

Talent Bootcamp Series -1

Day-3

NLP Lab

Lab created by: Hasan Shahid Kamal



Made with GAMMA

Core NLP Terminology

Corpus

A large, structured collection of text documents used for linguistic analysis and model training.

Token

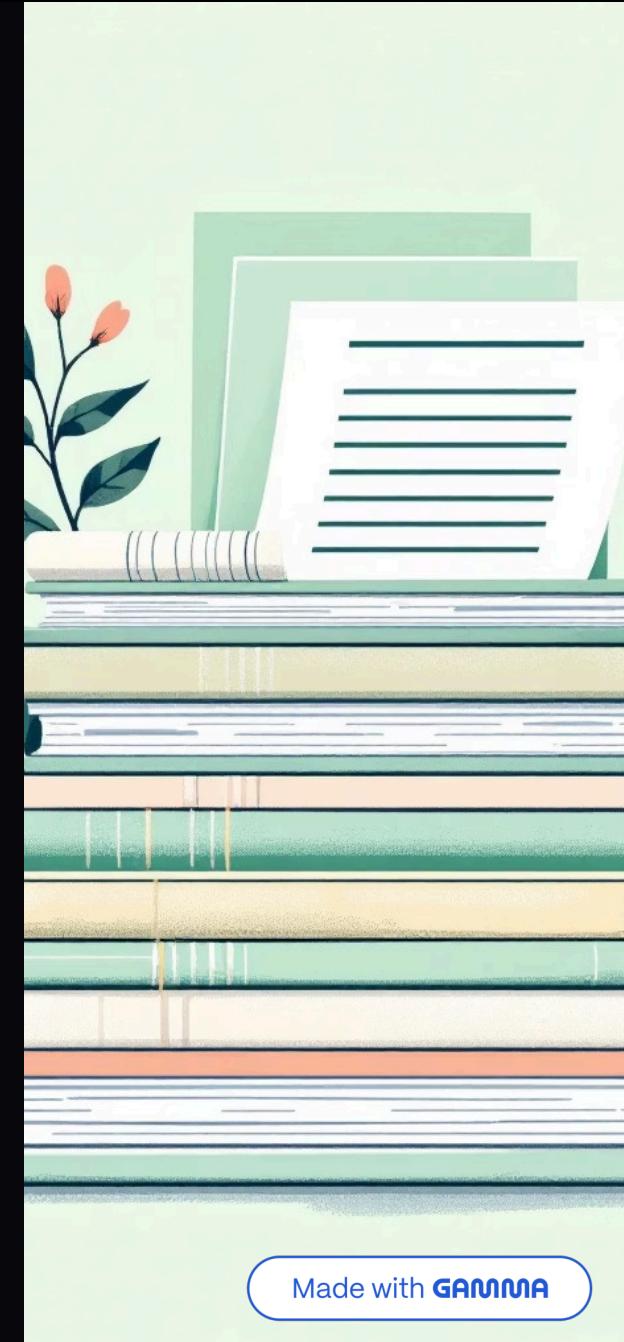
The smallest meaningful unit of text, typically a word or sentence, created through tokenization.

Vocabulary

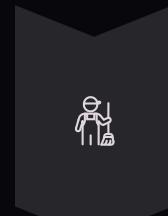
The complete set of unique tokens present in a corpus, forming the basis for text representation.

Document

A single text sample or instance within a corpus, such as an article, review, or message.



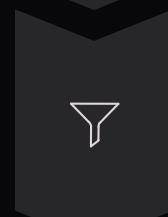
Text Pre-Processing



Clean Raw Text



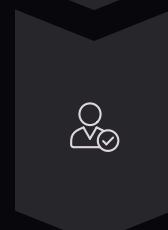
Removes unwanted characters, formatting issues, and inconsistencies from unstructured text sources.



Reduce Noise



Eliminates irrelevant information that could interfere with accurate analysis and model performance.

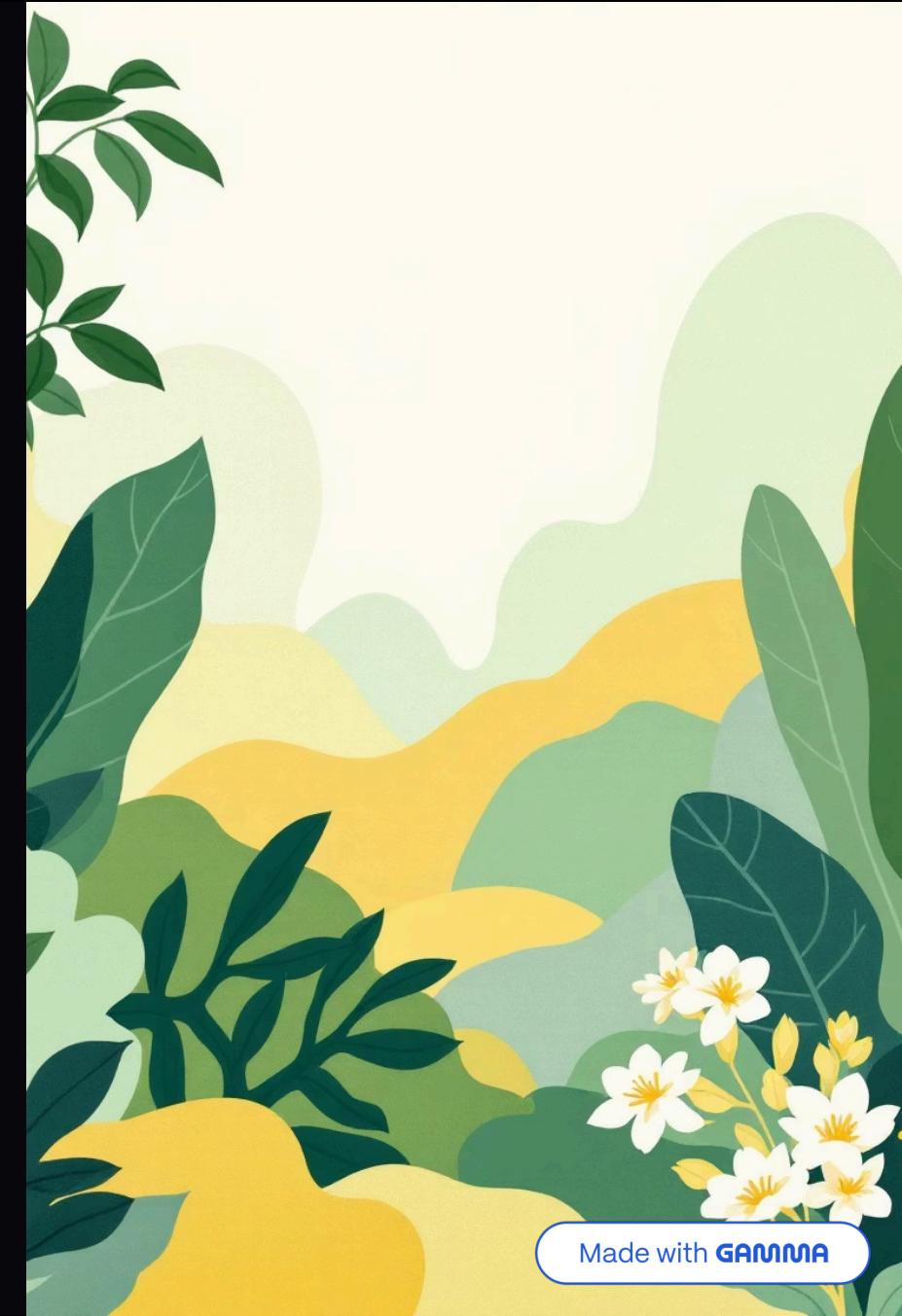


Prepare for Analysis



Transforms text into a standardized format that algorithms can effectively process and understand.

PRE-PROCESSING IS AN ESSENTIAL PREREQUISITE STEP THAT DIRECTLY IMPACTS THE QUALITY AND ACCURACY OF ALL DOWNSTREAM NLP TASKS.



Tokenization

Breaking Text into Meaningful Units

Tokenization is the foundational process of splitting continuous text into discrete, analyzable units. This critical step transforms unstructured text into structured data that machines can process.

Word Tokenization

Splits text into individual words: "I love NLP" becomes ["I", "love", "NLP"]

Sentence Tokenization

Divides text into complete sentences for document-level analysis and context preservation.

['Natural Language Processing is interesting.', 'It helps computers understand human language.', 'NLTK makes NLP easy.']

What is punkt in NLTK?

punkt is a pre-trained tokenization model in NLTK that is used to split text into sentences and words.

It is not a function, but a data package (model) that NLTK needs to perform tokenization correctly.

Why is punkt Needed?

Computers cannot automatically understand:

- Where sentences end
- Whether a dot (.) ends a sentence or is part of an abbreviation

punkt solves this by using language-aware rules and statistics.

Stop Words Removal

What Are Stop Words?

Common words that appear frequently but carry minimal semantic meaning for most NLP tasks. Removing them reduces dimensionality and computational cost.

Common examples include:
the, is, at, and, to, in, of, for, with, on

Improved Efficiency

Reduces vocabulary size and speeds up processing without losing critical information.

Better Focus

Helps models concentrate on content-bearing words that drive meaning and classification.



Stemming

Root-Based Approach

Reduces words to their base or root form by removing suffixes using rule-based algorithms.

Example: **running, runs, ran** all become run

Fast Processing

Computationally efficient and quick, making it ideal for large-scale text processing tasks.

Imperfect Results

May produce non-dictionary words (e.g., "running" → "run", "studies" → "studi") due to crude rules.

1 PorterStemmer

What it is

- A **class in NLTK**
- Implements the **stemming process**

What it does

- Cuts words using **rule-based suffix stripping**
- Does **NOT** check a dictionary

Lemmatization

1

Dictionary Lookup

Uses comprehensive vocabulary databases and morphological analysis to find proper base forms.

2

Grammar-Aware

Considers part of speech and grammatical context to determine the correct lemma.

3

Higher Accuracy

Produces valid dictionary words, ensuring better semantic preservation than stemming.

- ☐ While slower than stemming, lemmatization produces linguistically correct results: "better" → "good", "am/are/is" → "be"

2 WordNetLemmatizer

What it is

- A class in NLTK
- Implements the lemmatization process

What it does

- Uses WordNet (a dictionary)
- Considers word meaning and grammar

Part-of-Speech Tagging

Understanding Grammatical Roles

Part-of-speech (POS) tagging assigns grammatical categories to each word, revealing the syntactic structure and relationships within sentences. This enables deeper linguistic analysis and meaning extraction.



Noun (NN)

Person, place, thing, or idea



Verb (VB)

Action or state of being



Adjective (JJ)

Describes or modifies nouns



Adverb (RB)

Modifies verbs or adjectives

Application: POS tagging helps with named entity recognition, information extraction, and understanding sentence structure for machine translation and question answering systems.

Sentiment Analysis

Sentiment Analysis is an NLP technique used to **identify the emotional tone or opinion** expressed in a piece of text.

- Could be Positive, Negative, Neutral

How NLTK Performs Sentiment Analysis

NLTK uses **VADER (Valence Aware Dictionary for sEntiment Reasoning)**

- Rule-based model
- Designed for short texts
- Works well with informal language



VADER combines three main ideas:

1 Sentiment Lexicon (Dictionary-Based)

VADER has a **pre-built dictionary** where words are assigned sentiment scores.

- ✓ Positive words → positive score
- ✓ Negative words → negative score

2 Rule-Based Heuristics

VADER does not just count words — it applies **linguistic rules**.

a) Capitalization

- ALL CAPS increases intensity

Example:

- “*This is GOOD*” → more positive than “*good*”

3 Sentiment Aggregation

After analyzing all words and rules, VADER:

- Combines scores
- Normalizes them
- Produces **four final scores**

Score	Meaning
pos	Strength of positive sentiment
neg	Strength of negative sentiment
neu	Neutral proportion
compound	Overall sentiment (-1 to +1)

Example for Sentiment Analysis

```
from nltk.sentiment import SentimentIntensityAnalyzer  
sia = SentimentIntensityAnalyzer()  
text = "I really enjoyed this lab session!"  
scores = sia.polarity_scores(text)  
print(scores)
```

Output:

```
{  
'neg': 0.0,  
'neu': 0.3,  
'pos': 0.7,  
'compound': 0.7269  
}
```

Compound Score Interpretation

- $> 0.05 \rightarrow$ Positive 😊
- $< -0.05 \rightarrow$ Negative 😞
- Between -0.05 and $0.05 \rightarrow$ Neutral 😐

Named Entity Recognition (NER)

Identifies real-world entities

- Examples: Person, Organization, Location
- Used in information extraction



Dependency Parsing

Analyzes grammatical relationships

- Shows how words are connected
- Useful for syntactic analysis

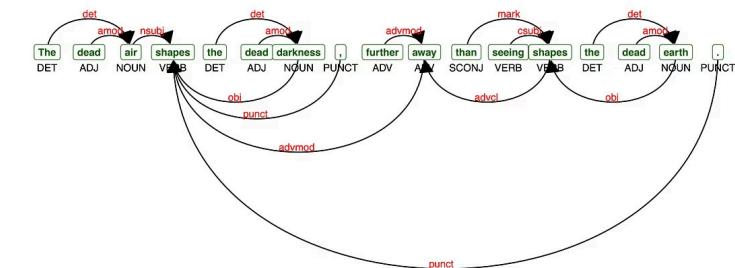
👉 It answers:

- Who did what?
- Which word depends on which?

Each word (token) is connected to another word (its **head**).

Dependency Parser

tokenisation, parts of speech tagging & dependency relations



Sentence:

“**Students** study NLP.”

Dependency Relationships:

- **study** → main verb (ROOT)
- **Students** → subject of *study*
- **NLP** → object of *study*

Students → nsubj → study ← dobj ← NLP



Text Similarity

Measures semantic similarity between texts

- Used in plagiarism detection and search
- Based on vector similarity

The logo features the word "spaCy" in a large, bold, blue sans-serif font. The letters are slightly overlapping and have a slight shadow effect. The background consists of several overlapping, curved bands of various colors, including orange, yellow, red, green, and purple, creating a dynamic, abstract pattern.

spaCy Library

- Industrial-strength NLP library
- Fast and efficient processing
- NER, dependency parsing, similarity

```
12     inYt cativatio lettiox,(0/c andw);
13
14     cinner/o/tener = (Manl);
15     ionar.dations = chargiation (sto tenxrox(iavemyin));
16     (iureatiq anaw ane (ustanction no erection untricel));
17     for our fensor(lf;
18         iant tin oysection - Tensor;
19     |));
20     eutspors/tenagw/(tintion));
21     tensor/tanto tiy;
22     endariowl(iectiol); f
```

Text Vectorization

- Converts text into numerical form
- Required for ML models
- BoW, TF-IDF, Word Embeddings

Bag of Words (BoW)

- Uses word frequency
- Ignores grammar and order
- Implemented using CountVectorizer

Bag of words (BoW)

Very good drama although it appeared to have a few blank areas leaving the viewers to fill in the action for themselves. I can imagine life being this way for someone who can neither read nor write. This film simply smacked of the real world: the wife who is suddenly the sole supporter, the live-in relatives and their quarrels, the troubled child who gets knocked up and then, typically, drops out of school, a jackass husband who takes the nest egg and buys beer with it. 2 thumbs up... very very very good movie.



('the', 8),
(',', 5),
('very', 4),
(',', 4),
('who', 4),
('and', 3),
('good', 2),
('it', 2),
('to', 2),
('a', 2),
('for', 2),
('can', 2),
('this', 2),
('of', 2),
('drama', 1),
('although', 1),
('appeared', 1),
('have', 1),
('few', 1),
('blank', 1)
....

TF-IDF

- Measures word importance
- Reduces impact of common words
- Implemented using TfidfVectorizer

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

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

$$\log \frac{1 + n}{1 + df(d, t)} + 1$$

Term frequency
Number of times term t appears in a doc, d

Inverse document frequency
 n ← # of documents

Document frequency of the term t

Text Classification

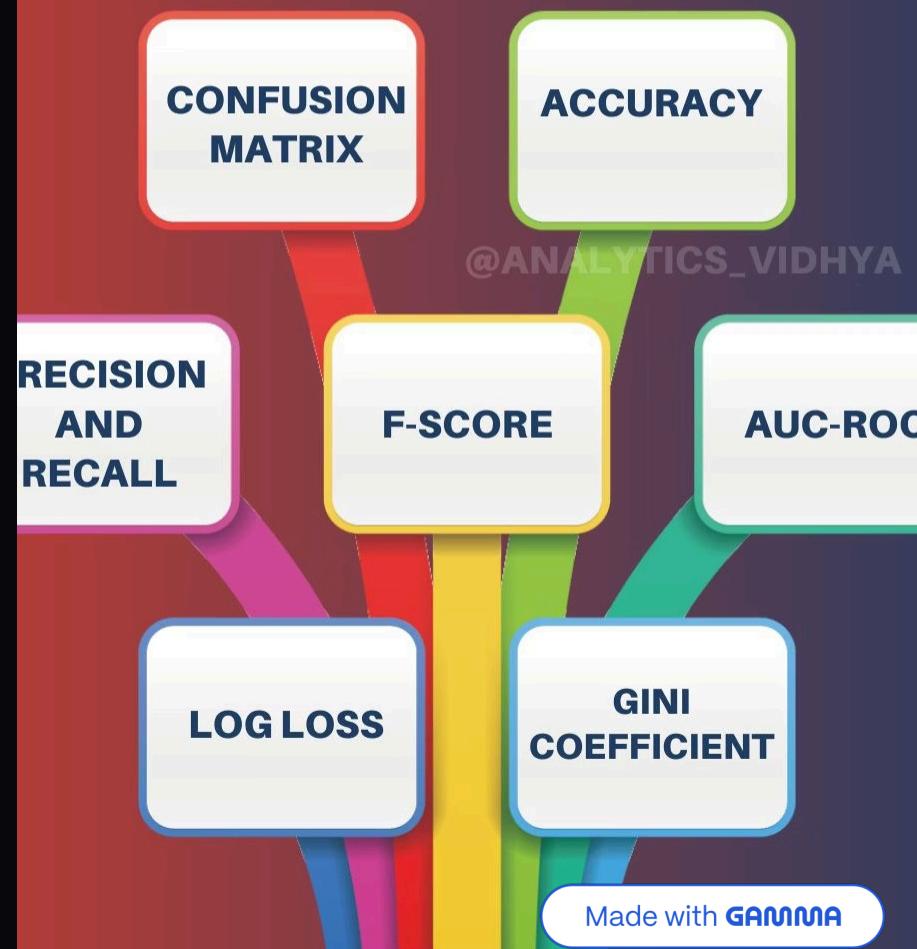
- Automatically categorizes text
- Spam detection, sentiment classification
- Uses ML algorithms



Model Evaluation Metrics

- Accuracy, Precision, Recall, F1-score
- Confusion Matrix for error analysis
- Ensures model reliability

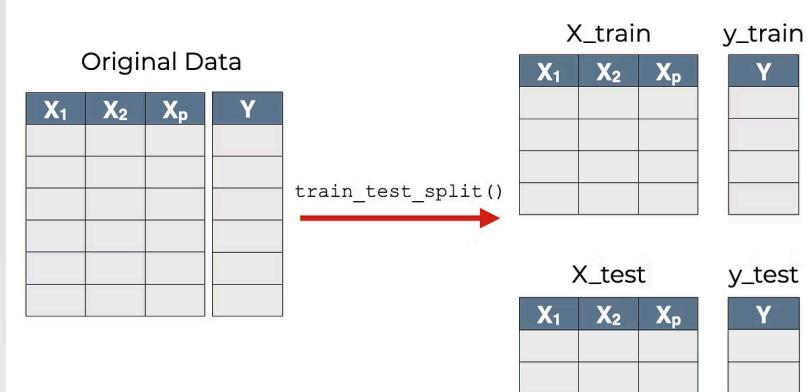
EVALUATION METRICS FOR CLASSIFICATION



Scikit-Learn Tools

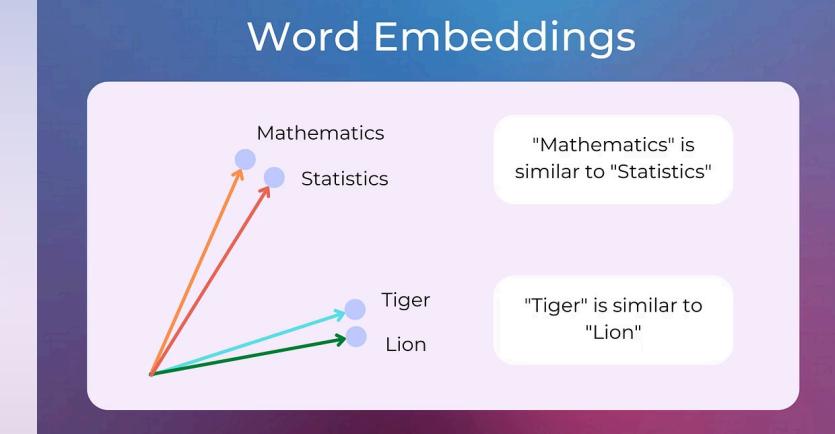
- CountVectorizer, TfidfVectorizer
- train_test_split, MultinomialNB
- classification_report, confusion_matrix

TRAIN_TEST_SPLIT SPLITS DATA INTO TRAINING DATA AND TEST DATA



Word Embeddings

- Dense vector representation of words
- Captures semantic meaning
- Similar words have similar vectors



Word2Vec

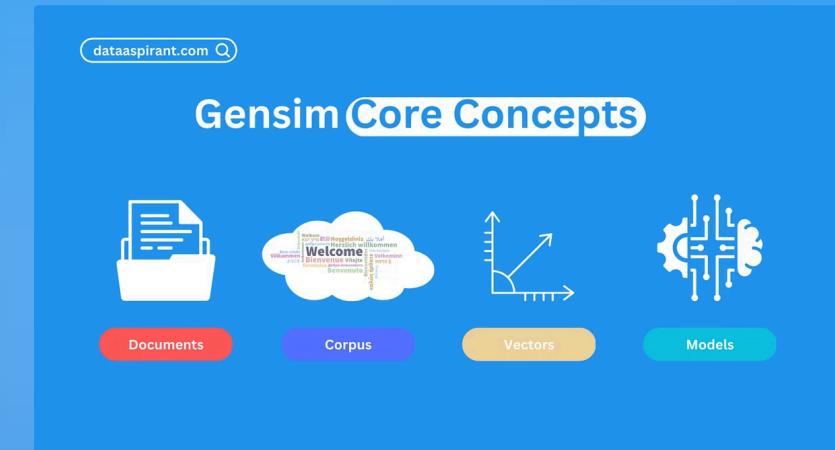
Skip-gram

Word2Vec

- Neural network-based embeddings
- CBOW and Skip-gram models
- Implemented using Gensim

Gensim Functions

- Word2Vec()
- most_similar()
- similarity()



NLTK vs spaCy vs Gensim

Feature	NLTK	spaCy	Gensim
Main Role	Pre-processing & learning	Industrial NLP	Semantic modeling
Tokenization	✓	✓	✗
POS Tagging	✓	✓	✗
NER	✗	✓	✗
Word Embeddings	✗	Limited	✓ (Best)
Topic Modeling	✗	✗	✓
Speed	Medium	Very fast	Very fast
ML Focus	Low	Medium	High

Key Takeaways from the Lab

- NLP converts **raw human language** into a form machines can understand
- **Text pre-processing** (tokenization, stop words, stemming, lemmatization) is a **crucial first step**
- **NLTK** is best for learning and basic NLP tasks, while **spaCy** is used for fast, advanced linguistic analysis
- **VADER** enables quick sentiment analysis without training data
- **Vectorization techniques** (BoW, TF-IDF, Word2Vec) allow text to be used in machine learning models
- **Word embeddings** capture semantic meaning better than frequency-based methods
- Proper **feature extraction and evaluation** directly impact model performance