

**POST GRADUATE  
PROGRAM IN  
GENERATIVE AI  
AND ML**

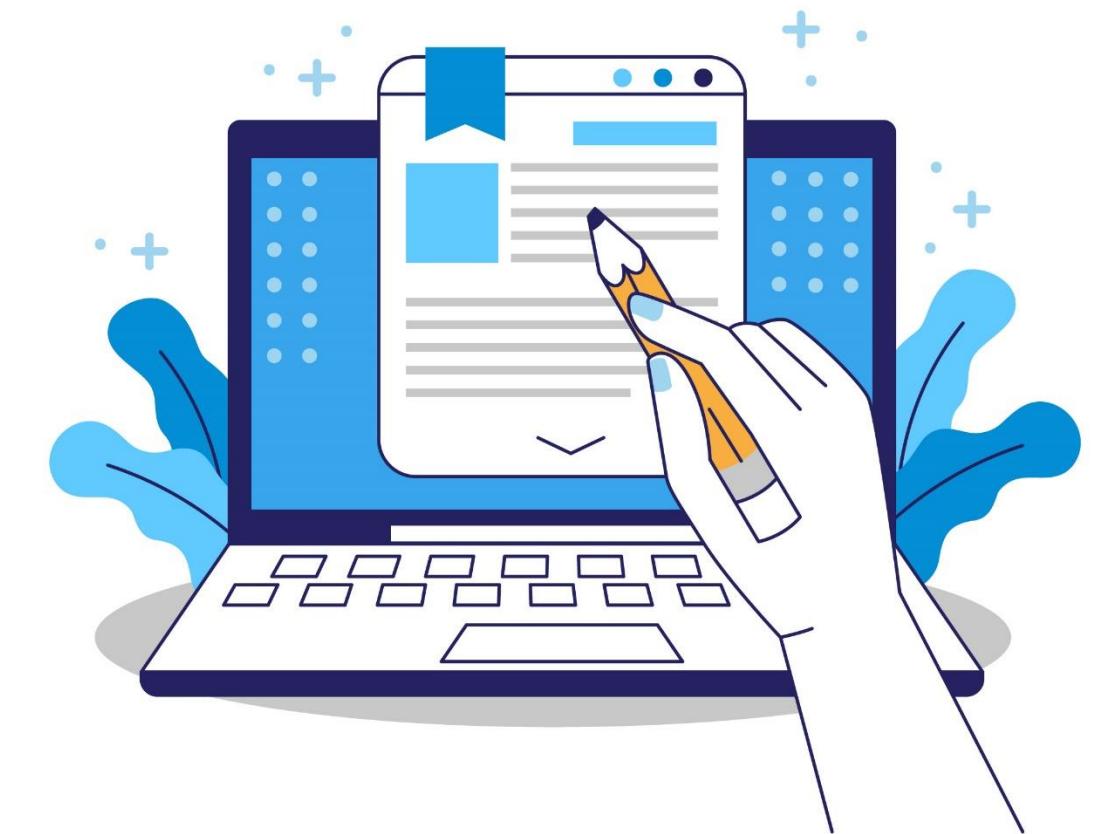
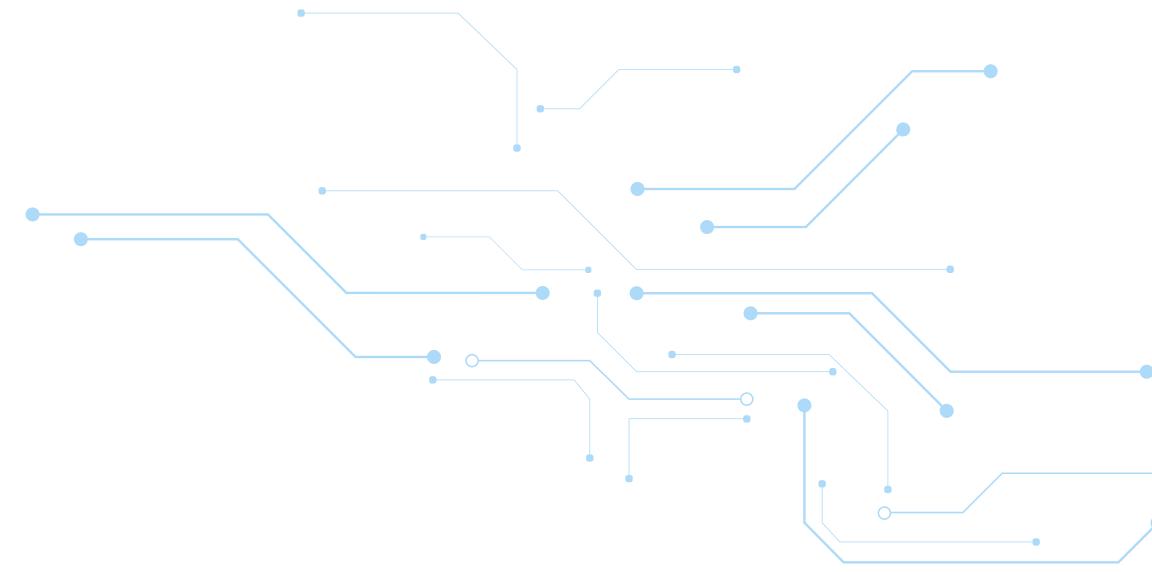
**Natural Language  
Processing**



# Text Processing and Feature Engineering

# Topics

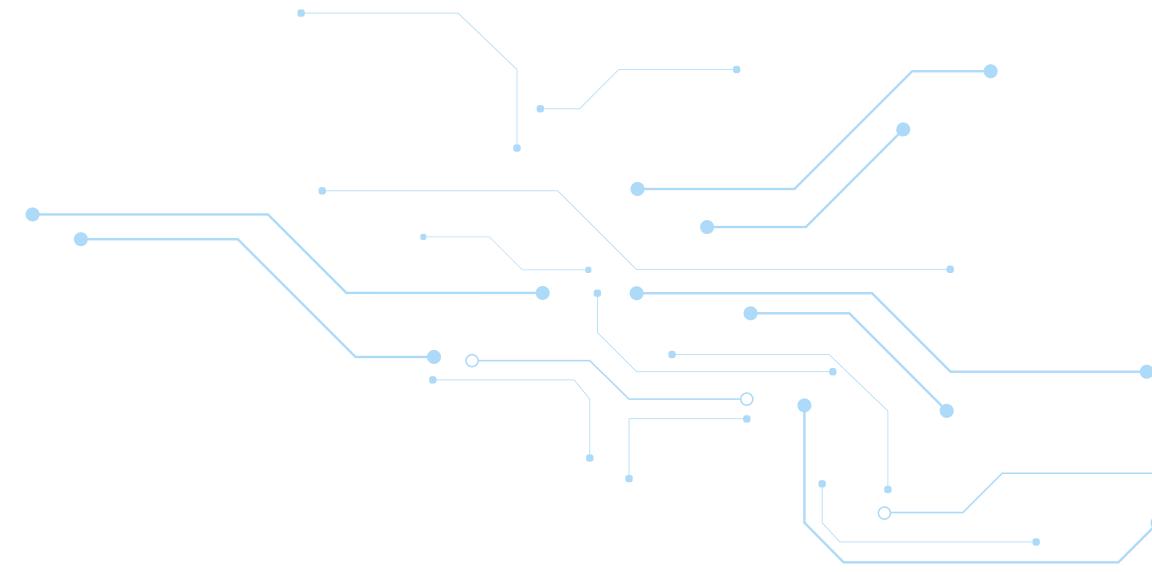
- e! Stemming Techniques and Applications in NLP
- e! Lemmatization Methods and Real-World Relevance
- e! Part-of-Speech Tagging and Its Role in NLP
- e! Text Representation Techniques: BoW and TF-IDF
- e! Word Embeddings: Word2Vec, GloVe, and FastText
- e! Processing Code-Mixed and Multilingual Text
- e! Text Classification and Feature Selection Methods



# Learning Objectives

By the end of this lesson, you will be able to:

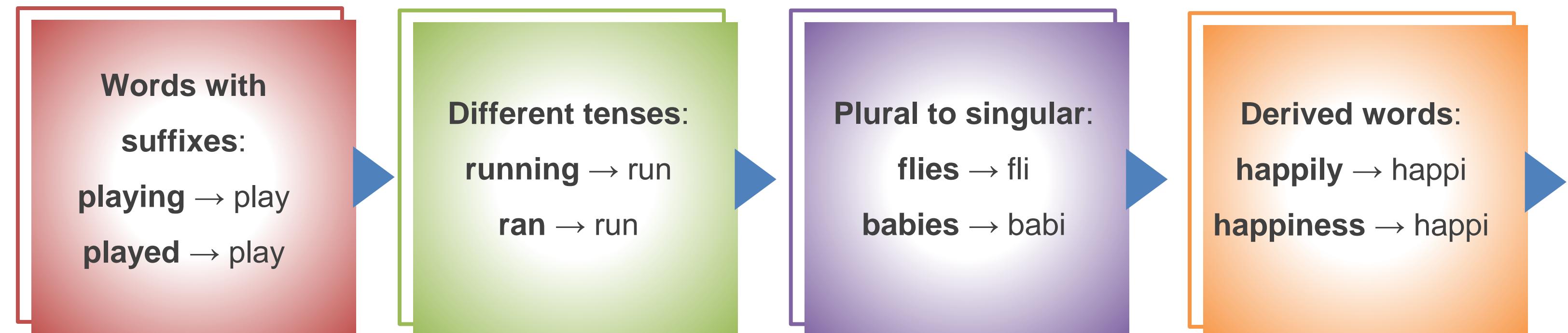
- e! Describe the principles of stemming, lemmatization, and part-of-speech tagging in NLP workflows.
- e! Demonstrate the use of traditional text representation techniques like Bag of Words and TF-IDF.
- e! Evaluate the effectiveness of word embedding models such as Word2Vec, GloVe, and FastText.
- e! Apply text classification techniques to multilingual data using appropriate feature selection methods.



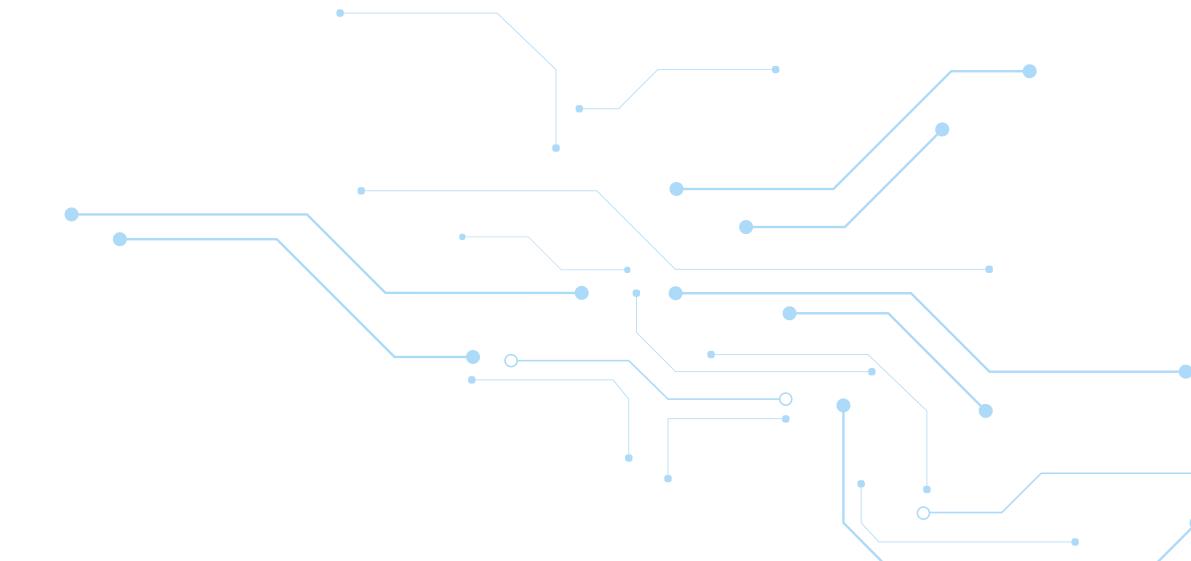
# Stemming

# What is Stemming ?

Stemming is a simpler alternative to lemmatization, using handcrafted rules to strip word endings and reduce them to a common base form known as stems.



# Types of Stemming



## Rule-based Stemming

Example:

"running" → "run"  
(remove "ing")

"happily" → "happi"  
(remove "ly")

## Porter Stemming

Example:

"running" → "run"  
"fishing" → "fish"

## Lancaster Stemming

Example:

"running" → "run"  
"loved" → "lov"

## Snowball Stemming

Example:

"happiness" → "happi"  
"playing" → "play"

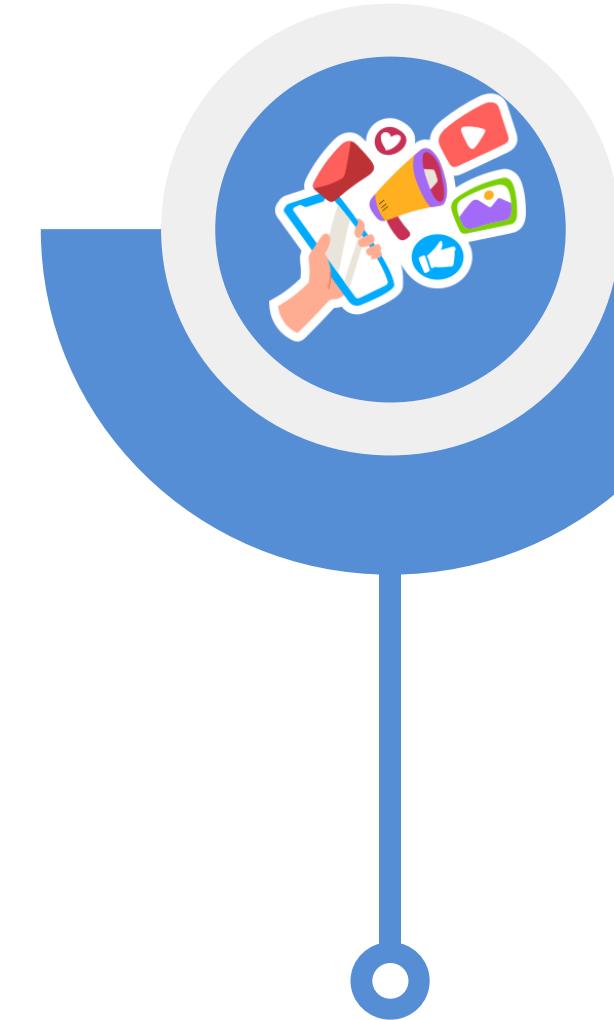
01

02

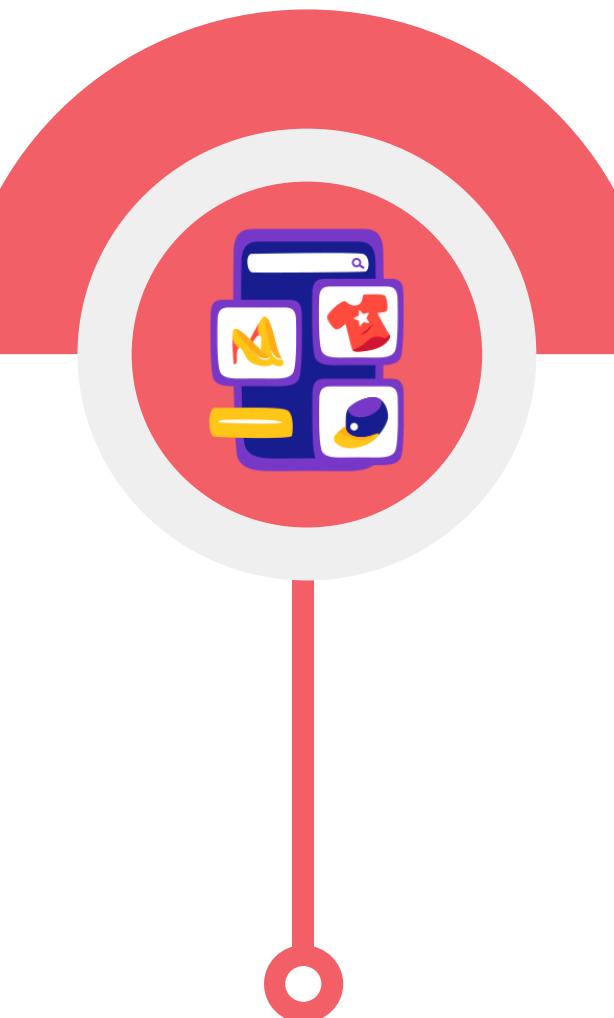
03

04

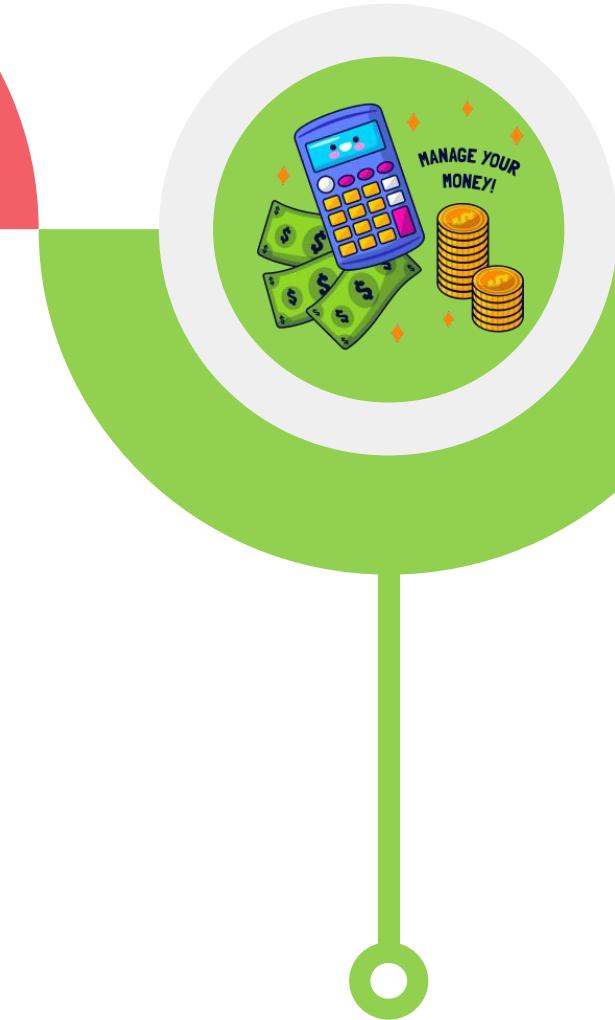
# Applications of Stemming



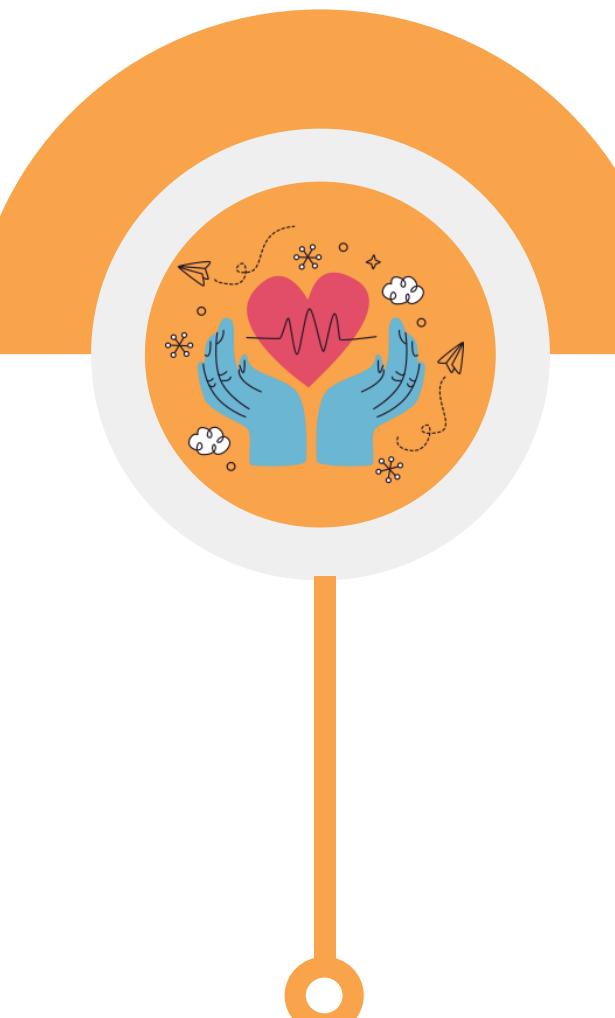
Social Media Analysis



E-commerce Industry



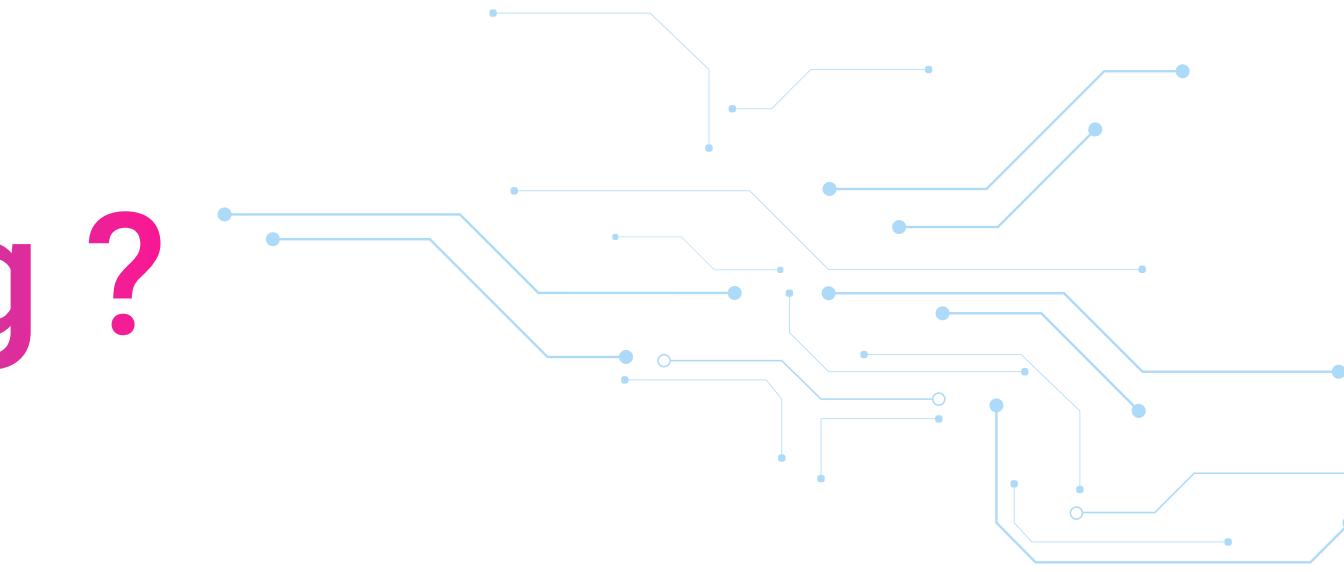
Finance Industry



Healthcare Industry



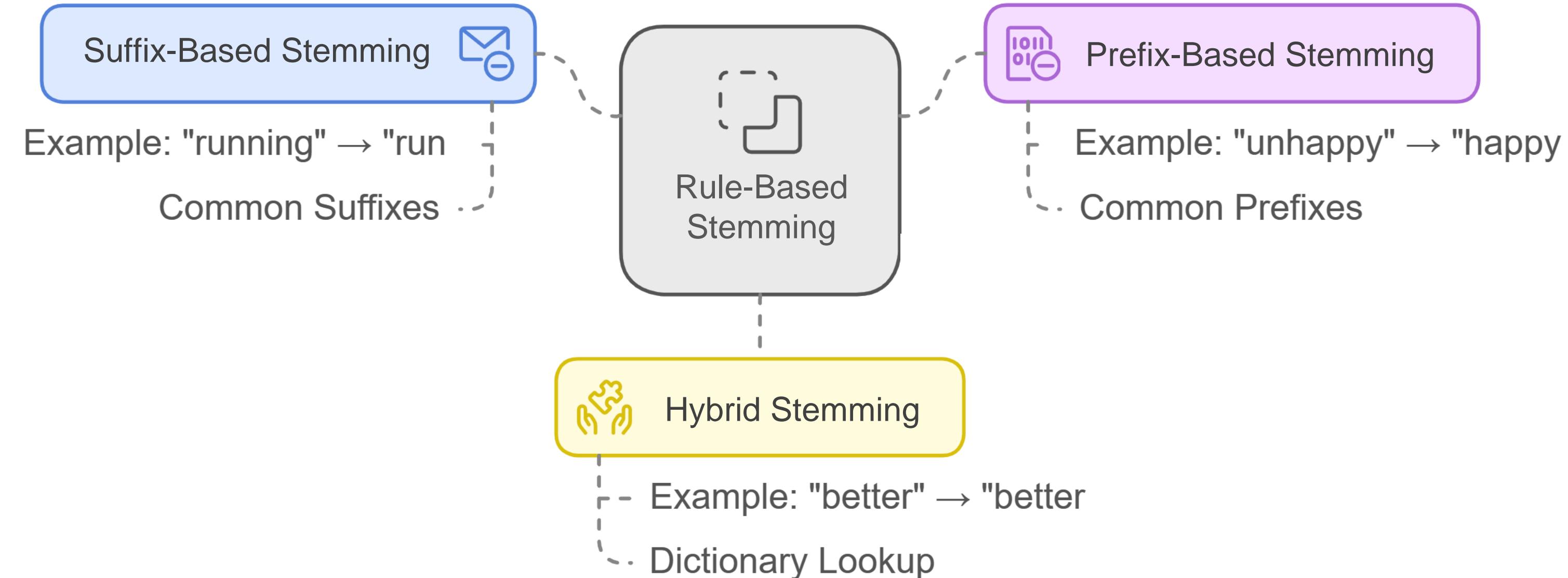
# What is Rule-Based Stemming ?



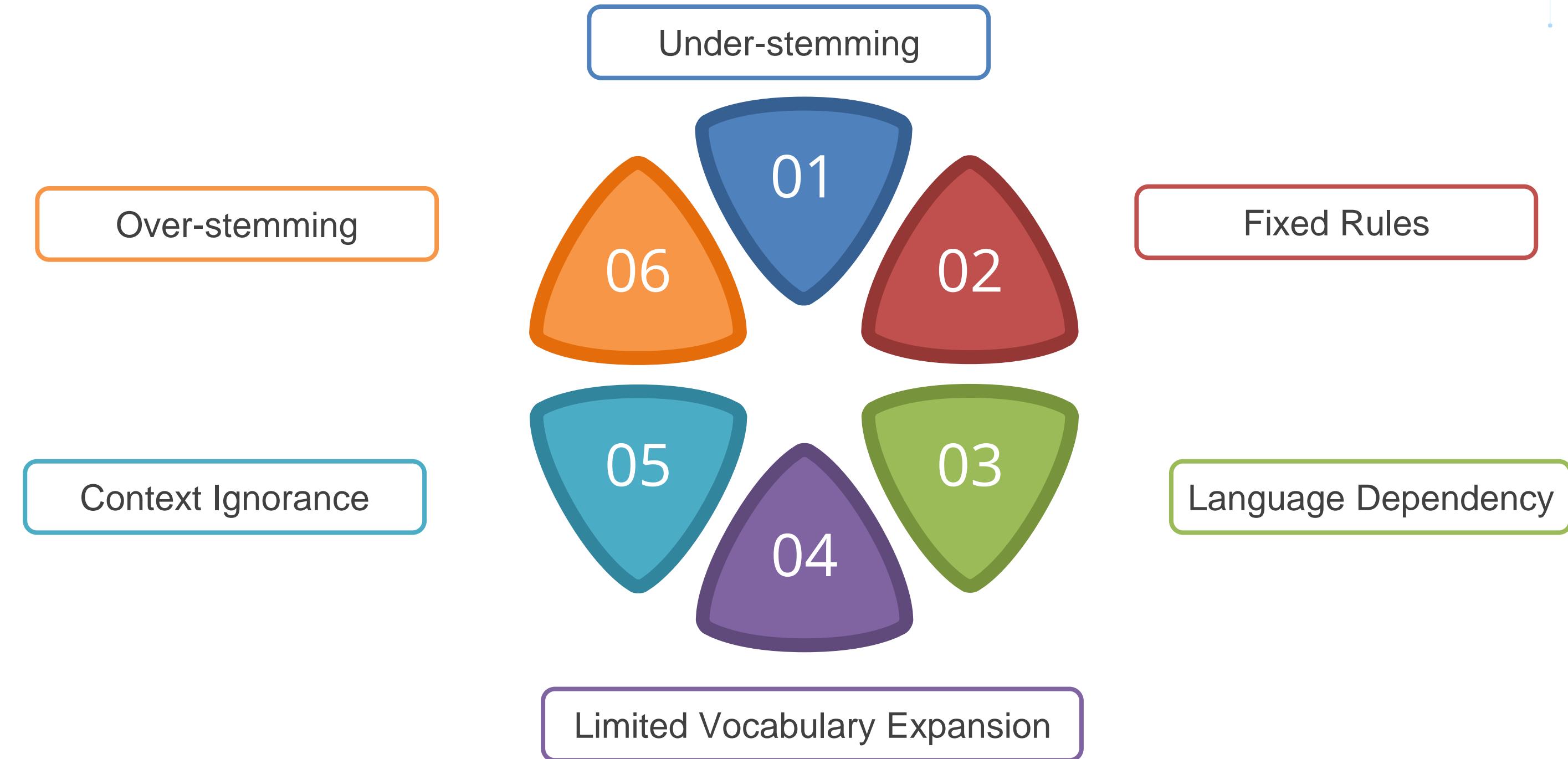
Rule-based stemming uses **predefined rules** to remove **suffixes or prefixes** from words to find their root (stem). These rules are often language-specific and focus on common affixes like **-ing**, **-ed**, **-s**, etc.



# Types of Rule-Based Stemming



# Limitations of Rule-Based Stemming

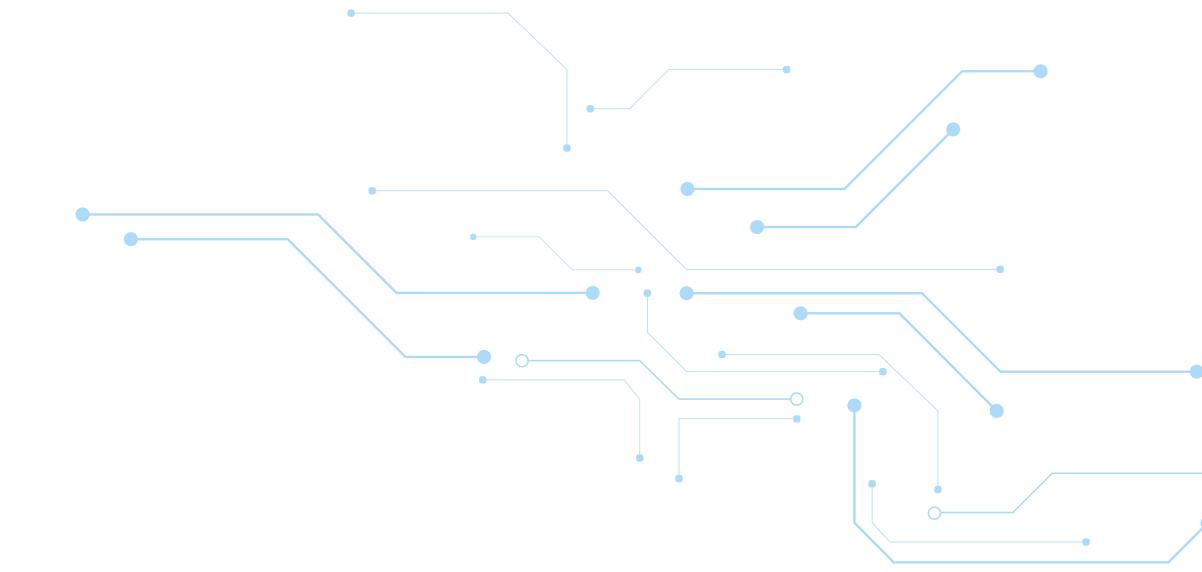


# Porter Stemmer



- e!** The Porter Stemmer is one of the most widely used stemming algorithms.
- e!** This algorithm applies a set of rules to identify and remove suffixes from words to extract their stems.
- e!** For instance, the Porter algorithm transforms the word "leaping" into "leap" by stripping the "ing" suffix.

# Snowball Stemmer



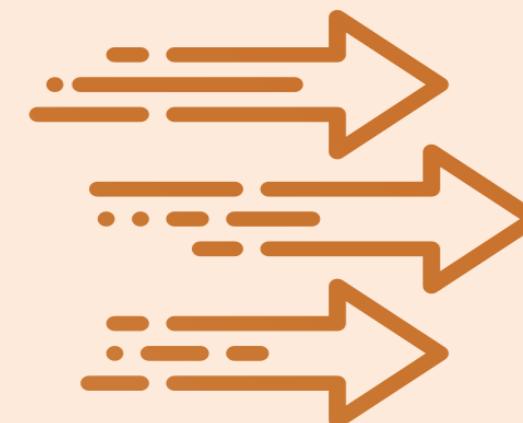
English: "hotels"

French: "*hôtel*"

Spanish: "*hoteles*"

German: "Hotels"

## Snowball Stemmer



English: "hotel"

French: "*hôtel*"

Spanish: "*hotel*"

German: "Hotel"

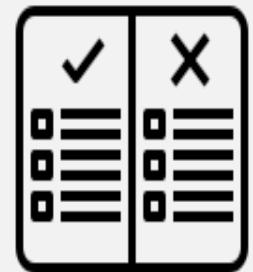
# Why Snowball Stem. over Porter Stem. ?

Aspect	Porter Stemmer	Snowball Stemmer
<b>Aggressiveness</b>	Less aggressive in reducing words.	More aggressive in stemming words.
<b>Fixes</b>	Contains some inconsistencies and issues.	Addresses many issues present in Porter.
<b>Working Differences</b>	Slightly older and simpler rules.	Updated and more refined rules.
<b>Example: 'fairly'</b>	Stemmed to 'fairli'.	Stemmed to 'fair'.
<b>Example: 'sportingly'</b>	Stemmed to 'sportingli'.	Stemmed to 'sport'.
<b>Accuracy</b>	Can produce less accurate stems.	Provides more accurate and usable stems.

# Lancaster Stemmer

01

The Lancaster Stemmer is a rule-based stemming algorithm that is known for its **simplicity and high aggressiveness** in reducing words to their base form.



It uses a recursive approach with a predefined set of **stripping rules** to iteratively remove suffixes from words.

02

03

While this makes it computationally efficient, it often results in **over-stemming**, producing very short stems that may **lose meaning**.



# LS, KS, and CAS

## Lovins Stemmer (LS)

The first stemming algorithm (1968), it removes suffixes in a **single step** using predefined rules, balancing speed and linguistic accuracy.

## Krovetz Stemmer (KS)

Validates stemmed words **against a dictionary** to ensure meaningful stems, avoiding over-stemming for improved accuracy.

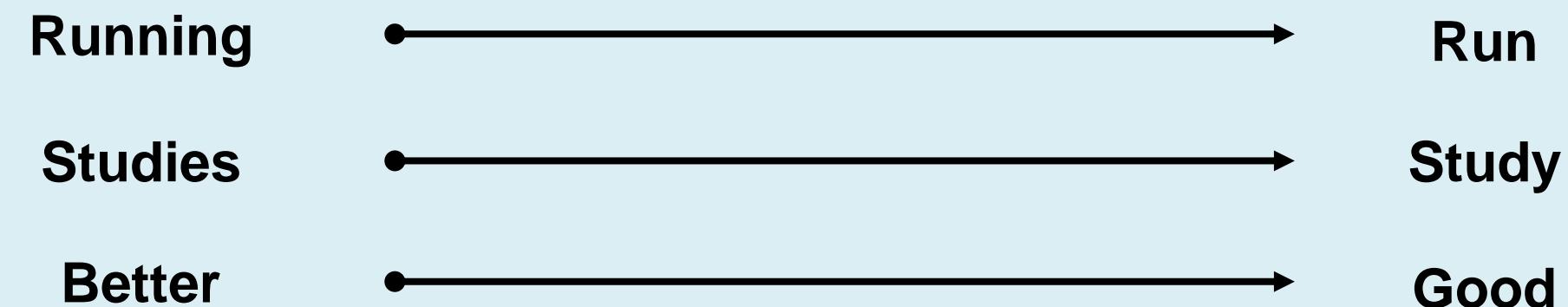
## Context-Aware Stemming (CAS)

Uses machine learning or rules to stem words based on **sentence context**, prioritizing precision but requiring high computational resources.

# Lemmatization

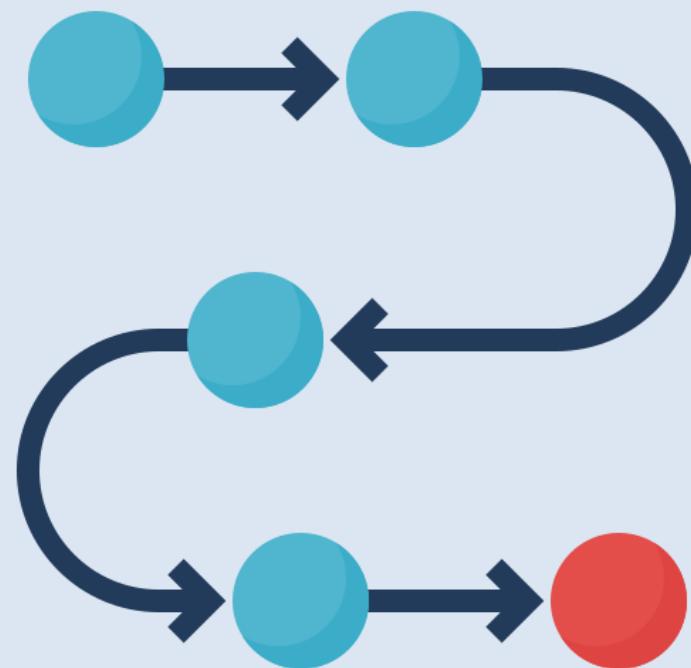
# Defining Lemmatization

Lemmatization is the process of reducing words to their dictionary form (lemma) while considering their context and part of speech.



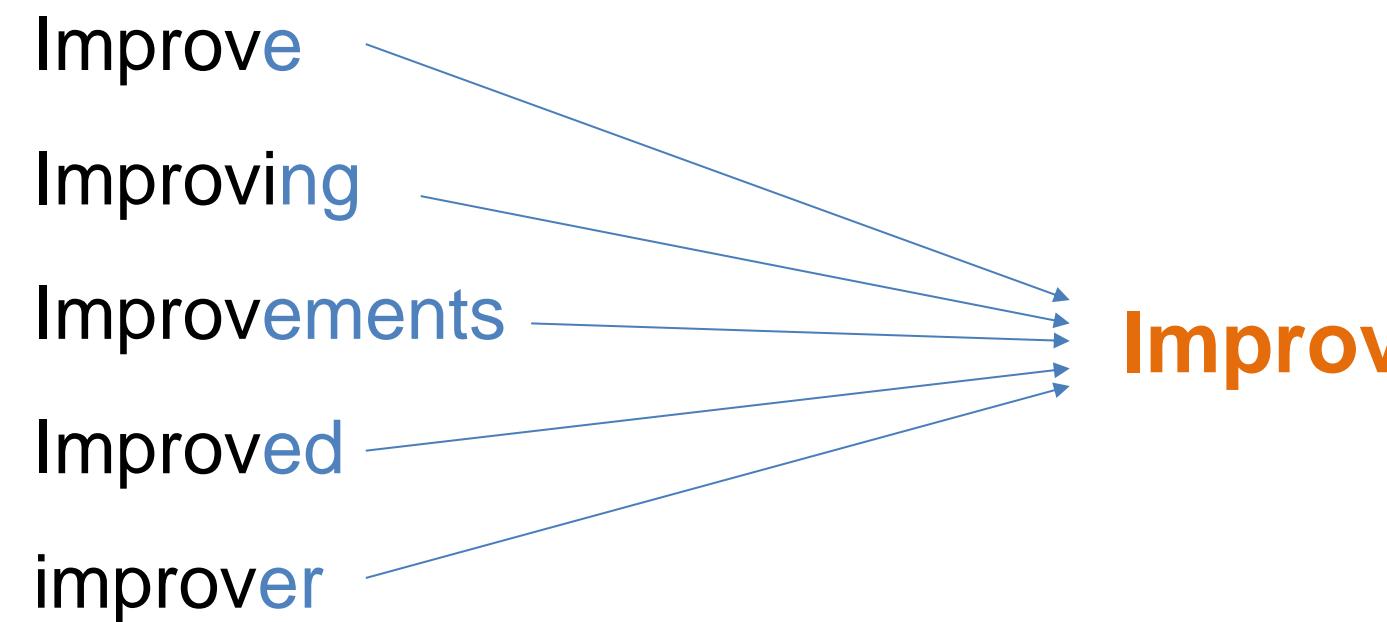
# Process of Lemmatization

- 01 Morphological Analysis
- 02 Part-of-Speech Tagging
- 03 Dictionary Lookup

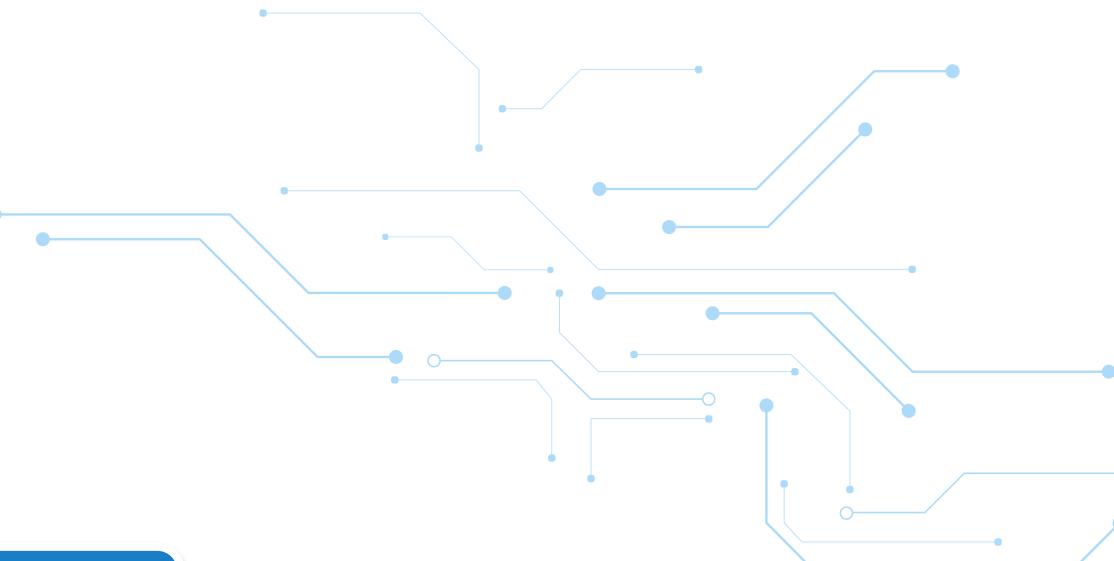
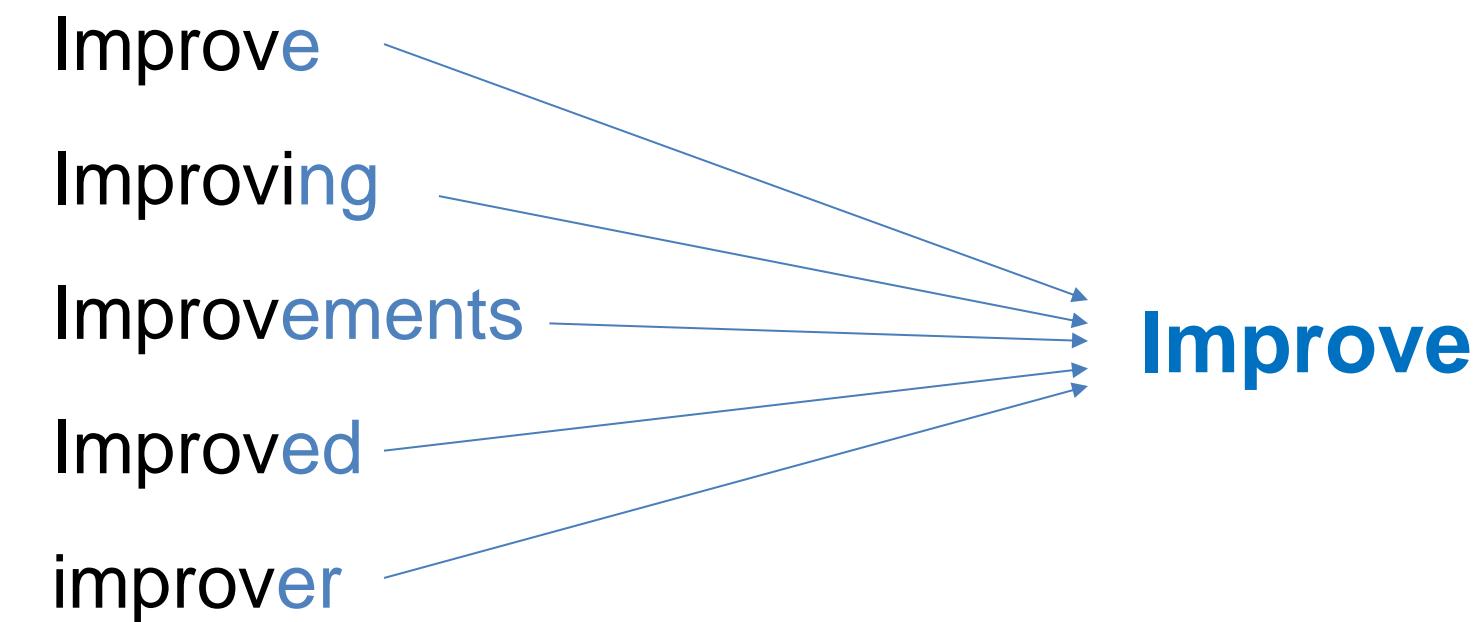


# Lemmatization vs. Stemming

Stemming



Lemmatization



# Lemmatization vs. Stemming Example

Word	Stemming Output	Lemmatization Output	Context
Caring	Car	Care	Verb (action)
Studies	Studi	Study	Plural noun
Running	Run	Run	Verb (action)
Better	Bet	Good	Comparative adjective
Children	Children	Child	Irregular plural noun

# Why Lemmatization is Important?

Standardizes Words

01.

Improves Search Results

Sentiment Analysis

03.

Machine Translation

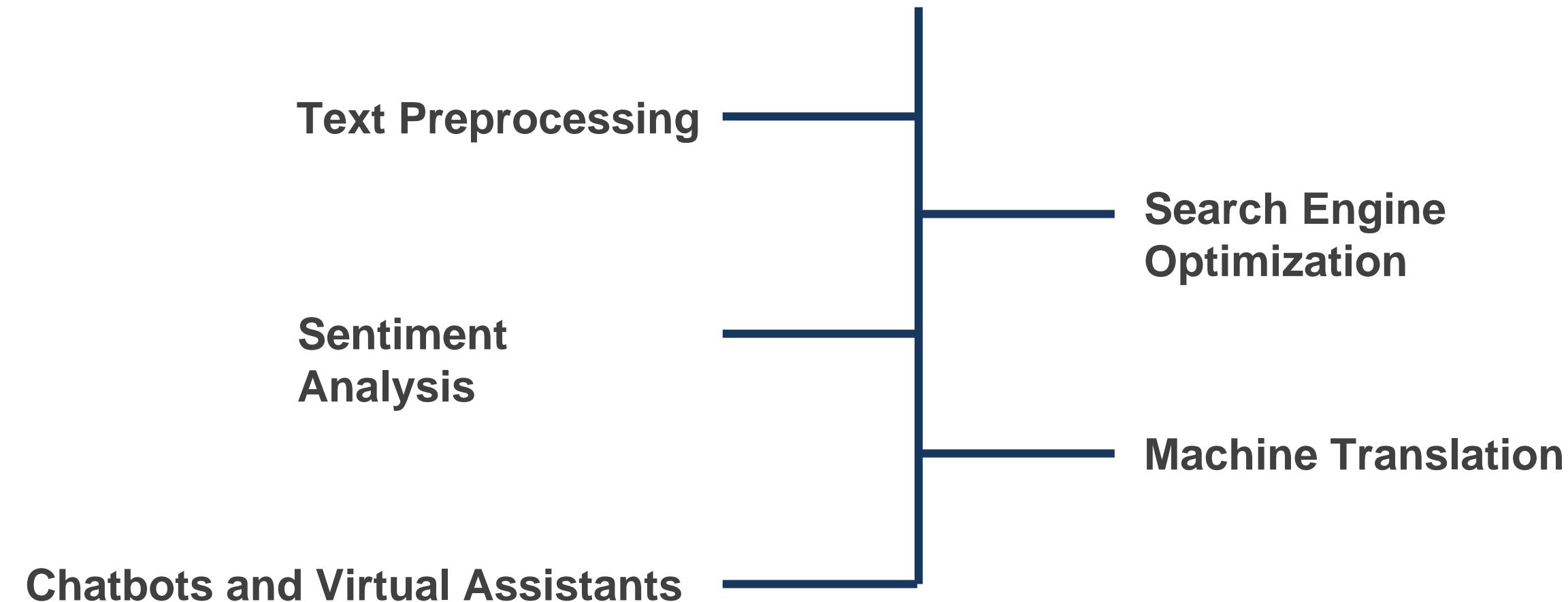
05.

Text Mining

# Lemmatization

# Overview of Lemmatization Applications

Lemmatization has wide-ranging applications across various fields of Natural Language Processing (NLP), search engines, and machine learning.



# Text Preprocessing in NLP

Lemmatization is a text normalization technique to prepare data for processing.

Reduces inflected words to their dictionary form (lemma).

Simplifies unstructured text for tasks like classification, clustering, and analysis.

**The children are playing happily.** —————→ **The child be play happy.**



# Search Engine Optimization

Lemmatization improves search engines by mapping words to their base forms.

User Queries: Variations in search terms are matched to a single lemma for better results.

**Query: “Running shoes” » Searches for “run shoes,” “ran shoes.”**

**Result: Search engines retrieve results for all related terms.**

# Sentiment Analysis

Lemmatization helps in sentiment analysis by reducing words to their base form for consistency.

Different forms of a word (e.g., happy, happier, happiness) are analyzed as a single entity.

Input Sentence	Lemmatized Words	Sentiment
"I loved the movie."	"I love the movie."	Positive
"He is running happily."	"He be run happy."	Positive

# Machine Translation

In machine translation, lemmatization simplifies word variations to improve translation.

Translates sentences by focusing on lemmas rather than inflected forms.

Input: “She studies daily.”



Lemmatized: “She study daily.”



Translation (to French):  
“Elle étudie tous les jours.”

# Chatbots and Virtual Assistants



Lemmatization enhances chatbot responses by normalizing user input.  
Helps identify user intent accurately by simplifying variations of words.

User Input: "I am running late."



Lemmatized: "I be run late."

Bot Response: "I understand. Do you need to reschedule?"



# Rule Based Lemmatization

Relies on predefined rules for word transformation.

Words are stripped of suffixes or prefixes based on grammatical rules.

Works well for inflectional morphology (tense, plural forms, etc.).

“ing” → Remove “ing” if preceded by a vowel (e.g., running → run).

“ed” → Replace with the root form (e.g., studied → study).



# Dictionary Based Lemmatization

Uses a predefined dictionary of words and their lemmas.

The input word is looked up in the dictionary to return its base form.

Input: "children"

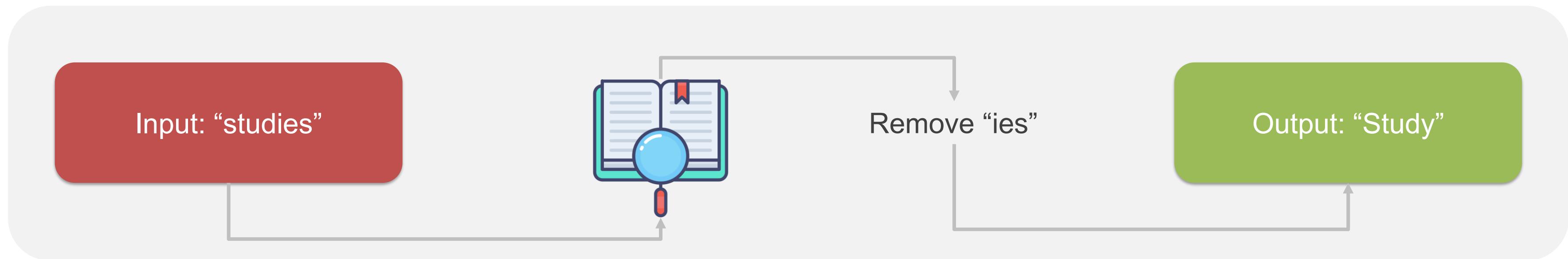
Dictionary Lookup: children → child

Output: child



# Hybrid Lemmatization

Combines rule-based and dictionary-based approaches for improved accuracy.



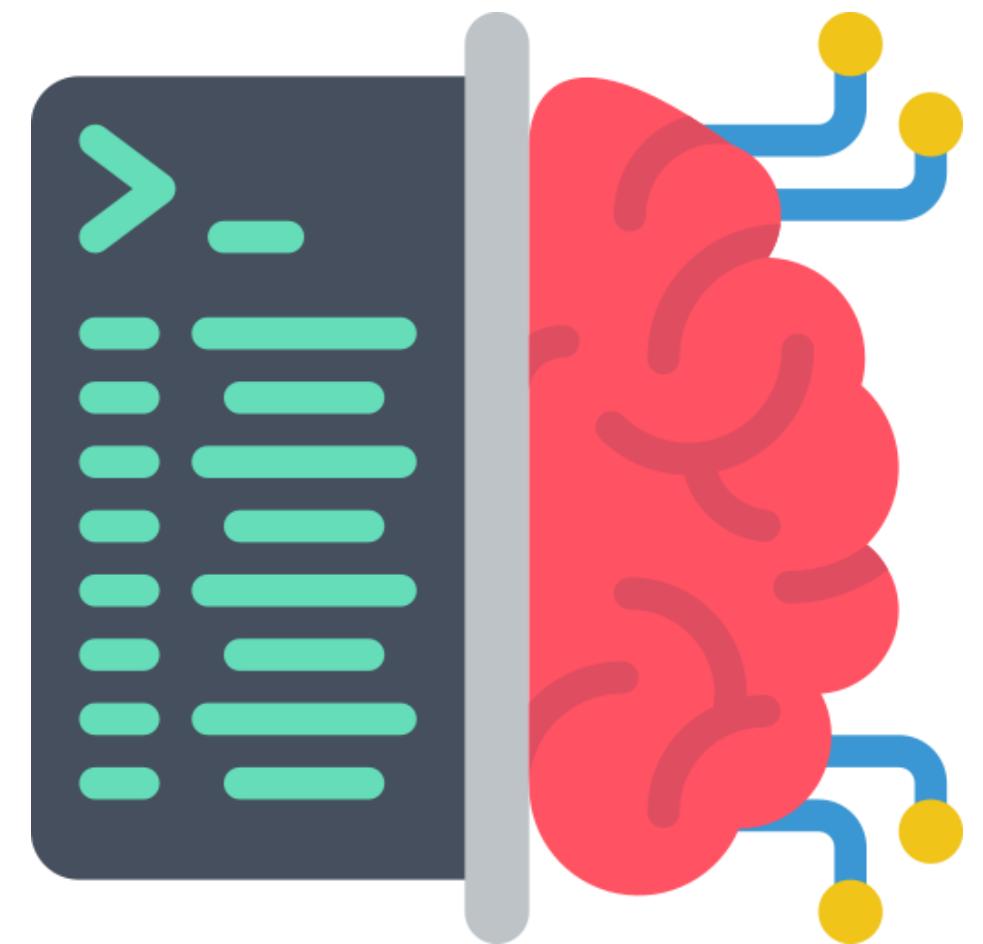
# Machine Learning Based Lemmatization

Uses machine learning models to predict lemmas based on training data.

Trained on large text corpora with words and their corresponding lemmas.

Studies → study

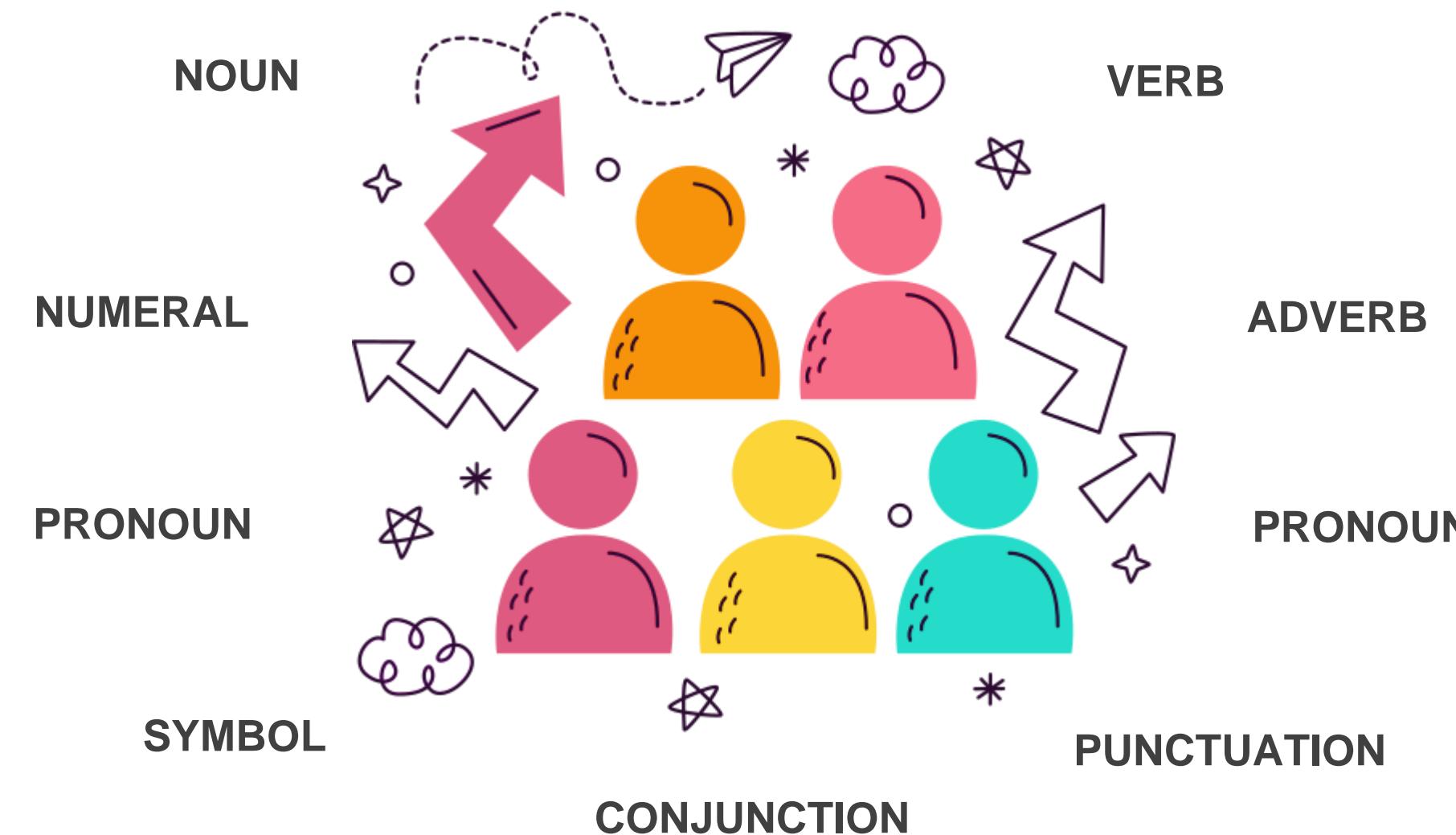
went → go



# Part-of-Speech (POS) Tagging

# What is POS ?

Part-of-Speech (POS) tagging is a fundamental technique in natural language processing that assigns grammatical labels—such as noun, verb, adjective, adverb, or pronoun—to individual words within a sentence.



# How POS Tagging Works?

For the sentence: "*The cat sleeps on the mat.*"

POS tags could be:

The → Determiner (DT)

cat → Noun (NN)

sleeps → Verb (VBZ)

on → Preposition (IN)

the → Determiner (DT)

mat → Noun (NN)



# Text Representation: Bag of Words (BoW), TF-IDF

# Introduction to Bag of Words - BoW

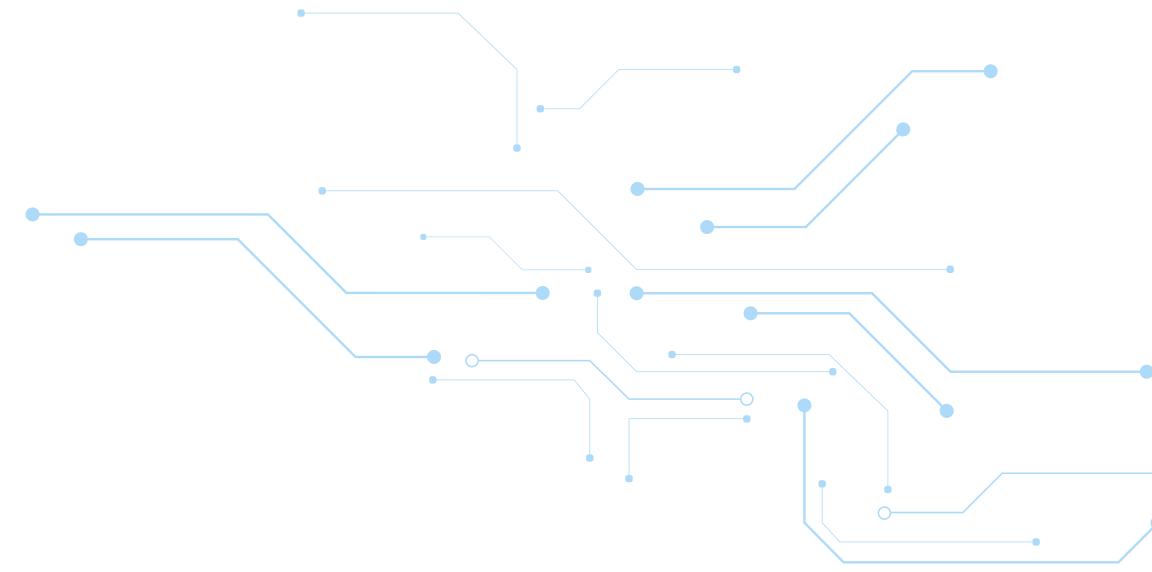
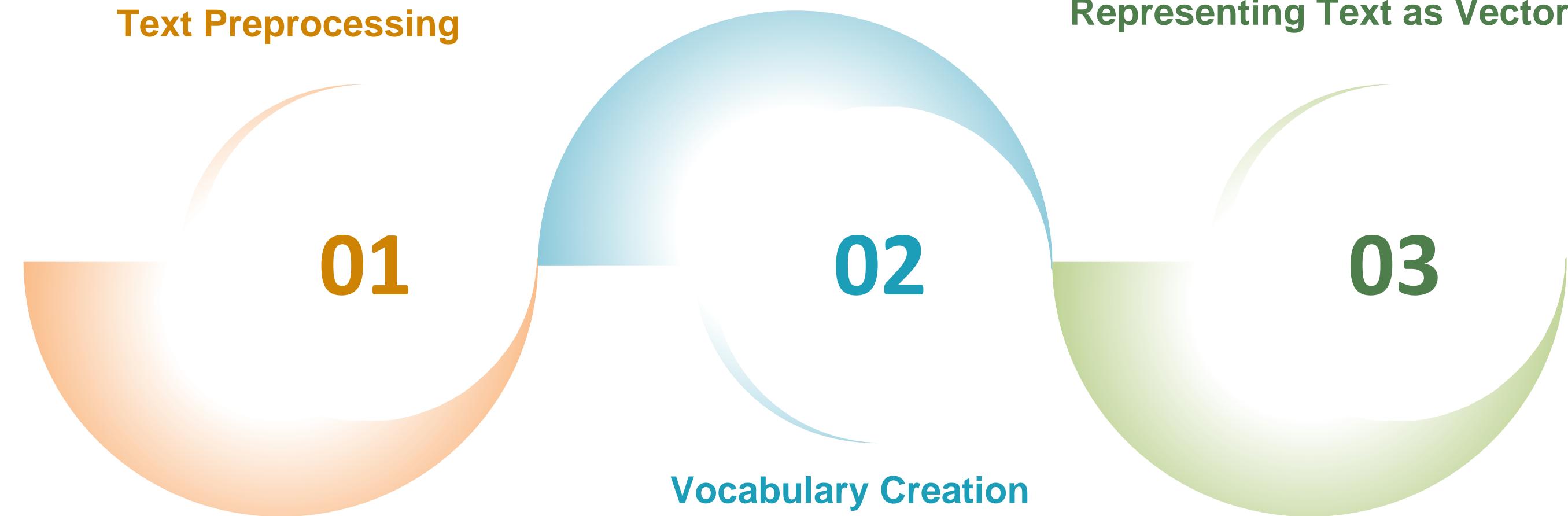
Bag of Words (BoW) is an NLP technique that represents text by converting it into a set of unique words, ignoring grammar and word order, and counting word occurrences.

The cat is on the mat.



the: 2, cat: 1, is: 1, on: 1, mat: 1

# Steps in Bag of Words - BoW



# Representation of Bag of Words

Suppose we have the following two sentences in a document:

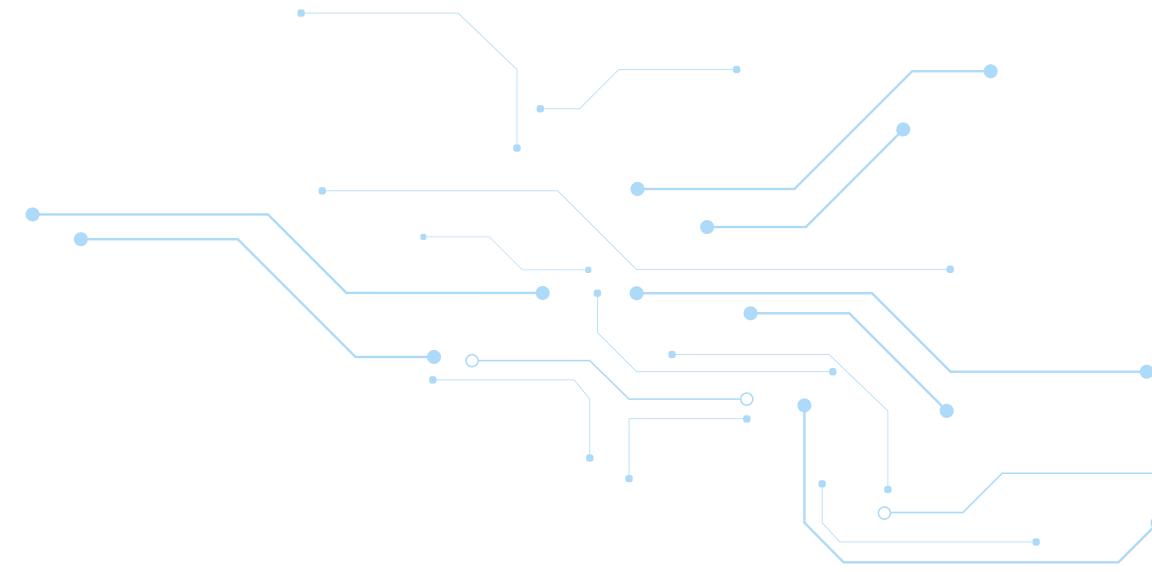
e! "The cat sat on the mat."

e! "The dog ate the cat."

["the", "cat", "sat", "on", "mat", "dog", "ate"]

Word	The	Cat	Sat	On	Mat	Dog	Ate
Sentence 1	2	1	1	1	1	0	0
Sentence 2	2	1	0	0	0	1	1

# Types of BoW Representation



## Binary BoW

Words are represented as binary values: 1 if present, 0 otherwise.

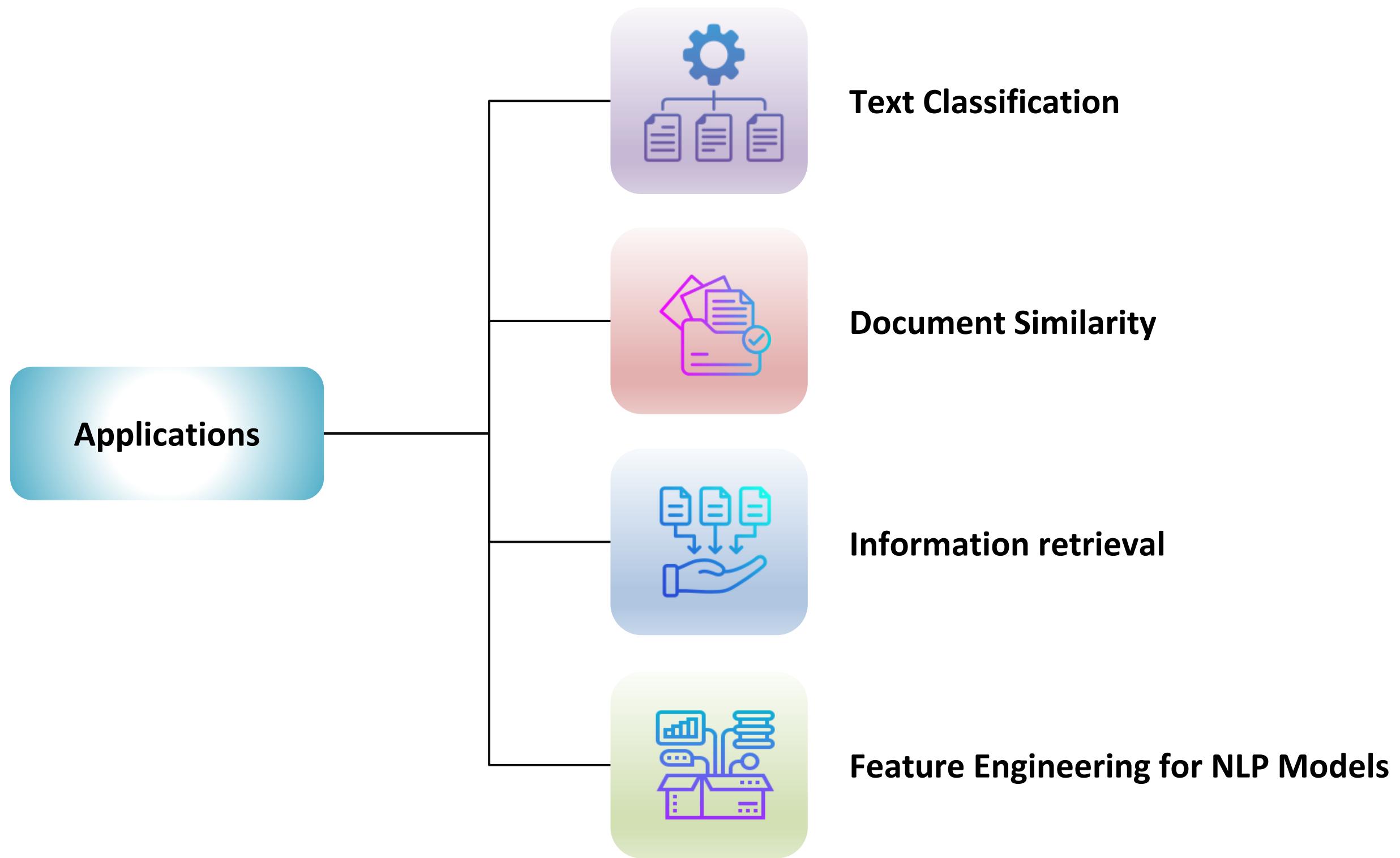
## Count-Based BoW

Words are represented by their raw occurrence counts.

## Weighted BoW

Assigns weights to words based on their importance

# Application of BoW



# Calculating TF-IDF

Document 1: "The quick brown fox jumped over the lazy dog."

Document 2: "The fast fox jumped over the quick dog."

Document 3: "The brown dog barked loudly."

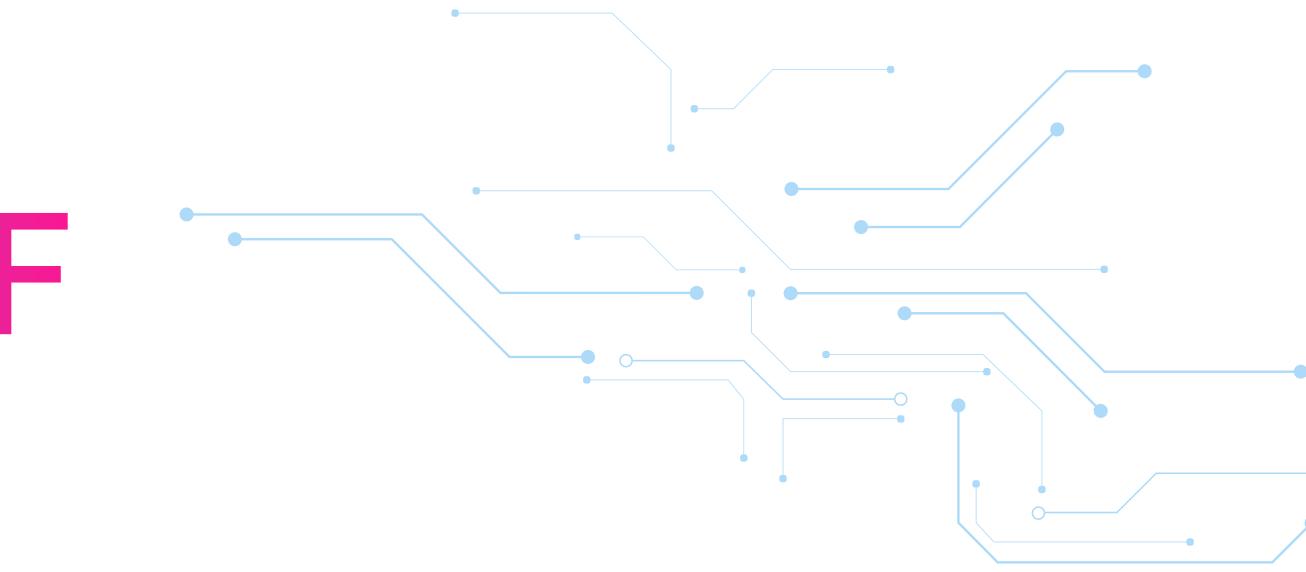
"fox" appears once in Document 1, which has a total of 9 words.TF for "fox" in Document 1 =  $1/9 = 0.111$ .

"fox" appears in Document 1 and Document 2 (2 out of 3 documents).IDF for "fox" =  $\log(3/2)\log(3/2) \approx 0.176$ .

$$\text{TF-IDF}(t,d) = \text{TF}(t,d) \times \text{IDF}(t)$$

TF-IDF for "fox" in Document 1 =  $\text{TF} \times \text{IDF} = 0.111 \times 0.176 \approx 0.0195$ .

# Key Steps in Computing TF-IDF



## Preprocess the Text

Tokenize, remove stop words, lowercase, and apply stemming/lemmatization.

## Calculate TF for Each Term

Compute term frequency for all words in the document.

## Compute IDF for Each Term

Use the entire corpus to calculate how rare each word is.

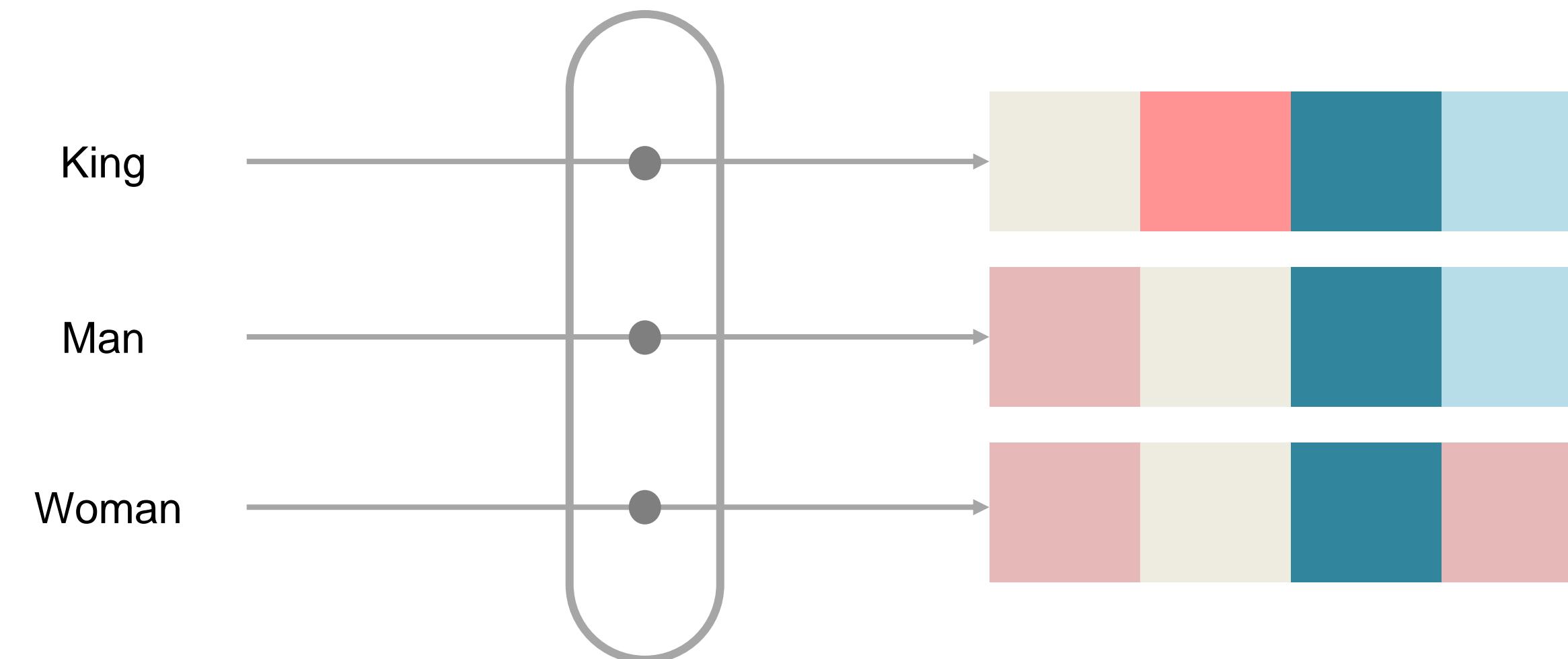
## Multiply TF by IDF

Generate the final TF-IDF score for each word in the document.

# Word Embeddings: Word2Vec, GloVe, FastText

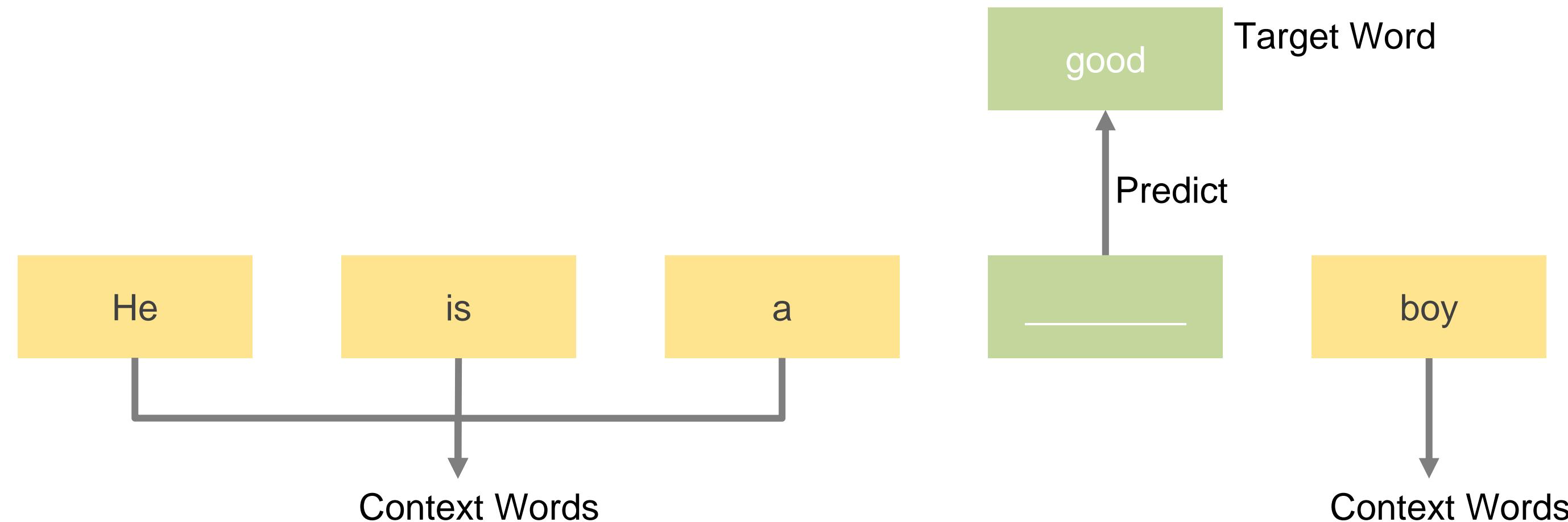
# Introduction to Word2Vec

Word2Vec generates word embeddings that capture semantic relationships



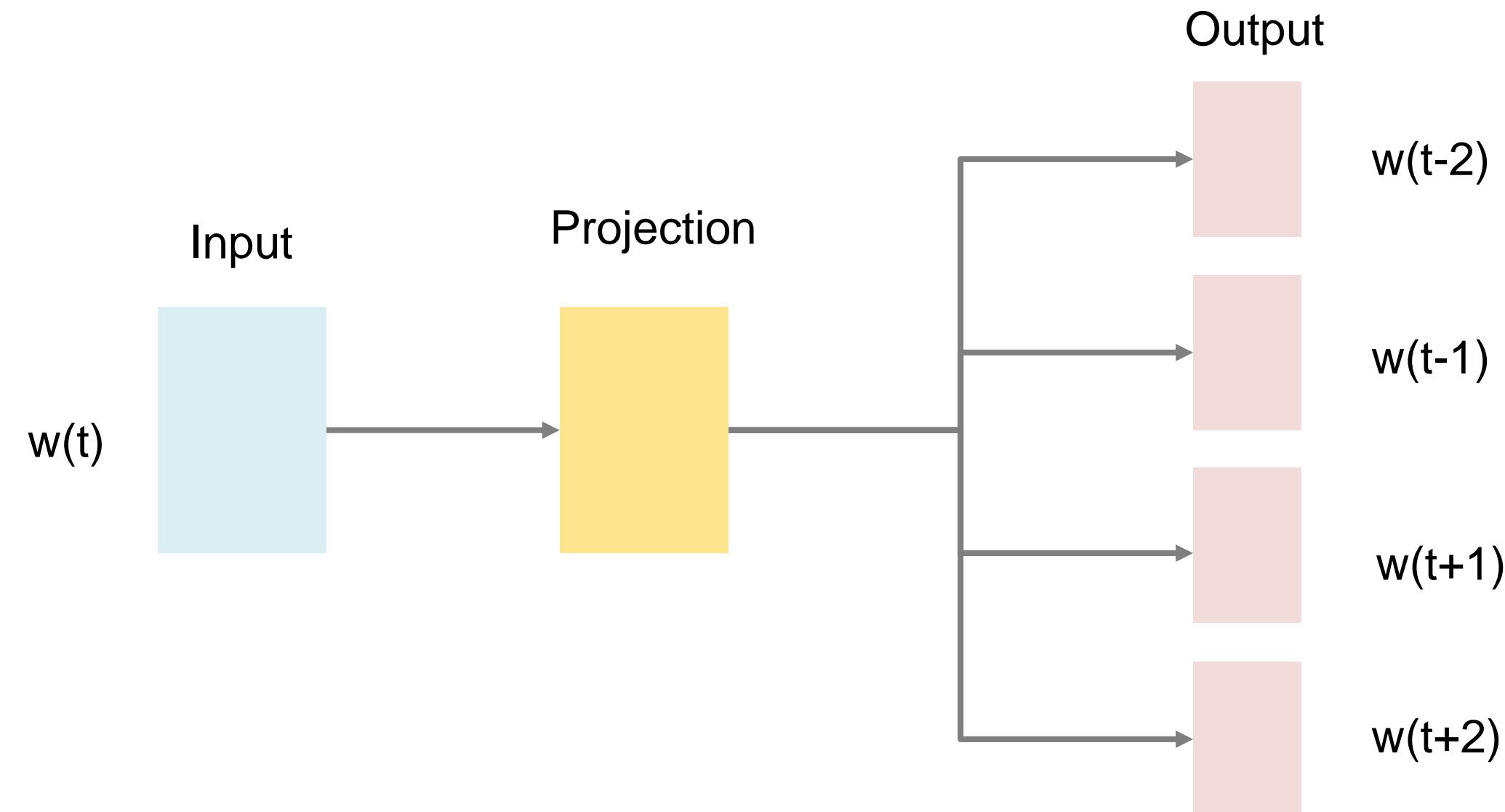
# CBOW Model

Continuous Bag of Words (CBOW) is a neural network model that learns word embeddings by predicting a target word based on its context words



# Skip-Gram Model

Skip-gram is a popular algorithm used in natural language processing (NLP) to create word embeddings



# Negative Sampling and Hierarchical Softmax

**Problem:** High computational cost of Softmax for large vocabularies.

## **Negative Sampling**

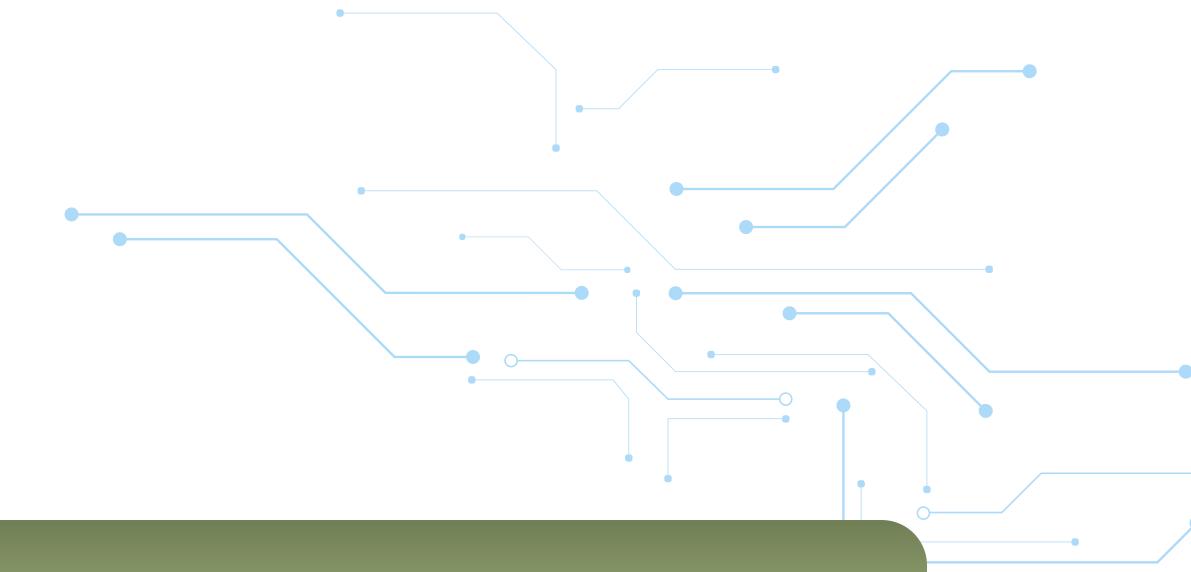
Samples a few "negative" examples (random words) to reduce computations.

## **Hierarchical Softmax**

Uses a binary tree structure for faster computation of probabilities.

# Introduction to GloVe

Global Vectors for Word Representation (GloVe), is an unsupervised learning algorithm that generates vector representations, or embeddings, of words.



## Limitations

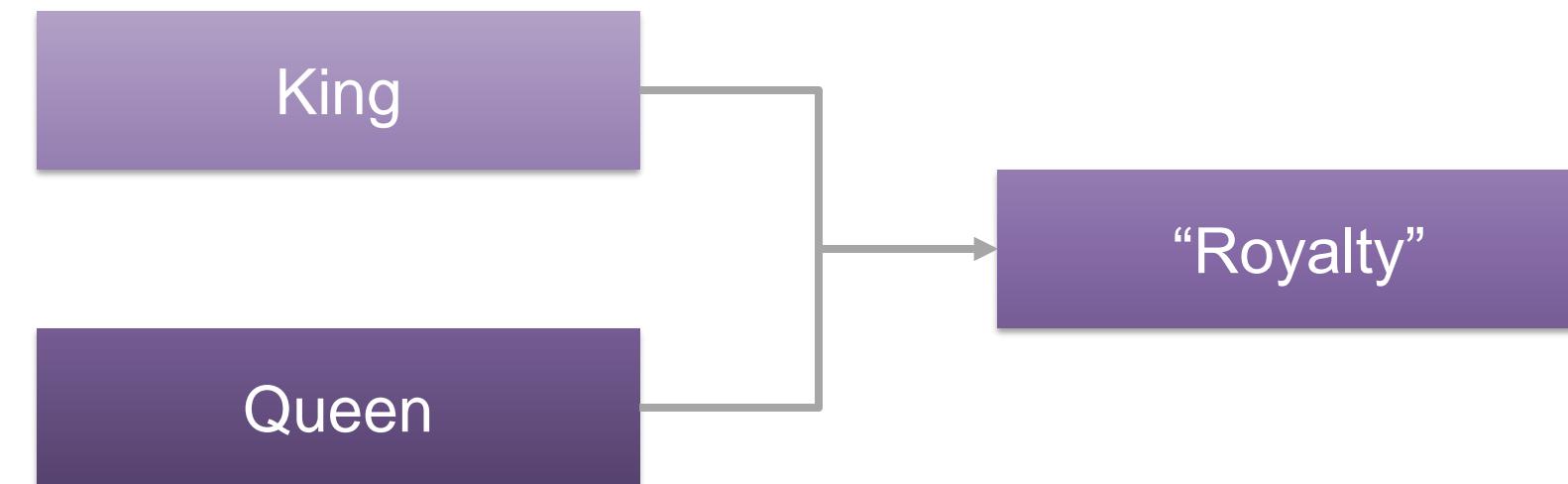
Focuses only on local context (neighboring words).

## GloVe's Approach

Combines local context with global statistical information.

# Role of Co-occurrence Matrix

A matrix where each entry  $M(i,j)$  represents how often word  $i$  co-occurs with word  $j$  in the corpus.

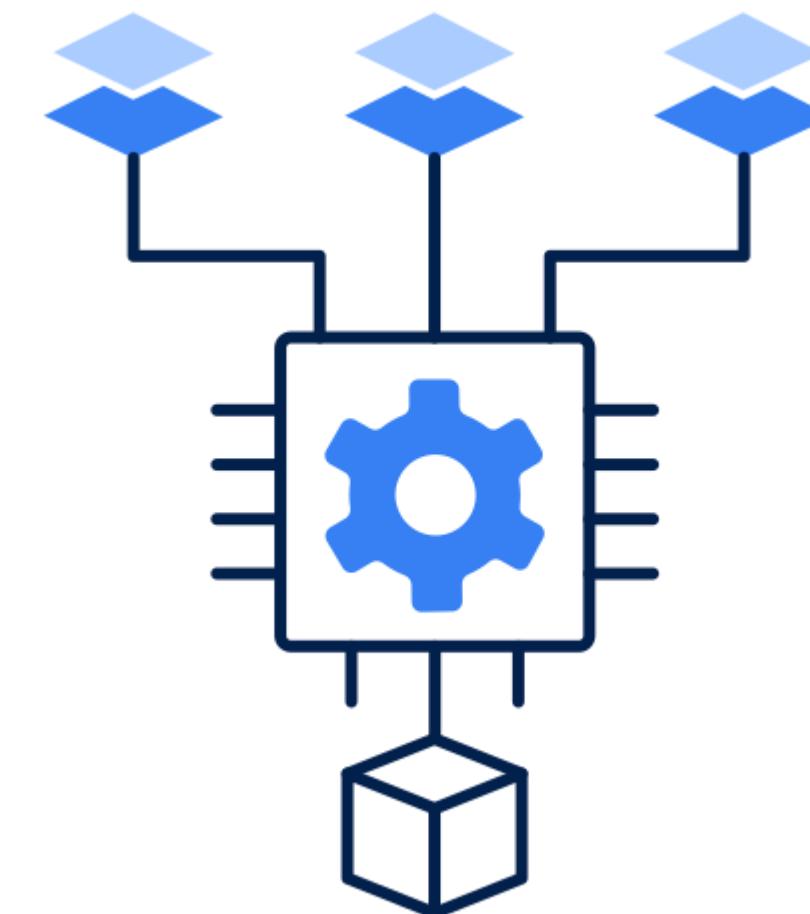
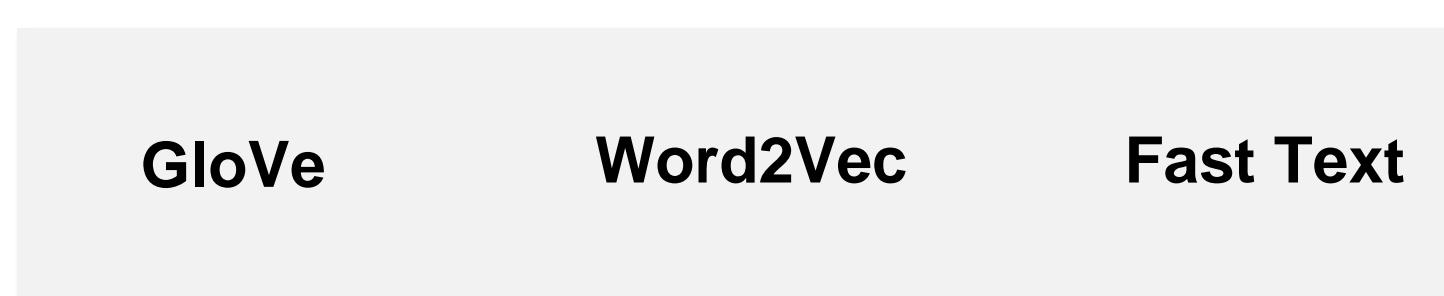


# Co-occurrence Matrix Representation

	first	is	second	sentence	the	this
first	0.0	1.0	0.0	1	1.0	1.0
is	1.0	0.0	1.0	3	2.0	2.0
second	0.0	1.0	0.0	2	1.0	1.0
sentence	1.0	3.0	2.0	2	3.0	3.0
the	1.0	2.0	1.0	3	0.0	2.0
this	1.0	2.0	1.0	3	2.0	0.0

# What are Pre-Trained Embeddings

Word embeddings trained on large corpora to capture semantic relationships and linguistic patterns.



# Overview of FastText Embeddings

Leverages sub word-level information, breaking words into character n-grams.

jumping → <ju>, <jum>, <jump>, <mpin>, <ping>

**fastText**

# Comparing FastText, Word2Vec and GloVe

Feature	fast Text	Word2Vec	GloVe
Handles OOV Words	Yes(via sub words)	No	No
Context Focus	Local + Sub words	Local Context Only	Global Statistics
Morphological Support	Strong(sub word embeddings)	Weak	Weak
Training Complexity	Higher (sub word-level model)	Moderate	Moderate
Pre-trained Models	Available	Available	Available

# Processing Code-mixed (multilingual) Text

# What is Code-mixed (multilingual) Text ?

Code-mixed (multilingual) text refers to language content where **two or more languages** are interwoven within the same sentence, phrase, or conversation.

*“I'm heading to the market to buy some **sabzi** and fruits.”*



# Why It Occurs ?

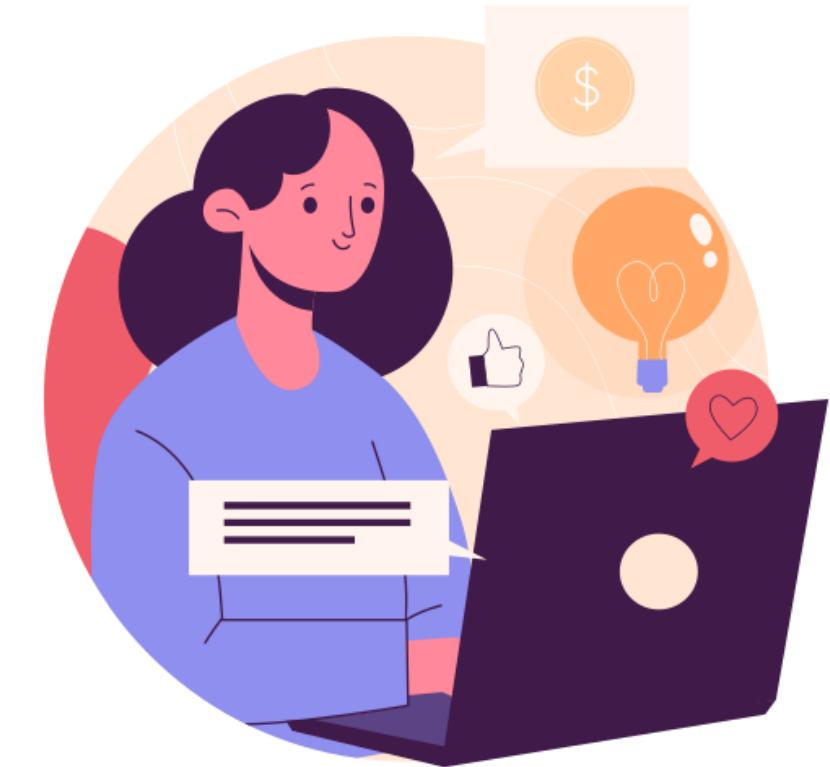


Social Identity

Practical Reasons

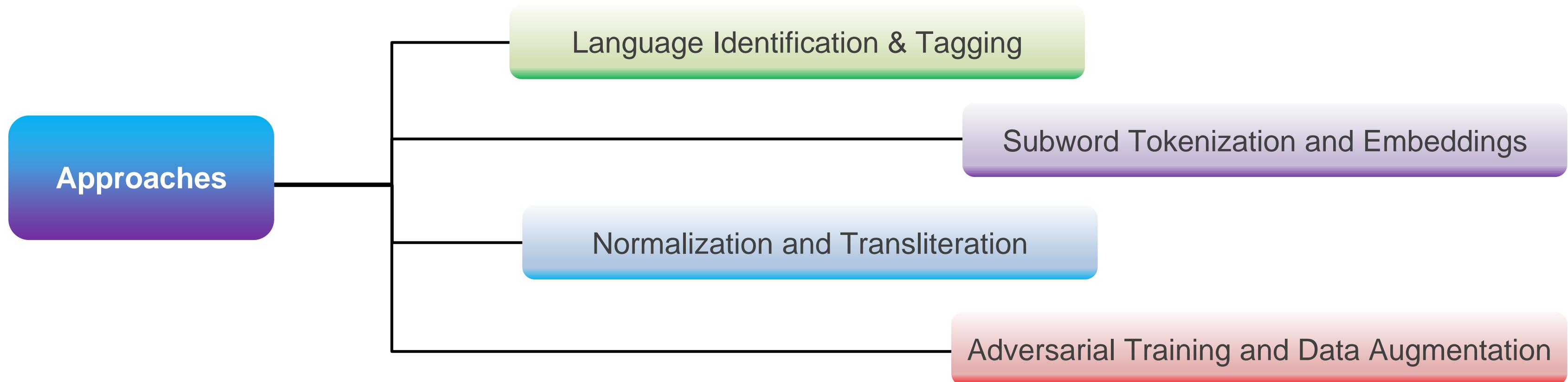


Expressive Nuance



# Tackling Code-mixed (multilingual) Text

To tackle code-mixed (multilingual) text using NLP, here are a few key approaches:



# Language Identification & Tagging

Multilingual pre-trained models (like **mBERT** or **XLM-R**) use **subword tokenization** (WordPiece or BPE) to break words into smaller units. This approach helps in handling:

**OOV words or intra-word mixing:**

**For example**, “chating” might be tokenized into “chat” and “ing”, allowing the model to infer meaning even when languages mix at the subword level.

“chating”

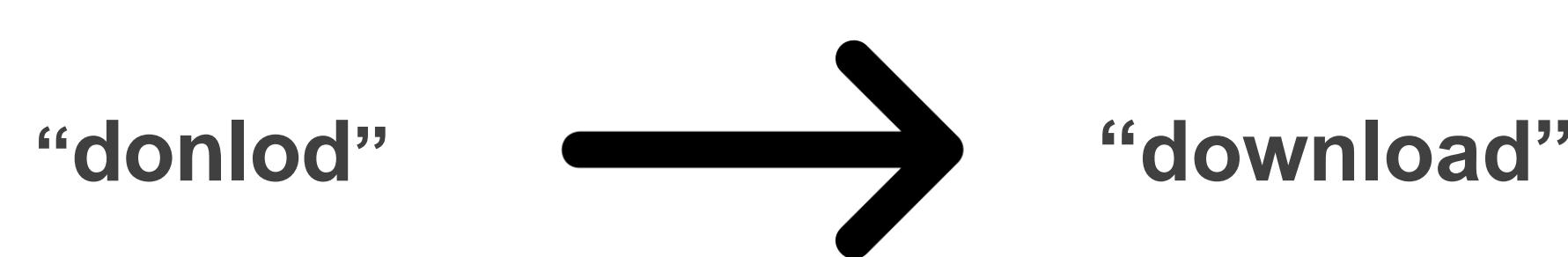


“chat” + “ing”

# Subword Tokenization and Embeddings

Code-mixed text is often informal and noisy. Normalization converts non-standard **spellings** into a **standard form**, and transliteration maps words written in one script to another.

**For instance:** A user might write “donlod” instead of “download” in a mixed sentence. A normalization module can correct this before further processing.



# Adversarial Training and Data Augmentation

When real code-mixed data is scarce, synthetic code-mixed examples can be generated using adversarial attacks.

These examples are used to augment the training data, making the model more robust.

Generating adversarial examples:  
Perturb monolingual sentences by substituting words with their counterparts from another language (e.g., replacing “help” with “ayuda”) to simulate real-world code-switching.



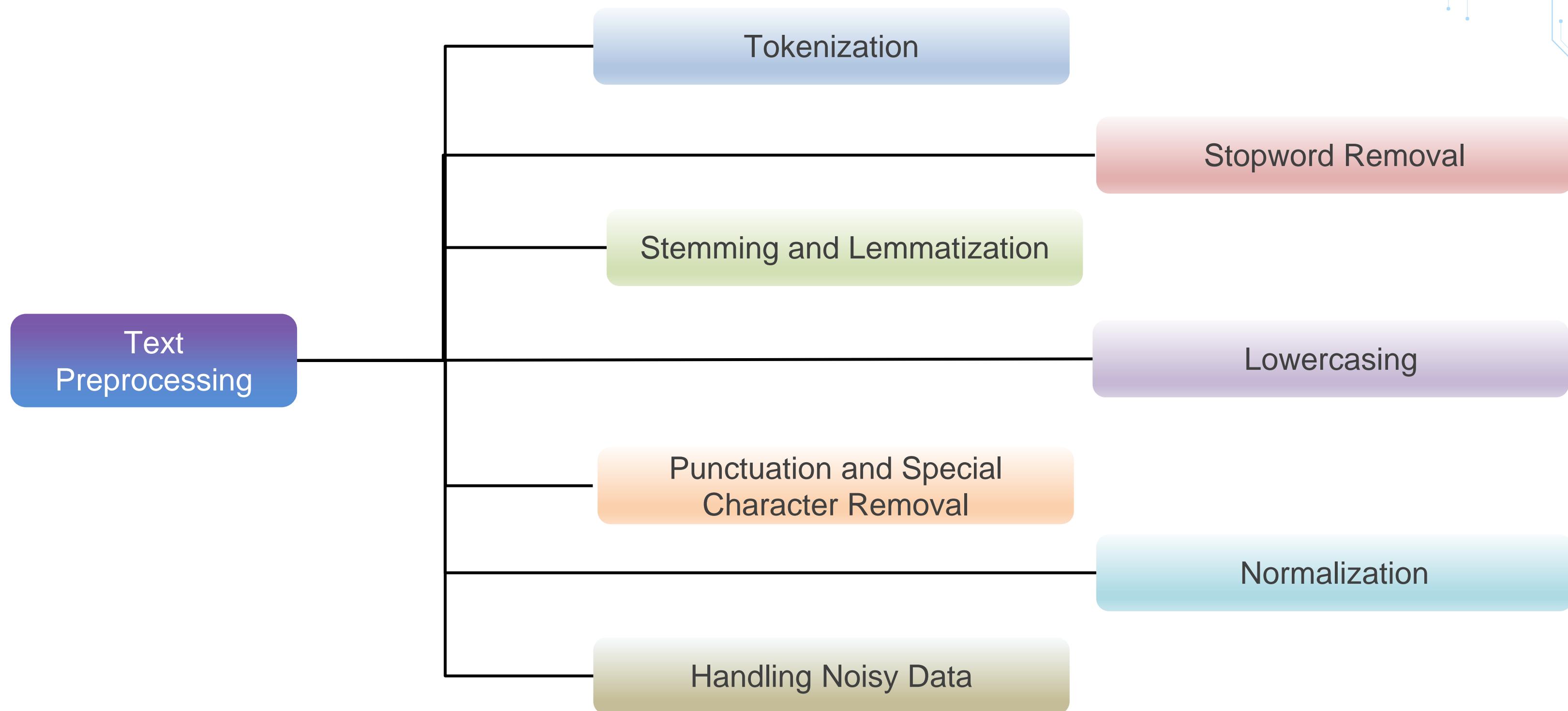
# Text Classification in NLP using Common ML Models

# What is Text Classification in NLP?

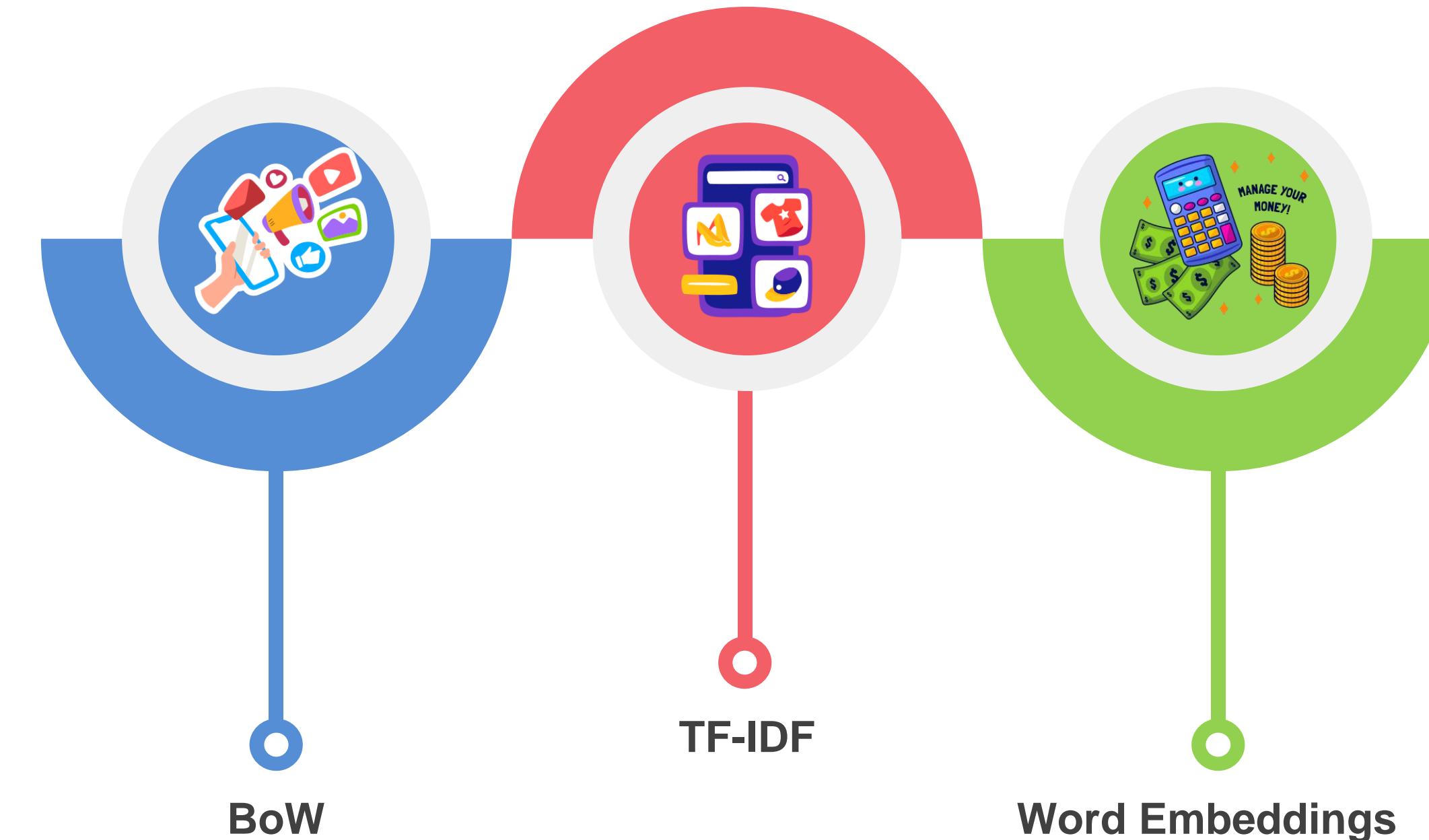


Text classification is a Natural Language Processing (NLP) task that involves automatically assigning predefined categories to text based on its content. It is widely used in spam detection, sentiment analysis, topic categorization, and more.

# Key Steps in Text Preprocessing



# Feature Extraction (Text into Numerical Form)



# Model Training for Text Classification



# Prediction and Model Evaluation

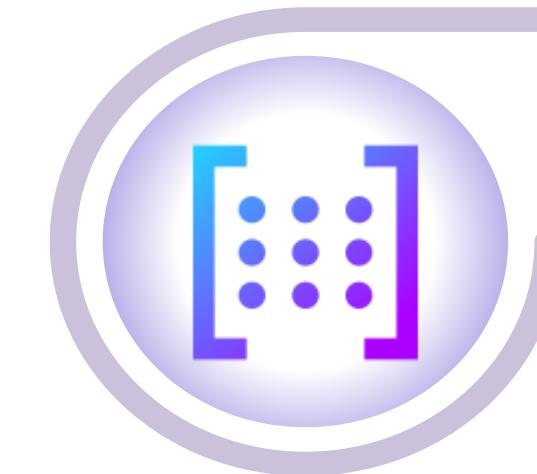
Once the model is trained, it is tested on **new text** to classify it into one of the predefined categories. The performance is evaluated using:



Accuracy



Precision, Recall,  
and F1-score



Confusion Matrix

# Feature Selection for Classification

# What is Feature Selection?

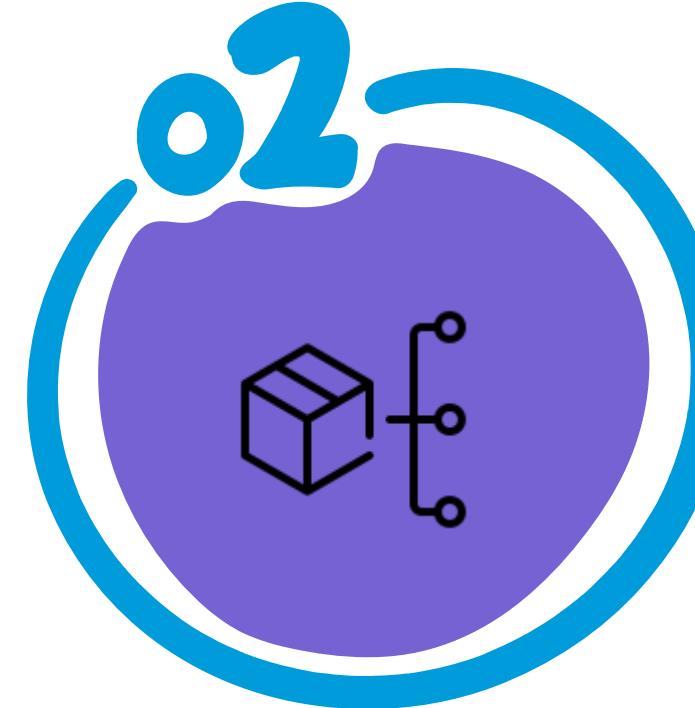
In a dataset, not all features contribute equally to the classification task. Some may be irrelevant or highly correlated, which can negatively impact the model. Feature selection helps identify the best subset of features that maximize classification accuracy.



# Feature Selection- Filter Methods



Correlation Coefficient  
(Pearson/Spearman/  
Kendall)



Chi-Square Test  
(For Categorical Data)



Mutual Information (MI)



Variance Threshold

# Wrapper Methods (Feature Subset Evaluation)

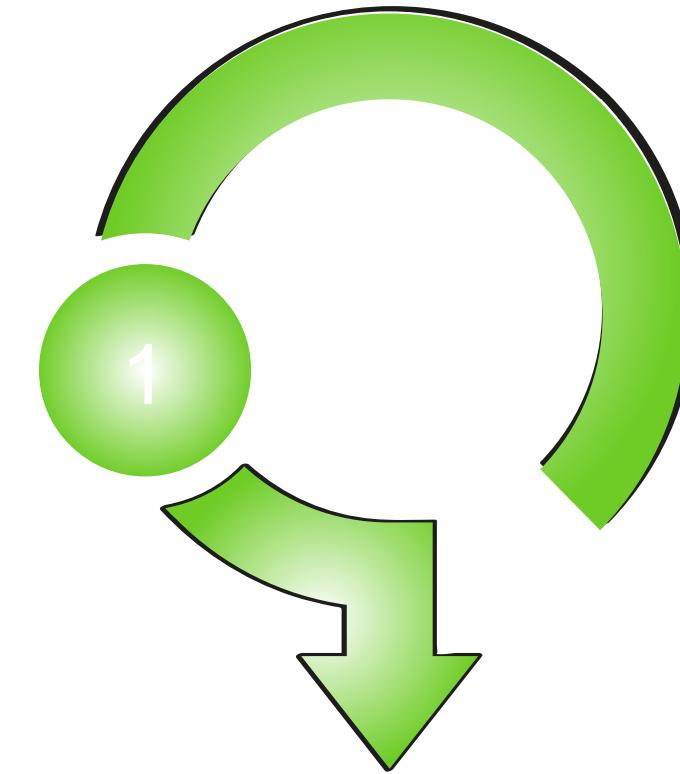
Recursive Feature Elimination (RFE)

Forward Selection

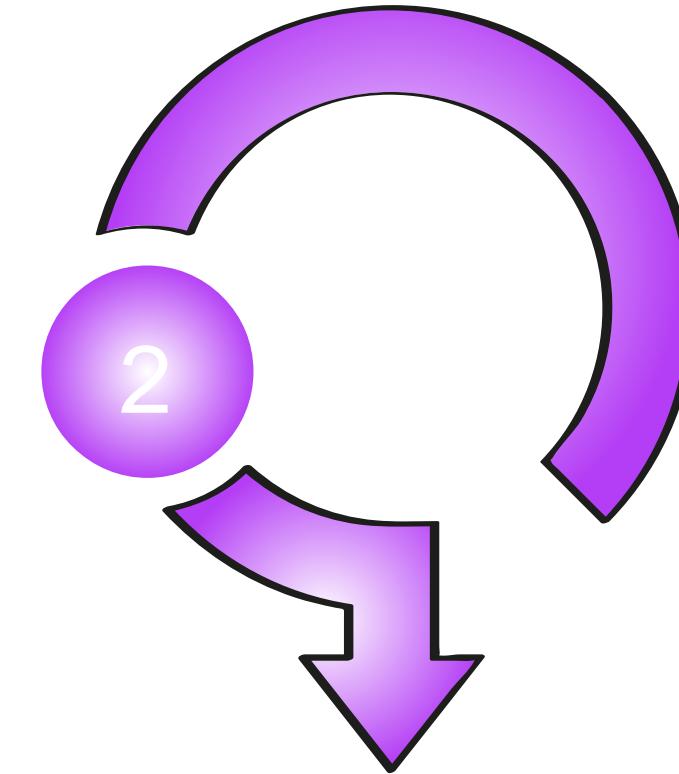
Backward Elimination

Exhaustive Feature Selection

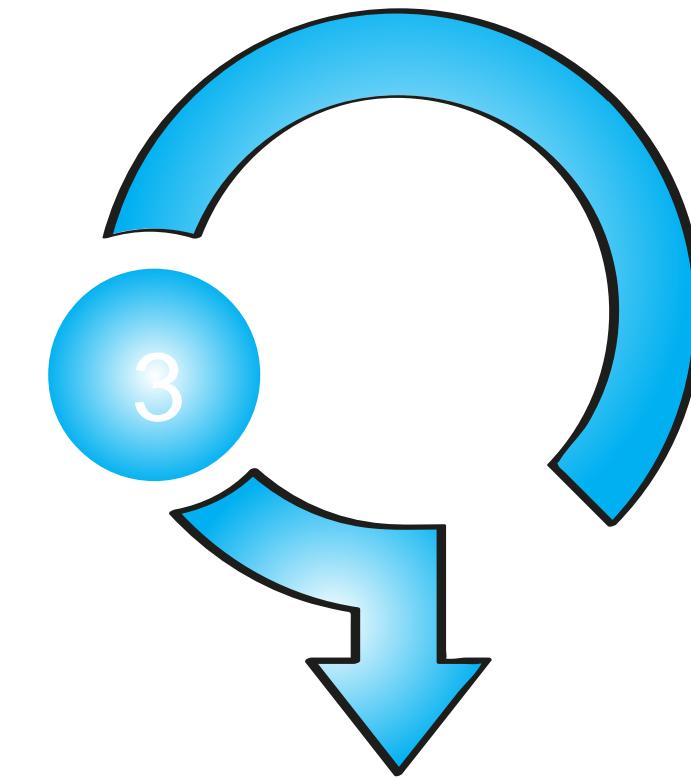
# Embedded Methods(Selection During Training)



Lasso (L1  
Regularization) in  
Logistic Regression



Decision Tree & Random  
Forest Feature  
Importance



Gradient Boosting  
Feature Importance  
(XGBoost, LightGBM,  
CatBoost)

# Feature Selection vs. Feature Extraction

Feature Selection	Feature Extraction
Selects a <b>subset</b> of existing features	Creates <b>new features</b> from existing ones
Removes redundant or irrelevant features	Transforms data into a different representation
Examples: RFE, Mutual Information, Lasso	Examples: PCA, Autoencoders, Word Embeddings
Keeps original feature meanings	Can lose interpretability (e.g., PCA reduces dimensions)

# BoW, TF-IDF and Word Embeddings (Demonstration)

**Note:** Refer to Module 2: Demo 1 on LMS for detailed steps.

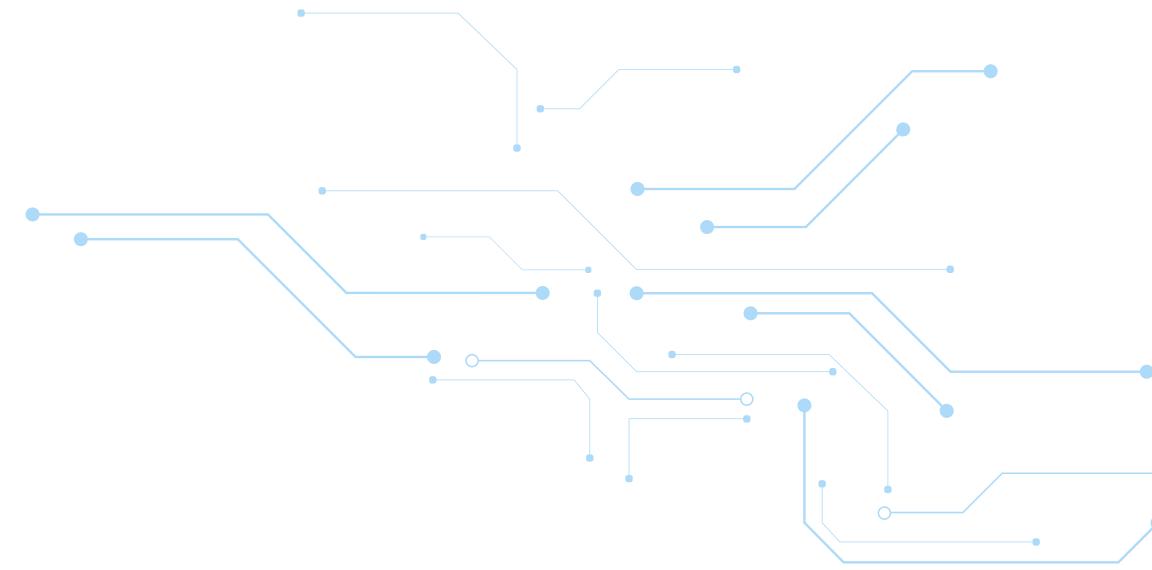
# Evaluating Stemming and Lemmatization Methods for Text Preprocessing (Demonstration)

**Note:** Refer to Module 2: Demo 2 on LMS for detailed steps.

# Summary

In this lesson, you have learned to:

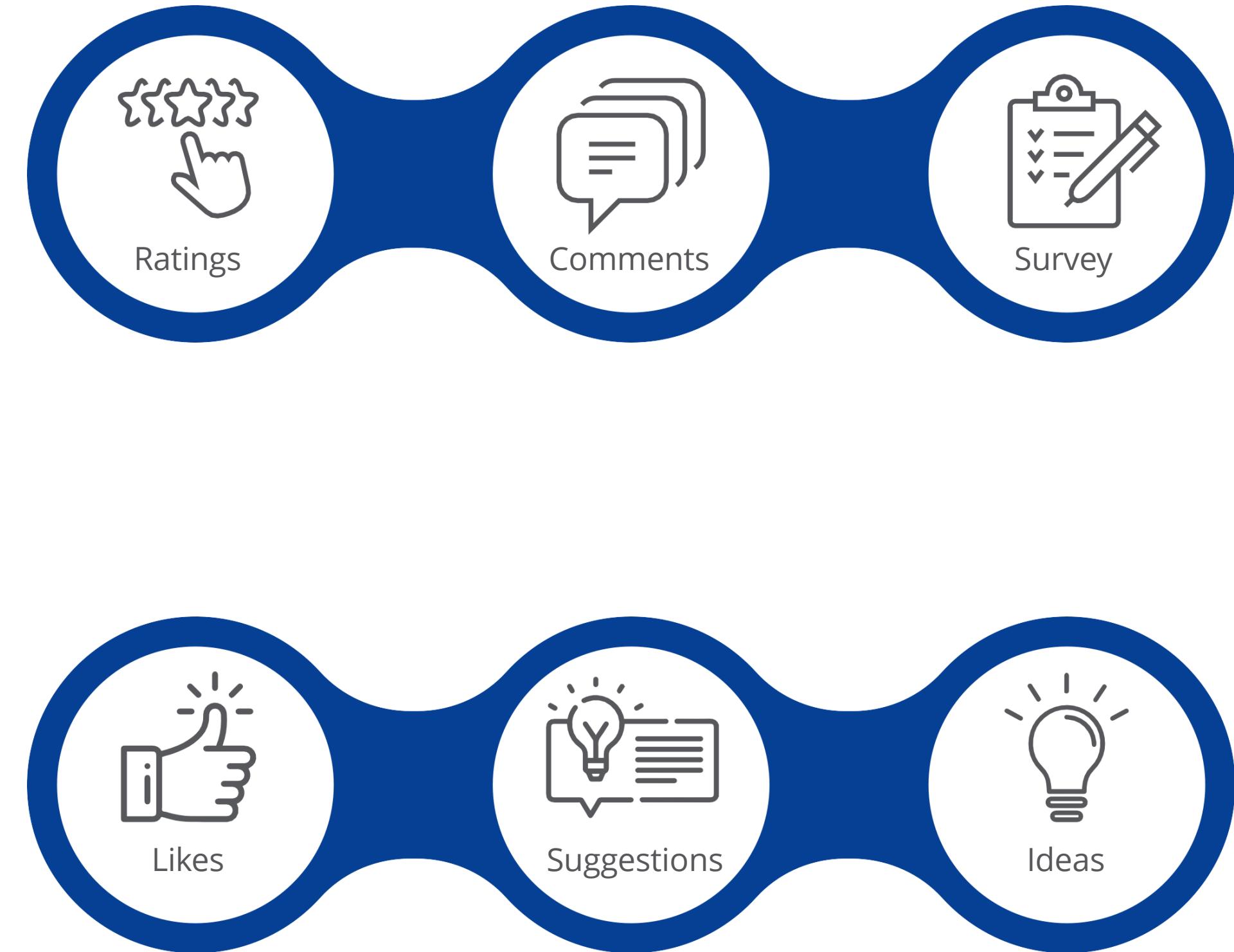
- e! Explain key techniques like stemming, lemmatization, and POS tagging used in NLP preprocessing.
- e! Illustrate how text can be represented using methods such as Bag of Words.
- e! Describe different word embedding models including Word2Vec, GloVe, and FastText for language representation.
- e! Implement text classification strategies and feature selection techniques for multilingual and code-mixed data.



# Questions



# Feedback



# Thank You

For information, Please Visit our Website  
[www.edureka.co](http://www.edureka.co)

