The Probabilistic Cocke-Younger-Kasami (CYK) algorithm extends the standard CYK parsing method to handle PCFGs, selecting the most likely parse instead of just determining if a sentence is valid.

3.11.1 Algorithm Steps

1. Convert the grammar into Chomsky Normal Form (CNF).
2. Use a dynamic programming table, where $T[i, j, A]$ stores:

$$P(A \Rightarrow w_i \dots w_j)$$

i.e., the probability of deriving a substring using a non-terminal.

3. Fill the table bottom-up using:

$$P(A \rightarrow BC) = P(A \rightarrow BC) \cdot P(B \Rightarrow w_i \dots w_k) \cdot P(C \Rightarrow w_{k+1} \dots w_j)$$

4. Compute the highest probability parse tree.

3.11.2 Example of Probabilistic CYK

Consider parsing the sentence "the cat sleeps" using a PCFG. The CYK table stores probabilities at each step:

	the		cat		sleeps	

S	0.0		0.0		0.4	
NP	0.6		0.6		0.0	
VP	0.0		0.0		0.7	

3.11.3 Advantages of Probabilistic CYK

- Finds the most probable parse tree rather than just any valid parse.
- Handles ambiguities by ranking parses based on probability.
- Used in speech and handwriting recognition where multiple interpretations exist.

3.12 Probabilistic Lexicalized Context-Free Grammars (PLCFG)

A Probabilistic Lexicalized CFG (PLCFG) extends PCFG by incorporating lexical (word-level) information into rules. This enhances parsing accuracy, especially for complex languages.

3.12.1 Lexicalization in CFGs

Standard CFGs treat words as atomic symbols, ignoring their dependencies. Lexicalized CFGs overcome this by:

- Associating head words with non-terminals.
- Encoding word dependencies to improve parsing.
- Using probabilities to select the best parse.

3.12.2 Example of Probabilistic Lexicalized CFG

A normal PCFG rule:

$$VP \rightarrow VNP \quad (0.7)$$

Becomes lexicalized as:

$$VP(chases) \rightarrow V(chases)NP(cat) \quad (0.7)$$

Where *chases* is the head of VP and *cat* is the head of NP.

3.12.3 Parsing with PLCFG

The parsing process involves:

- Identifying head words during parsing.
- Assigning probabilities based on both syntactic and lexical features.
- Using dependency information for better accuracy.

3.12.4 Advantages of PLCFG

- **Improves Accuracy:** Captures real-world language dependencies.
- **Better Ambiguity Resolution:** Context-aware parsing improves disambiguation.
- **Used in Probabilistic Parsing Models:** Incorporated in dependency parsers and probabilistic phrase structure parsers.

Probabilistic methods enhance traditional CFGs by incorporating statistical likelihoods, improving parsing efficiency and ambiguity resolution. PCFG provides a probabilistic framework for syntactic parsing, Probabilistic CYK extends CYK parsing for PCFGs, and PLCFGs further refine parsing by incorporating lexical information.

3.13 Feature Structures

Feature structures are a formalism used in computational linguistics and natural language processing (NLP) to represent linguistic information in a structured manner. Unlike traditional context-free grammar (CFG) rules, which operate on atomic symbols, feature structures encode syntactic, semantic, and morphological properties as attribute-value pairs.

3.13.1 Definition

A feature structure is a set of attribute-value pairs, where:

- Attributes (features) specify properties like category, number, tense, case, etc..
- Values can be atomic symbols (e.g., singular, past) or nested feature structures.

Formally, a feature structure is represented as:

$$\{\text{category} : \text{NP}, \text{number} : \text{singular}, \text{case} : \text{nominative}\}$$

3.13.2 Example of Feature Structure

Consider the noun phrase "the cat":

$$\left\{ \begin{array}{ll} \text{category} & : \text{NP} \\ \text{number} & : \text{singular} \\ \text{determiner} & : \text{the} \\ \text{head} & : \text{cat} \end{array} \right\}$$

3.13.3 Advantages of Feature Structures

- **Expressiveness:** Captures fine-grained linguistic properties.
- **Constraint Satisfaction:** Feature structures allow enforcing agreement constraints (e.g., subject-verb agreement).
- **Extensibility:** Supports complex linguistic representations, including lexical-functional grammar (LFG) and head-driven phrase structure grammar (HPSG).

3.14 Unification of Feature Structures

Unification is a fundamental operation in computational linguistics that combines two feature structures by merging their information while ensuring consistency.

3.14.1 Definition

Unification takes two feature structures F_1 and F_2 and produces a new structure F such that:

$$F = F_1 \cup F_2$$

if and only if no conflicting attributes exist.

3.14.2 Example of Unification

Consider two feature structures:

Feature Structure 1 (for "the"):

$$\left\{ \begin{array}{ll} \text{category} & : \text{Det} \\ \text{number} & : \text{singular} \end{array} \right\}$$

Feature Structure 2 (for "cat"):

$$\left\{ \begin{array}{ll} \text{category} & : \text{N} \\ \text{number} & : \text{singular} \\ \text{head} & : \text{cat} \end{array} \right\}$$

The unification result:

$$\left\{ \begin{array}{ll} \text{category} & : \text{NP} \\ \text{number} & : \text{singular} \\ \text{head} & : \text{cat} \\ \text{determiner} & : \text{the} \end{array} \right\}$$

3.14.3 Unification Algorithm

1. Compare attributes in both feature structures.
2. If an attribute is present in both, ensure their values are compatible:
 - If values are identical, keep them.
 - If one value is more specific, prefer the more specific value.
 - If values are inconsistent, unification fails.
3. Merge the resulting attributes into a new structure.

3.14.4 Applications of Feature Structure Unification

- **Parsing:** Used in unification-based grammars such as HPSG.
- **Machine Translation:** Ensures consistent agreement and feature matching.
- **Speech Processing:** Helps in phonological and morphological analysis.

Feature structures provide a flexible mechanism for representing linguistic information, and unification serves as a powerful tool for ensuring consistency in syntactic and semantic processing. These concepts play a crucial role in NLP frameworks, particularly in grammar formalisms like LFG and HPSG.

Syntactic analysis is crucial in natural language processing, enabling applications such as machine translation, speech recognition, and syntactic parsing.