

E) Optical character Recognition (OCR)

→ Reading printed Text.

- * Converts scanned text into editable text.

Input: → A scanned book page.

Output: → Digital text format.

Used in: → Digitizing old books, invoice scanning.

→

Tokenizer

A Tokenizer is a tool that breaks text into smaller parts called tokens.

These tokens can be words, subwords (or) characters.

I Love NLP.

Ex: "I" "Love" "NLP"

Types

1. Word Tokenization → splits text into words

Ex :- AI is powerful

"AI" "is" "powerful".

2) Subword Tokenization :- Breaks words into smaller meaningful parts.

Ex :- "Unhappiness" \rightarrow ["un", "happiness"]

Used in :- GPT, BERT, and Transformers models.

3) Character Tokenization :- splits text character by character.

"Hello" \rightarrow ["h", "e", "l", "o", ...]

useful for handling unknown words -

4) Sentence Tokenization :- splits text into sentences

Ex :- "Hello world. How are you ?

"Hello world"., "How are you".

Techniques TOKEnization.

1. Rule-Based TOKEnization.

- * uses grammar rules (like spaces, punctuation) to split text.

Ex :- splitting words by spaces in English.

- * Doesn't work well in Chinese or Japanese.

2) Regular Expression (Regex) TOKEnization

- * uses patterns to find words.

Ex :- `\w+` → Find all words in a sentence.

3) Byte pair Encoding (BPE) TOKEnization.

- * used in GPT and BERT models.

* merges common letters pairs to form subwords.

Ex: "lowest" ["low", "est"]

4) unigram Tokenization

* Breaks words into multiple subword options and picks the best one.

Ex: "Internationalization" → ["inter", "nation al", "ization"]

challenges in Tokenization

1. Handling multiple meanings (Ambiguity)

Ex: I saw the man with a tele
scope.

2. Dealing with compound words.

Ex: New York → should it be

["New", "York"] (or) ["New York"]

3. Languages without spaces

* Chinese.

4) Handling special characters

& emojis.

Ex: "python is awesome! 😊"

Detecting and Correcting spelling errors

When we type, we sometimes make spelling mistakes, computers use spell checkers to detect and correct these errors automatically. This process is used in word processors, search engines, and chatbots.

1. Types of spelling Errors

1. Typographical Errors

Mistakes from fast (or) incorrect typing.

Eg: 'hte' → 'the'.

2. Phonetic Errors → When a word is spelled as it sounds.

friend → "friend".

3. Omission Errors : Leaving out letters.

"receive" — "receive"

4. Insertion Errors : Adding extra letters.

Ex: "happenn" → "happen"

5. Transposition Errors : swapping two letters.

Ex: "ten" → "the"

6. Homophone Errors : words that sound the same but are spelled differently.

"there" vs. "their"

2) Techniques for Detecting spelling Errors

- * Dictionary - Based Approach.
 - + compares words to a predefined word list.
 - * If a word isn't found, it's marked as incorrect.

Ex :- If "recive" is not in the dictionary, it is flagged as a mistake.

* Edit Distance.

- * Measures how many changes (insert, delete, replace). are needed to fix a word.

Ex :- "hte" \rightarrow "the" (1 change; swap 'h' to 't')

* N - gram Analysis

- * Breaks words into small

parts (n-grams) and checks

common patterns,

- 4) probabilistic Language models.
- * uses models like n-grams (or)
deep learning to check if a word fits well in a sentence.

- 5) Neural networks & AI Based approaches.
- * uses deep learning models (like GPT, BERT) to predict the correct word.

"I wnet to teh market"
→ I went to the market.

Techniques for correcting spelling
errors.

1. Rule-Based correction.
- * uses grammar rules to fix mistakes.

Ex: "I before e except after c"

(currents "recieve" → "receive")

2) Edit Distance correction.

- * Suggests words that need the fewest changes.

Ex : "hte" → "the" (1 change).

3) Context - Based correction.

Looks at the entire sentence.
to fix mistakes.

Ex : "Their going to the park"
"They're going to the park".

4) Auto - complete & Auto - correct.

- * Suggests and fixes words while typing.

"definately" → "definitely".

Minimum Edit Distance

minimum edit distance is a way to measure how different two words (or text) are counting the smallest number of changes needed to turn one word into another. They include.

1. Insertion \rightarrow Adding a letter.

Ex : "aple" \rightarrow "apple" (insert "p").

2. Deletion \rightarrow Removing a letter

Ex : "hello" \rightarrow "heho" (delete "l").

3. Substitution \rightarrow changing one letter to another.

Ex : "cat" \rightarrow "bat" (change "c" to "b")

Advantages of Minimum Edit Distance:

* Simplicity :- Easy to understand and implement.

* Versatility :- Applicable across multiple NLP domains.

optimality :- Guarantees the minimum number of edits.

challenges with minimum edit Distance

* computational complexity :-

$O(m \times n)$ where m and n are the lengths of the strings.

* context Ignorance :- Does not consider word Semantics.

* handling Large Datasets :- Inefficient for real-time spell correction in large corpora.

Applications of Minimum Edit Distance in
NLP

* spell checking :- Suggesting corrections based on edit distance.

2) Machine Translation :- Aligning phrases between languages.

3) Plagiarism Detection :- Comparing text similarity.

4) Speech Recognition :- Aligning predicted and transcribed words.

5) Bioinformatics :- Comparing DNA sequences.

Regular Expressions (Regex).

- * A Regular Expression (Regex) is like a smart search tool. that helps find and match patterns in text..
- * It is used for searching, replacing and validating text in programming and data processing .

Basic Regex patterns.

Patterns	Meaning	Example.
.	any character.	"c.t" → Cat, cot, cat.
^	start of a line	^Hello → matches "Hello" only if it's at the start.
\$	End of a line.	end\$ → matches "end" only at the end..
\d	Any digit(0-9).	" I have 3 apples" → Matches "3"
\w	Any letter, digit or underscore.	" hello_123" → Matches everything.
\s	Any space, tab, or newline.	" hello world " → Matches the space.
*	Repeats 0 or more times	"god*" → "qd", "god", "good", "goood*?"

+	Repeat 1 (or) more times	"god" → "good", "good", "goood" (not "gd").
?	optional (0 or 1 time)	"colour?" → Matches "colour" and "Colour".
{n}	Repeat exactly n times	"a{3}" → Matches "aaa".
{n,m}	Repeat between n and m times	"a{2,4}" → Matching "aa", "aaa", "aaaa".

Python Regex

python has a built-in module called `re` that lets you use Regular Expressions (Regex) to search, Match, and Manipulate text easily.

2) Basic Regex Functions in Python.

re. search ()

* Finds the first match in a string.

re. search ("cat", "The cat sleeps").

re.findall()

* Finds all matches in a string.

re.findall ("dt+", "price : 100, Tax: 20")

re.match()

* Checks if a string starts with a pattern.

re.match("Hello", "Hello.world")

re.sub()

* Replaces a pattern with new text.

re.sub("blue", "red", "blue sky")

re.split()

* Splits a string using a pattern.

re.split ("st", "Hello world")

UNIT - 2.

① Understanding N-grams.

N-grams are essential for understanding language patterns and structure. They help in predicting the next item in a sequence.

i.e.

② N-grams comparison.

1. unsmoothed N-Grams (No smoothing),

- * uses raw counts from the training data.
- * If an N-gram wasn't seen before, it gets a probability of zero, meaning the model fails.

Ex :- (Bigram model).

Training data :- "I love NLP"

If we ask for "I love AI" - Fails.

(Since "love AI" was never seen).

② smoothed N-Grams (with smoothing).

- * Adjusts probabilities so unseen word combinations don't get zero probability.
- * Helps make better predictions.

Ex :- Laplace Smoothing :- Adds a small value

(like +1) to all word counts.

- * Backoff & Interpolation :- uses smaller n-grams (bigram - unigram) when needed.

N-Grams Models.

An N-Gram model is a type of language model that predicts the next word based on the previous words using N-grams.

1. Unigram Model ($N=1$)

~~~~~

- \* predicts the next word based only on individual word frequency.

Ex : "I love" → predicts next word only based on the most frequent word in the dataset.

## ② Bigram model ( $N=2$ ).

- \* predict the next word based on the previous word.

Ex : If "I love" is common, the next word might be "cybersecurity".

"love cybersecurity" appears frequently.

## ③ Trigram model ( $N=3$ ).

- \* predict the next word using two previous words.

Ex : "I love cybersecurity" → The model predicts the next word based on.

"love cybersecurity".

## \* Train and Test corpora.

In when building a machine learning (or) NLP model (like an N-gram model), we need data to train and evaluate it. This data is split into two main parts.

1. Training corpus (Train Data)
2. Testing corpus (Test Data).

### 1. Training Corpus (Train Data).

- \* This is the main dataset used to teach the model.
- \* The model learns patterns from this data.

Ex :- If you're training an autocomplete system, the training data would be books, articles, or chats.

Think of it as :- studying for an exam using text books.

## 2) Testing Corpus (Test Data)

- \* This is a separate dataset used to check how well the model learned.
- \* The model does not see this data during training.

Ex :- If you trained an N-gram model (or) news articles, the test set might be new articles from a different website

Think of it as :- Taking an exam to see how well you studied.

## Unknown words.

When we test an N-gram model (or) any NLP model, sometimes the test data contains words that were never seen in the training data. These called unknown words (out of vocabulary words),

## Handle unknown words?

### 1. special TOKEN ("")

- \* During training, replace rare words with a placeholder like `<unk>`.
- \* When an unknown word appears in the test set, the model treats it as `<unk>`.

### 2. smoothing Techniques.

- \* Techniques like Laplace Smoothing (or) Backoff help assign small probabilities to unknown words.

## Perplexity

Perplexity is a way to measure how good (or) bad a language model (like an N-gram model) is at predicting text.

- + Low perplexity : The model is confident and makes good predictions.
- + High perplexity : The model is confused and makes bad predictions.

### smoothing

☞ smoothing is a technique used in N-gram models to handle zero probability problems.

### Techniques

1. Laplace smoothing (. Add-one smoothing)
  - \* Adds 1 to every word count to prevent zero probability.
2. Add-K smoothing
  - \* Similar to Laplace, but instead of adding 1, we add a small value  $k$  (like 0.01 or 0.001).

### 3.) Backoff smoothing.

- \* If a higher-order N-gram (like a trigram) is missing, it "back off" to a lower N-gram (like a bigram @ unigram).

### 4) Interpolation Smoothing.

- \* combines different N-grams (trigram, bigram, unigram) instead of just back off.
- \* more advanced and balances short term and long-term dependencies.

### 5) Kneser-Ney smoothing

- \* one of the most powerful smoothing techniques.
- \* Instead of just adjusting probabilities. It also considers how diverse the word usage is.

## POS Tagging

POS (part-of-speech) Tagging is the process of labeling each word in a sentence with its grammatical category. (noun, verb, adjective etc.).

sentence : The cat sleeps on the mat.

The → Determiner (DET)

cat → Noun

sleeps → Verb.

on → preposition.

the → Determiner

mat → noun.

## Types of POS Tagging.

\* Rule based POS Tagging.

\* Statistical POS Tagging

\* Hybrid POS Tagging.

## ① Rule - Based pos Tagging.

- \* uses predefined grammar rules and lexical dictionaries to assign pos tags.

Ex: If a word ends in -ing, it's likely a verb.

### Types of Rule - Based pos Tagging

#### \* Lexicon - Based Tagging.

- \* uses a dictionary of words with their pos tags.

Ex: run is a verb, "beautiful" is an adjective.

#### \* Context - Based Tagging.

- \* uses rules to analyze context and decide the correct tag.