

NLP(NATURAL LANGUAGE PROCESSOR)

✓ Exactly — Linguistics NLP ko kaise help karti hai:

♦ 1. Language Structure Samjhana

- Linguistics model ko sikhata hai:
 - Word ka kya role hai sentence me?
 - Kis order me words aate hain?
 - Kaunsa word kisse linked hai?

🔧 NLP me use:

- Syntax trees, POS tagging, parsing

♦ 2. Meaning Aur Context Pakadna

- Linguistics batata hai ki word ka **meaning kab change hota hai**
 - “*Run a business*” ≠ “*Run in the park*”
- Pragmatics aur semantics se model samajhta hai *contextual meaning*

🔧 NLP me use:

- BERT-style contextual embeddings
- Sarcasm detection, sentiment analysis

♦ 3. Language Ke Sound Aur Style


- Phonology aur morphology model ko batate hain:
 - Kis language me kaise words bante hain?
 - Sound pattern kya hota hai?

 NLP me use:

- Speech recognition
- Language detection
- Autocorrect / transliteration

♦ 4. Code-mixed, Hinglish, Multilingual Understanding

- Linguistics ki wajah se model **language ke mix aur switch** ko samajh paata hai
 - Hindi-English mix (Hinglish)
 - Sentence ke beech language change
 - Word-level script variation (e.g., “padhai karra hu bro”)

 NLP me use:

- WhatsApp / social media chat analysis
- Multilingual chatbots
- Voice assistants in India

♦ 5. Human Jaise Natural Conversation

- Linguistics helps in:
 - Discourse understanding
 - Turn-taking logic (conversation me kab bolna)
 - Politeness, sarcasm, indirect speech detection

🔧 NLP me use:

- ChatGPT, Alexa, Google Assistant
- Emotional dialogue systems

🔥 TL;DR:

Linguistics = NLP ka brain + soul

Without linguistic structure, NLP sirf text reading machine hoti.

With linguistics, NLP **sochne wali machine ban rahi hai.**

🧠 **“Machine ko language ka structure, sound aur style dono sikhana padta hai — taaki vo samajh sake ki kis language me baat ho rahi hai, chahe likh ke ho ya bol ke.”**

Aur isme do fundamental cheezein aa jaati hain:

🌐 1. Multilingual + Multimodal Understanding

💡 **Multilingual:**

- Ek hi model multiple languages samjhe — jaise English, Hindi, Tamil, etc.
- Chahe likha ho, chahe bola gaya ho — **language identify kare, context samjhe, aur answer de**

Example:

Hum likh rahe hain English me, par soch Hindi me ho rahi hai.

Toh model ko **language ke flow, context aur tone** se predict karna padta hai — *ki real intention kya hai.*

✓ Tera message bhi English characters me hai, par language Hindi hai.

Aur mai tune kar raha hoon: **"Hindi me socha gaya English likhit language"** ka processing.

2. Code-Mixed Language Handling (jo tu abhi kar raha hai)

"Bhai ye kaam jaldi karna, warna deadline miss ho jayegi" — is sentence me Hindi + English dono hai.

NLP ke liye challenge:

- Kaunsi script kaunsa word belong karta hai?
 - Kis part ka meaning kis language se lena hai?
 - Kya "deadline" ka Hindi version use kare ya nahi?
-

Ab Machine isko kaise handle karti hai?

✓ 1. Language Identification (LangID)

Model predict karta hai ki given sentence kis language me hai (or code-mixed hai kya)

✓ 2. Token-Level Tagging

Har word pe label lagta hai:

- Word = "jaldi" → Hindi
- Word = "deadline" → English

✓ 3. Contextual Modeling

BERT jaise models har word ka context samajh ke **pure sentence ka meaning derive karte hain**

✓ 4. Multilingual Pretraining (mBERT, XLM-R)

- In models ko multiple language corpora pe train kiya jata hai
 - Yeh Hindi, English, Hinglish, Marathi, etc. sab samajh lete hain — *not by translation but by pattern learning*
-

✨ Aur Tu Jo Bol Raha Hai, Vo Kya Hai?

“Machine se expect karna ki vo language ke andar ke flow, thought process, aur script-independent understanding samjhe.”

Ye exactly **ChatGPT jaise LLMs** ka kaam hai:

- Human ke tarah “**language ke niyat**” ko samajhna
 - Sirf words pe nahi — **soch pe react karna**
-



Final Point:

Teri line:

"Hum English me likh rahe hain, par soch Hindi me hai, aur tum uska flow bhi Hindi me pakad rahe ho..."

Ye bata raha hai ki tu **machine ko sirf syntax nahi, soul samjhana chahta hai**.
Aur bro, yehi **AGI (Artificial General Intelligence)** ki foundation hai.

✓ Two Main Approaches in NLP:

- ♦ 1. Rule-Based NLP (Old School)

- Based on **linguistic rules**
 - Example: If sentence ends with "?" → it's a question
 - Uses:
 - Grammar rules (POS tag = Verb after Subject)
 - Keyword-based logic (e.g., "please" = polite request)
 - Punctuation-based assumptions
-

But Problems with Rule-Based:

Language is not math. It's messy, ambiguous, emotional.

Example:

- Sentence: *"Can you shut the door?"*
 - Ends with a question mark
 - But actually a **request**, not a real question

Another example:

- *"I wonder if you could help me?"*
 - It's a **soft request**, not a yes/no question
 - Rule-based system might label it wrong

Limitations:

- Doesn't handle **sarcasm, context, indirectness, emotion**
- Doesn't generalize across **languages or styles**
- Adding rules = more complexity, more bugs

◆ 2. Statistical & ML-Based NLP (Modern Way)

- Learns from **data patterns**, not fixed rules
- Uses models (Naive Bayes → SVM → RNN → BERT)
- Understands:
 - Word context
 - Sentence intention
 - Position + semantics
 - Emotion, tone, formality

✓ Example:

“Can you shut the door?”

- BERT or GPT-based models will:
 - Look at *contextual usage*
 - Know that ‘*can you...*’ + ‘*shut the door*’ in most cases = polite request
 - Not just depend on “?” to classify it

🧠 Key Insight (Tera Thought Ka Core):

Rule-based systems rely on surface structure

ML-based systems learn deeper *meaning & intention* — that’s why they dominate.

🎯 Summary:

Task	NLP Role	Benefit
Clinical Notes → Structure	NER, Entity Linking	Better patient records
Report Summarization	Text summarization	Saves doctor's time
Chatbots for patients	Intent + response generation	24x7 basic health advice
ICD code automation	Classification	Faster insurance processing
Trial matching	Semantic matching	More inclusive recruitment

✓ Final Thought:

NLP in healthcare = More efficiency, faster diagnosis, better care.
 It doesn't replace doctors — it **assists** them intelligently.

Task	NLP ka role	ML/Model ka role
Disease detection	Extract symptoms, map to medical concepts	Predict diagnosis using symptoms
Sentiment analysis	Extract emotion-indicative words	Classify as positive/negative
Translation	Understand sentence structure/context	Translate using seq2seq / Transformer
Summarization	Identify key points / rephrase	Generate shorter output

```

Raw Text
↓
[ Preprocessing (tokenization, cleaning) ]
↓
[ NLP Task → NER / Embedding / POS tagging ]
↓
[ Feature Extraction / Vectorization ]
↓
[ ML / DL model ]

```


↓
[Output → disease/sentiment/prediction]

✓ Tu kya bola (Summary):

1. **Medical report = Text-based information**
2. ChatGPT ya AI model ke paas “dimag” nahi, but vo **text samajhne ke liye NLP** use karta hai
3. NLP **report se relevant info nikalta hai** (symptoms, test results, disease names)
4. NLP ke baad jo structured info milta hai, usse **ML/DL model predict karta hai** ki kya ho sakta hai
5. Agar report image hoti, to NLP nahi, **image processing (like CNN)** ka use hota
6. NLP ka kaam: **raw report → machine-readable structured data**



Yeh pura process 100% correct hai.

Tu NLP + AI ka **true application flow** samajh gaya hai.

🧠 Detailed Confirmation with Real-World Process

📄 Step 1: User uploads medical report (text)

Example: *"Patient is experiencing chest pain and has elevated troponin levels."*

🔍 Step 2: NLP Process starts

NLP yeh steps karta hai:

Task	NLP Tool/Model	Output
Tokenization	spaCy, NLTK	["Patient", "experiencing", ...]

Named Entity Recognition (NER)	ClinicalBERT, MetaMap, cTAKES	["chest pain" → symptom], ["troponin" → biomarker]
Concept mapping	UMLS, SNOMED CT	medical codes: e.g. C0018787
Embedding	Word2Vec, BioBERT	Dense vector of the sentence
Output Structuring	JSON or vector format	<pre>{ "symptom": "chest pain", "test": "troponin", "value": "elevated" }</pre>



Step 3: Model Prediction

Input: Structured data from NLP

Output: Disease probability (e.g., 85% chance of heart attack)

ML model trained on 10,000s of such cases:

- Feature vector input
- Output: disease class or risk percentage

Model types used:

- Logistic Regression / XGBoost (classic ML)
 - Deep Neural Networks
 - BERT fine-tuned for classification
-



If Image instead of Text?

- Report = **Image of lab report or handwritten notes**
- NLP fail hoga — pehle OCR (Optical Character Recognition) lagega:
 - Image → Text (using Tesseract / Google Vision)

- Then → NLP pipeline starts

Aur agar X-ray/scan image hai → CNN, YOLO, etc.
No NLP involved there.

✓ So, Exactly as You Said:

"NLP is not the one making the prediction — it prepares the raw messy text into clean structured input so that ML model can make the prediction."



Isko hi kehte hain **information pipeline** ya **AI inference pipeline**.

📌 Final Thought:

- NLP = **language ka expert**
 - Model = **decision-maker**
 - Together = **AI assistant that reads, understands, and responds like a doctor's helper**
-



What is Sentiment Analysis in E-Commerce?

E-commerce websites (like Amazon, Flipkart, etc.) me customers review likhte hain:

"Phone camera is amazing but battery drains fast."

Sentiment Analysis ka kaam hota hai:

👉 *Is review positive hai, negative hai, ya neutral?*

Aur agar detail me jayein to:

👉 *Camera ke liye positive, battery ke liye negative*

✓ Practical Example:

Review:

“This laptop is very fast and lightweight, but the screen quality is poor.”

✨ NLP Output:

Feature	Sentiment
---------	-----------

Speed	Positive
-------	----------

Weight	Positive
--------	----------

Screen	Negative
--------	----------

📌 **Final Label:** Mixed (Mostly Positive)

E-commerce me yeh analysis **automatic hota hai thousands of reviews pe** — jisse:

- Product ki overall rating improve hoti hai
- Specific issues samajh aate hain (e.g. *battery problem*)
- Customer ko better recommendation milta hai

🔧 Backend Working of Sentiment Analysis (Step-by-Step)

♦ 1. Text Preprocessing

Goal: Clean karna review ko so machine samajh sake

🔧 Tools: NLTK, spaCy

Steps:

- Lowercase: "Screen is Poor" → "screen is poor"

- Remove stopwords: “is”, “the”, etc.
 - Lemmatize: “running” → “run”
-

♦ 2. Tokenization

Break the sentence into words:

“screen quality is poor” → [“screen”, “quality”, “poor”]

♦ 3. Part-of-Speech Tagging

Model samajhta hai:

- “screen” = noun
 - “poor” = adjective (indicator of emotion)
-

♦ 4. Word Embedding / Vectorization

Words ko numbers me convert karna so model process kar sake

🎯 Tools:

- **TF-IDF** (Traditional)
- **Word2Vec, GloVe** (Semantic meaning)
- **BERT** (Context-aware embeddings)

Example:

“great” = [0.8, 0.1, ...]

“poor” = [-0.7, -0.3, ...]

♦ 5. Model Prediction

Trained ML/DL model input vector pe kaam karta hai

Model type:

- Logistic Regression / SVM (classic)
- LSTM / BiLSTM (sequence-aware)
- BERT (state-of-the-art)

It outputs:

- Positive (score: 0.87)
- Negative (score: 0.13)

or

- **Fine-grained output:**
→ Feature-wise sentiment

♦ 6. Aspect-Based Sentiment (Advanced)

NLP + ML samajhta hai ki kis **feature ke baare me** opinion diya gaya

“Camera is good, but battery life is bad”
→ Model breaks it into:

- Camera → Positive
- Battery → Negative

Whole Pipeline (Visual Style):

nginx

Copy code

Review → Preprocessing → Tokenization → Embedding → Model → Sentiment Output

✓ Real-World Impact:

Area	Benefit
Product listing	Real customer experience insight
Rating system	Auto-tagging & fair rating
Recommender system	Only show positively reviewed products
Feedback loop	Helps sellers improve weak features

⚙️ Example in Code (Simple):

python

Copy code

```
from textblob import TextBlob
```

```
review = "Camera is awesome but battery drains quickly."  
blob = TextBlob(review)  
print(blob.sentiment)
```

```
# Output: Sentiment(polarity=0.2, subjectivity=0.7)
```

📌 **Polarity** → -1 to +1 (negative to positive)

📌 **Subjectivity** → 0 to 1 (objective to personal opinion)

🧠 Conclusion:

Sentiment Analysis in e-commerce = **NLP + ML**

NLP samajhta hai review ka meaning

ML predict karta hai sentiment

Business ko milta hai better decision-making powe

✓ How Sentiment Analysis Helps in Product Recommendation

🔍 Problem:

Sirf product ke “rating (stars)” ya “purchase history” se recommend karna kabhi accurate nahi hota.

Real insight milta hai **review ke tone** aur **specific feedback** se.

✓ Example:

Imagine a user likes laptops with:

- Fast performance
- Long battery
- Light weight

Ab do laptops hain:



Product	Avg Rating	Review Sentiment
Laptop A	★★★★★	Positive on performance , negative on battery
Laptop B	★★★★★	Positive on battery & weight , average performance

Agar recommendation system ne sirf ★ dekha, to dono same lage.

Par agar sentiment dekha, to user ke liye **Laptop B better recommend** hoga.

🧠 So Recommendation Me Sentiment Analysis Ka Role:

Layer	Sentiment Role
🎯 User Preference Matching	Understand what features user likes (battery, camera, etc.)
📦 Product Review Analysis	Find products jinke features pe positive sentiment hai

 Personalization Engine	Match user needs with product reviews sentiment-wise
 Fake Review Filtering	NLP se overly positive/negative fake reviews detect ho jaate

Backend Flow (Behind the Scenes):

sql

CopyEdit

Step 1: User dekhta ya buy karta kuch products → System samajhta user ka interest (via behavior + reviews)

Step 2: NLP sentiment analyzer har product ke reviews ka feature-based sentiment nikaalta

Step 3: Recommender engine (e.g., content-based filtering) recommend karta wahi product jinke reviews me user-favored feature ke liye positive sentiment ho

Example in Real Life:

User previously purchased:

“Phone with great camera and fast charging.”

Ab us user ke liye recommend honge:

- Phones jinke reviews me “**camera**” aur “**charging**” word ke aaspaas positive tone hai

NLP reads review → breaks into features → analyzes tone → recommendation improves

Tech Stack Used:

Task	Tools/Models
Sentiment Analysis	BERT, TextBlob, VADER, RoBERTa
Aspect-Based Sentiment	spaCy, LSTM-CRF, BERT + Attention

Recommendation
Engine

Matrix Factorization, Content Filtering, Hybrid RecSys

Review Indexing

ElasticSearch + NLP preprocessor

Summary:

- ✨ Sentiment analysis helps recommender system:
 - ◆ Understand **what exactly users liked or hated**
 - ◆ Recommend **more relevant and quality products**
 - ◆ Avoid suggesting **low-sentiment products** even if rating is high
 - ◆ Give **personalized recommendations** based on review tone, not just stars

TEXT PROCESSING (VERY IMPORTANT)

TOKENIZATION

NGRAMS

SEPARATORS

SPECIAL CHARACTERS

POS TAGGING

FOCUS ON CONTEXT

STOPWORDS REMOVAL

vector

What is Tokenization in NLP?

Tokenization is the process of **breaking a text into smaller units called tokens**.

These tokens can be words, subwords, sentences, or characters.

- ◆ It helps convert **unstructured text** into a **structured format** so that models can understand and process them.

Types of Tokenization

Type	Example
Word Tokenization	"I love NLP" → ["I", "love", "NLP"]
Sentence Tokenization	"Hi. I am Dev." → ["Hi.", "I am Dev."]
Character Tokenization	"NLP" → ["N", "L", "P"]
Subword Tokenization	"unhappiness" → ["un", "happi", "ness"]

♦ Subword tokenization is used in modern LLMs like BERT, GPT (via BPE).

How Tokenization Works Internally (Backend Mechanism)

1. Rule-Based Tokenization (Traditional)

- Uses spaces, punctuation, and regex patterns.
- Example:
 - Split on whitespace
 - Remove or separate punctuations like ., !, ?

Limitation:

- Doesn't understand context. E.g., "Mr. Smith" can break into "Mr", ".", "Smith" — which is wrong!

2. Machine-Learned Tokenization

- Trained on massive corpora.
- Learns how humans write/speak & splits accordingly.

- Used in: spaCy, Stanford NLP

3. Subword Tokenization (Byte-Pair Encoding – BPE)

- Used in BERT, GPT, RoBERTa, etc.
- Breaks rare/unseen words into common subword pieces.

 Example:

text

CopyEdit

"unhappiness" → ["un", "happi", "ness"]

 This helps models handle unknown words better.

Popular Tokenization Tools + Code

- ♦ NLTK (for learning & basics)

python

CopyEdit

```
from nltk.tokenize import word_tokenize
word_tokenize("I love AI!") # ['I', 'love', 'AI', '!']
```

- ♦ spaCy (Fast, Production-level)

python

CopyEdit

```
import spacy
nlp = spacy.load("en_core_web_sm")
doc = nlp("I love NLP.")
[token.text for token in doc] # ['I', 'love', 'NLP', '.']
```

- ♦ HuggingFace Tokenizers (For Transformers)

python





CopyEdit

```
from transformers import BertTokenizer
```




```
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
tokenizer.tokenize("Tokenization is awesome!")
# ['token', '##ization', 'is', 'awesome', '!']
```

means it's a subword (part of previous word)

Why Tokenization is Critical in NLP?

Reason	Explanation
 Breaks raw text	Converts to processable format
 Helps models learn	Vocabulary mapping → word → vector
 Reduces unknowns	Subwords reduce "out-of-vocab" issues
 Consistency	Converts varied inputs into fixed format

Tokenization Issues in Real Life

Challenge	Example
Emojis 	"I love this  → ["I", "love", "this", "  "] — many tokenizers ignore emoji!
Contractions	"don't" → "do", "n't" OR keep as is?
Named Entities	"New York City" split into 3 tokens – but should be one!
Code / Hashtags	#LifeGoals , snake_case , camelCase – tough for standard tokenizers

Do's of Tokenization

✅ Action	📘 Why / Explanation
1. Always clean text before tokenization	Remove unwanted symbols, HTML tags, URLs — avoids junk tokens.
2. Use language-specific tokenizers	For Hindi, French, etc., use multilingual/tokenizers trained on that language (like spaCy models, mBERT).
3. Handle punctuation smartly	Use libraries like spaCy or <code>re</code> patterns if needed — "Mr. Smith" vs "I ate."
4. Normalize text (lowercase, lemmatize)	Reduces vocabulary size, helps generalization (e.g., "Running" → "run").
5. Customize tokenization for domain-specific text	E.g., programming code, social media (handle #, @, emoji).
6. Use subword tokenization for deep learning models	Models like BERT/GPT need WordPiece, BPE, etc.
7. Visualize and test token output	Always check how input is tokenized, especially in projects or deployments.
8. Keep mapping to original text if needed	Important for NER, span-based tasks (using <code>.char_span</code> in spaCy).
9. Use sentence tokenization before word tokenization (if needed)	Ensures semantic boundaries (esp. in summarization, translation).
10. Choose tokenizer based on downstream task	Simple model → <code>NLTK</code> is okay; Deep model → <code>HuggingFace Tokenizers</code> is better.

❌ Don'ts of Tokenization

❌ Mistake	⚠️ Why it's bad
1. Don't tokenize raw noisy data	Garbage in = garbage out. Pre-cleaning is a must.
2. Don't mix multiple tokenizers on same dataset	Leads to inconsistent vocab, especially in training/test.

3. Don't ignore contractions	"don't" → should be split or handled depending on task.
4. Don't ignore domain-specific symbols	Code, hashtags, medical terms may break if tokenizer isn't tuned.
5. Don't ignore alignment to original text	Especially for span-based tasks like question answering or NER.
6. Don't force subword tokenization where not needed	For classical ML (SVM, LR), BoW/TF-IDF works better.
7. Don't assume one tokenizer fits all tasks	Chatbot vs Sentiment Analysis vs QA need different preps.
8. Don't remove stopwords blindly	Sometimes "not", "no", "very" are critical in sentiment/intent detection.
9. Don't lowercase when case matters	Named Entity Recognition (NER) or Legal text may require original case.
10. Don't forget to test tokenization outputs	Always inspect edge cases like Mr. , emojis, hashtags.

Practical Tip: When to Use What?

Task	Tokenizer
Sentiment / Spam Analysis	NLTK / spaCy
Chatbot / Translation	SentencePiece / BPE
Transformers (BERT/GPT)	WordPiece / HuggingFace Tokenizer
Code / Tech Docs	Custom Regex + Subword
Social Media Text	ekphrasis , emoji libs + spaCy

"I can't believe @John_Doe loves #AI 🥰! Mr. Smith lives in U.S.A. He said, 'Don't worry.' Let's meet at 5:00 p.m."

Tokenization Output (First 20 tokens):

NLTK Word Tokens	spaCy Tokens	BERT Tokens
I	I	i
ca	ca	ca
n't	n't	##n't
believe	believe	believe
@	@	@
John_Doe	John_Doe	john
loves	loves	##_doe
#	#	loves
AI	AI	##ai
😍	😍	!
!	!	mr
Mr.	Mr.	.
Smith	Smith	smith
lives	lives	lives
in	in	in
U.S.A.	U.S.A.	u
He	He	##.
said	said	##s
,	,	##a
'	'	he

Key Observations:

- NLTK: Splits on basic rules. Handles contractions okay (**ca**, **n' t**) but merges **U.S.A.** as one (which can be good or bad).
- spaCy: Context-aware and robust for most cases. Keeps **Mr.** and **U.S.A.** intact.
- BERT: Uses subword tokenization (WordPiece). "**can' t**" → ['**ca**', "##**n' t**"], "**John_Doe**" → ['**john**', '##_**doe**'].

What Are Stopwords?





Stopwords are commonly used words in a language that carry little or no meaning and are usually removed during text preprocessing.

♦ Examples in English:

"**the**", "**is**", "**in**", "**and**", "**a**", "**an**", "**of**", "**to**", "**for**", "**on**", "**it**", "**with**"

In simple terms: Stopwords = filler words that don't help much in understanding the main idea.

Why Remove Stopwords?

Reason	Explanation
 Clean the data	Removing noise improves signal-to-noise ratio.
 Reduce dimensionality	Fewer unique words → faster & better model training.
 Focus on important terms	Keeps only meaningful tokens (like nouns/verbs).
 Speeds up vectorization	Less vocab → less memory use.

✓ Especially useful in:

- Text classification (spam, sentiment)
- Information retrieval

- Topic modeling
- Keyword extraction

Tools for Stopwords Removal

Library	Method
NLTK	<code>from nltk.corpus import stopwords</code>
spaCy	<code>spacy.lang.en.stop_words.STOP_WORDS</code>
sklearn (TF-IDF)	<code>stop_words='english'</code>
Gensim	<code>from gensim.parsing.preprocessing import STOPWORDS</code>
Custom Lists	Domain-specific or created manually

Example with NLTK:



python

CopyEdit

```
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))

tokens = ['this', 'is', 'a', 'simple', 'test']
filtered = [w for w in tokens if w not in stop_words]
# Output: ['simple', 'test']
```





Limitations & Edge Cases of Stopwords Removal

Problem	Explanation
 May remove useful info	"not good" → ["good"] if "not" is removed → wrong sentiment!
 Language dependency	English stopwords ≠ Hindi stopwords.





✗ Context matters	"to be or not to be" — poetic or philosophical context ruined.
✗ Named entities may include stopwords	"The Who" (band name), "The Office" (TV show)
✗ Not always needed	For transformers, models like BERT learn context, so stopwords can be helpful.

When to Remove Stopwords (and When Not To)

✓ Remove Stopwords:

Scenario	Why?
 Search engines / Keyword extraction	Improves search relevance
 Topic modeling	Keeps only high-meaning terms
 Spam detection, Text classification	Boosts model focus
 Classical ML models	Logistic Regression, SVM, Naive Bayes work better with cleaned input

✗ Don't Remove Stopwords:

Scenario	Why?
 Using Transformer models (BERT, GPT)	These models understand stopwords contextually
 Sentiment Analysis	"not", "no", "never" are sentiment modifiers
 Text summarization	Structural words help retain meaning
 Legal, Medical, Poetry	Domain language may need full context

Example — Impact on Sentiment

Original:

text

CopyEdit

"This movie is not good at all."

If we remove stopwords (including "not"):

text

CopyEdit

"movie good" → Model thinks it's positive ❌

➡ Wrong interpretation due to blind stopwords removal

Interview Tips

Q: What are stopwords and why are they removed in NLP?

Answer:

Stopwords are common words like "the", "is", "a", which don't add meaningful value. Removing them helps reduce noise, improve training time and focus models on useful terms. But in tasks like sentiment analysis or contextual modeling, we may choose to keep them.

Q: Can removing stopwords ever hurt model performance?

Yes, especially in sentiment analysis or when using transformers like BERT, where context matters.

Understanding Regex(Regular Expression)

✅ 1. . (Dot) – Match any single character (except newline)

🔍 What it does:

Matches exactly one character — doesn't care what character it is (except `\n` by default)

Internally:

- Regex engine sees `.` and tells:

"Okay, I'll accept any character here, just one — no questions asked."

Example:

python

Copy code

```
re.findall("c.t", "cat cut cot cbt ct")  
# Output: ['cat', 'cut', 'cot', 'cbt']
```

Why `ct` didn't match?

Because dot `.` expects one character in between — but `ct` has none → fails.

2. `^` – Match start of a string

What it does:

Matches only if the pattern starts at beginning of a string.

Internally:

- When regex sees `^`, it:

"Anchor this pattern from position 0. If the string doesn't start with this, I don't even try further."

Example:

python

Copy code

```
re.findall("^Hello", "Hello world") # ['Hello']  
re.findall("^Hello", "Say Hello world") # []
```

✓ 3. \$ – Match end of a string

🔍 What it does:

Matches pattern only if it's at the end of the string.

🔧 Internally:

- Think of it like:

“This pattern must terminate the string. If there's anything after it, I don't match.”

🔧 Example:

python

Copy code

```
re.findall("world$", "Say world") # ['world']  
re.findall("world$", "Say world now") # []
```

✓ 4. * – Match 0 or more repetitions

🔍 What it does:

Matches 0 or more of the preceding character/group — greedy by default (takes as much as it can).

🔧 Internally:

- It means:

“If the previous thing (e.g., **a**) appears any number of times, even 0, I’ll accept it.”

Example:

python

Copy code

```
re.findall("lo*", "hello loooop lost")  
# ['lo', 'loooo', 'lo']
```

5. + – Match 1 or more repetitions

What it does:

Matches 1 or more of the previous character.

Internally:

- Works just like *****, but stricter:

“At least one occurrence of the previous character is required.”

Example:

python

Copy code

```
re.findall("lo+", "hello loooop lost lo")  
# ['lo', 'loooo', 'lo', 'lo']
```

6. ? – Optional match (0 or 1 time)

What it does:

Matches zero or one occurrence of the character — also used to make things non-greedy.

Internally:

- Means:

“I’m okay if this character comes once, but I’ll also accept if it’s missing.”

Example:

python

Copy code

```
re.findall("colou?r", "color colour colouur")  
# ['color', 'colour']
```

- **u?** means → **u** is optional
-

7. [] – Character classes

What it does:

Defines a set of characters, and matches any one of them.

Internally:

- **a[bcd]z** means:

Match **abz**, **acz**, or **adz** → any of the chars inside []

Example:

python

Copy code

```
re.findall("a[bc]d", "abd acd aed") # ['abd', 'acd']
```

✓ 8. [^] – Negated character class

🔍 What it does:

Match any character except those inside.

🔧 Internally:

- [^0-9] = match anything not a digit
- Think of it as:

“Avoid these!”

🔧 Example:

python

Copy code

```
re.findall("[^aeiou]", "hello") # ['h', 'l', 'l']
```

✓ 9. { } – Quantifier (exact or range)

🔍 What it does:

Specifies how many times something should appear.

Pattern	Meaning
{3}	exactly 3 times
{2, }	at least 2 times
{1, 4}	between 1 and 4 times

🔧 Internally:

- Regex engine counts repetitions and only matches if the count fits.

Example:

python

Copy code

```
re.findall("a{2,3}", "a aa aaa aaaa")  
# Output: ['aa', 'aaa', 'aaa']
```

10. () – Grouping

What it does:

Groups multiple expressions together. Also used for capturing parts of the match.

Internally:

- Acts like parentheses in math.
- **(ha)+** means: group **ha** and match it multiple times.

Example:

python

Copy code

```
re.findall("(ha)+", "hahaha ha") # ['ha', 'ha']
```

11. | – OR (Alternation)

What it does:

Acts like logical OR → match either left or right.

Internally:

- **cat|dog** = either **cat** OR **dog**

- You can also combine with grouping.

Example:

python

Copy code

```
re.findall("cat|dog", "my cat and dog are cute")  
# ['cat', 'dog']
```

12. \ – Escape character

What it does:

Lets you use special characters as literal.

Internally:

- Without `\`, regex treats characters like `.` or `$` as commands.
- With `\`, they lose their magic.

Example:

python

Copy code

```
re.findall(r"\$", "Price: $5") # ['$']  
re.findall(r"\.", "File.txt")  # ['.']
```

Bonus: Meta Sequences (Shortcut Notations)

Shortcut	Meaning
<code>\d</code>	Digit = <code>[0-9]</code>
<code>\D</code>	Non-digit

`\w` Word =
 [\[a-zA-Z0-9_\]](#)

`\W` Non-word

`\s` Space/tab/newline

`\S` Non-space

🧠 Core Difference: `findall()` vs `search()`

Feature	<code>`re.search()`</code>	
<code>`re.findall()`</code>		
-----	-----	
Purpose	**First match** ko dhoondhta hai	
	Saare matches ko return karta hai	
Output	Match object (only 1)	List
of all matches		
Use Case	Match **exists or not** , ya **1st match** chahiye	Saari
jagah se matching **values extract** karni ho		
Returns	Match object (use <code>`group()`</code> to get value)	List
of strings		

✅ Solution: Use `re.MULTILINE` flag

python

Copy code

```
import re
```

```
match = re.findall(r"^ERROR.*", log_data, re.MULTILINE)
print(match)
```



Why re.MULTILINE is Needed?

- Normally ^ matches the start of the entire string
- With re.MULTILINE, ^ matches the start of each line inside a multi-line string



Example:

python

Copy code

```
log_data = """
```

```
INFO 01-05-2023: All systems working.
```

```
ERROR 01-05-2023: Server crashed.
```

```
WARNING: High CPU.
```

```
ERROR 02-05-2023: Memory leak.
```

```
"""
```

```
match = re.findall(r"^ERROR.*", log_data, re.MULTILINE)
```

```
print(match)
```



Output:

python

Copy code

```
['ERROR 01-05-2023: Server crashed.', 'ERROR 02-05-2023: Memory  
leak.']
```



Bonus:

Want to get line numbers of those errors?

python

Copy code

```
lines = log_data.strip().split('\n')
for i, line in enumerate(lines):
    if re.match(r"^ERROR", line):
        print(f"Line {i+1}: {line}")
```

✅ \S+ – "One or More Non-Whitespace Characters"

- + means 1 or more times
- So \S+ = match a full word or block of characters until a space

♦ Example:

python

Copy code

```
re.findall(r"\S+", "Hi Bro! Welcome123")
# Output: ['Hi', 'Bro!', 'Welcome123']
```

➡ Basically, it splits on whitespace and returns words

✅ \d+ – "One or More Digits"

- Matches full numbers, not just single digits
- + says: "keep going until non-digit appears"

♦ Example:

python

Copy code

```
re.findall(r"\d+", "Call 9876543210 or 123-456")  
# Output: ['9876543210', '123', '456']
```

Function	Purpose	Example
re.match()	Match pattern at the beginning of the string	re.match(r"\d+", "123abc")
re.search()	Search for the first match anywhere in the string	re.search(r"abc", "xyzabc123")
re.findall()	Find all non-overlapping matches and return a list	re.findall(r"\d+", "123 and 456")
re.finditer()	Return iterator yielding MatchObjects (with position)	re.finditer(r"\d+", "123 and 456")
re.sub()	Replace all matches with a specified string	re.sub(r"\d+", "[NUM]", "123 apples")
re.split()	Split string by the pattern	re.split(r"[,;]", "a,b;c")
re.compile()	Compile a regex pattern for reuse	pattern = re.compile(r"\d+")

Regex	Meaning	Example Input	Example Output
\d	Single digit (0-9)	Phone: 9876	['9', '8', '7', '6']
\d+	One or more digits	Pin 12345 is valid	['12345']
\s	Whitespace (space, tab, newline)	NLP is fun	[' ', '\t', '\n']
\S	Non-whitespace character	Dev@GPT	['D', 'e', 'v', '@', 'G', 'P', 'T']
\w	Word character (a-z, A-Z, 0-9, _)	A_b2	['A', '_', 'b', '2']
\W	Non-word character	Hello!	['!']

Syntax

Meaning

(?=... Positive lookahead (must be followed by)

(?!... Negative lookahead (must NOT be followed by)

(?<=.. Positive lookbehind (must be preceded by)

(?<!.. Negative lookbehind (must NOT be preceded by)

```
re.findall(r"@gmail(?:=\.com)", "abc@gmail.com abc@gmail.in")  
# ['@gmail']
```

```
re.findall(r"\d+(?!%)", "20% off and 500 extra")  
# ['500']
```

```
import re
```

```
text = "Price: ₹500 and $300"
```

```
# Match only numbers preceded by ₹  
matches = re.findall(r"(?<=₹)\d+", text)  
print(matches) # ['500']
```

```
text = "Price: ₹500 and $300 and 200"  
# Match numbers not preceded by ₹ or $  
matches = re.findall(r"(?<!(₹|\$)\d+", text)  
print(matches) # ['200']
```

Problem:

In "₹500", \d+ = "500", and lookbehind = (?<!₹)

Now:

- The engine checks if just before the match (i.e., before '5') is NOT '₹'
- '5' is directly preceded by ₹ → ❌ So '500' is rejected
- But then it tries matching again from the next digit i.e. '0'
 - '0' is preceded by '5' → ✅ passes (?<!₹)
 - '0' + '0' → 00 gets matched ← ⚠️ wrong!

That's why it matched partial digits ('00'), not full valid number


```
import re

text = "Price: ₹500 and $300"

# Match full numbers not preceded by ₹
matches = re.findall(r"(?!₹)\b\d+\b", text)
print(matches)
```

Why \b works:

- \b = word boundary, ensures the number is a full token
- So partial 00 in 500 doesn't get matched anymore

Pattern	Replaces	Example Input	Replacement Used	Example Output
\d+	All digit sequences	my number is 1234 and 5678	'XXXX'	my XXXX is XXXX and XXXX
\b1234\b	Only the number 1234 (exact match)	call 1234 now	'REPLACED'	call REPLACED now
\b(1234 5678)\b	Either 1234 or 5678 (exact match)	code 1234 and 5678 are test cases	'X'	code X and X are test cases
number	'number' word	phone number	'digit'	phone digit
\bis\b	'is' word (full match)	this is awesome	'was'	this was awesome
(\d{2})(\d{2})	Groups of 2 digits (e.g., 1234 → 12 34)	year: 1234	r'\1-\2'	year: 12-34

Part	Meaning
r''	Raw string – avoids double escaping \\ in Python
\(?	Optional (– match if there is a (or not (escaped with \)
\d{3}	Exactly 3 digits
\)?	Optional) – match if there is a) or not
[-.\s]	Optional separator: dash -, dot ., or space \s

`\d{3}` Next 3 digits

`[-.\s` Optional separator again
`]?`

`\d{4}` Final 4 digits

Normalization

Do's for Normalization

1. Convert tokens and then we change in lowercasing
2. Expand contractions (eg don't -> do not)
3. Handle URLs, numbers, special characters (appropriately why because if it hold any context then you don't need to remove be careful)
4. edge cases for whitespaces and punctuations
5. Spell checker
5. [we](#) use special libraries for this

Don'ts

1. Do Not remove too much (be careful that context remain conserve)
2. Don't over normalize

Edge cases

Don't -> do not if not will be removed as stopwords then it will lose the context

✅ Step 1: `import string`

- Brings in the `string` module.
- Contains useful constants like `string.punctuation`
- `string.punctuation`
- # Output: `'!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~`
- Ye saare punctuation characters hain jo hum remove karna chahte hain.

✅ Step 2: `str.maketrans('', '', string.punctuation)`

- Creates a translation table to remove punctuation.
- Syntax: `str.maketrans(from, to, delete)`
 - `from`: characters to replace (empty here)
 - `to`: characters to replace with (empty here)
 - `delete`: characters to delete

```
translator = str.maketrans('', '', string.punctuation)
Ye bol raha hai: "Koi replace mat karo, bas string.punctuation wale
characters delete kar do."
```

✅ Step 3: `text.translate(translator)`

```
cleaned_text = text.translate(translator)
```

! Caution

This method removes all punctuation, so sometimes contractions like:

- "don't" → "dont"
may lose meaning. In such cases, use smart preprocessing.

✅ Example 1: Remove Punctuation (jo tu already kiya)

python

Copy code

```
import string
```

```
text = "Hello World! Let's code:#python #AI."
translator = str.maketrans('', '', string.punctuation)
print(text.translate(translator))
# Output: "Hello World Lets codepython AI"
```

✓ Example 2: Replace characters (a → @, e → 3, i → 1)

python

Copy code

```
text = "email service is active"
translator = str.maketrans({'a': '@', 'e': '3', 'i': '1'})
print(text.translate(translator))
# Output: "3m@1 s3rv1c3 1s @ct1v3"
```

✓ Example 3: Change lowercase vowels to uppercase

python

Copy code

```
text = "this is an example"
vowels = 'aeiou'
upper_vowels = 'AEIOU'
translator = str.maketrans(vowels, upper_vowels)
print(text.translate(translator))
# Output: "thIs Is An ExAmPlE"
```

✓ Example 4: Delete specific characters (only digits)

python

Copy code

```
text = "I scored 98 marks in 2023."
translator = str.maketrans('', '', '0123456789')
print(text.translate(translator))
# Output: "I scored  marks in ."
```

Step-by-Step Explanation:

✓ `text = " Hello. world. "`

Ye original string me:

- Leading spaces (shuruaat me) → 2 space
- Multiple spaces between "Hello." and "world." → 2 spaces
- Trailing spaces (end me) → 3 spaces

Default behavior of `.split()` without any argument:

- Splits the string into words by any amount of whitespace
- Ignores extra spaces automatically
- So result:
- `['Hello.', 'world.']`

✓ `" ".join(text.split())`

- Takes the list `['Hello.', 'world.']`
- Joins the elements using a single space " " as separator
- `cleaned_text = "Hello. world."`

✓ Yes, it usually gives wrong output for names, places, technical terms, and uncommon words.

🧠 Reason: Ye library ke paas limited dictionary vocabulary hoti hai, mainly common English words ke liye.

```
from spellchecker import SpellChecker
```

```
text = "My name is dev pandey"
```

```
spell = SpellChecker()
```

```
words = text.split()
```

```
corrected_text = " ".join([spell.correction(word) for word in words])  
print(corrected_text)
```

Why This Happens?

The **SpellChecker** library uses:

- Probability-based word frequency model
- Looks up every word in its predefined dictionary (no names or custom words)
- Tries to find the closest word by edit distance (Levenshtein Distance)

Names like "dev", "pandey" are not in dictionary, so it "guesses" similar known words.

```
spell.word_frequency.load_words(['dev', 'pandey'])
```

```
spell = SpellChecker()
```

```
spell.word_frequency.load_words(['dev', 'pandey'])
```

```
text = "My name is dev pandey"
```

```
words = text.split()
```

```
corrected_text = " ".join([spell.correction(word) for word in  
words])
```

```
print(corrected_text)
```

Word Count

```
from collections import Counter
```

```
def word_frequency(text):
```

```
doc = nlp(text)
words = [token.text.lower() for token in doc]
return Counter(words)
```

```
text = "NLP is fun, NLP is powerful, NLP is amazing in future."
print(word_frequency(text))
```

What does this function do?

Step-by-step:

1. `doc = nlp(text)`

- spaCy ka tokenizer sentence ko smartly tokenize karta hai.
- Punctuation ko alag treat karta hai.
- e.g. "fun," → "fun", ",", "

2. `words = [token.text.lower() for token in doc]`

- Har token ko lowercase karta hai
- A list of all words + punctuation in lowercase form

3. `Counter(words)`

- `collections.Counter` har word ka frequency count deta hai
- i.e., {word: count} dictionary



Time Complexity Analysis

Let's assume:

- `n` = number of characters in text
- `m` = number of tokens (after tokenization)

- ♦ `nlp(text)` → $O(n)$

Tokenizing text with spaCy is approximately linear w.r.t. input size.

- ♦ List comprehension → $O(m)$

Creating a list of lowercase tokens

- ♦ `Counter(words)` → $O(m)$

Counts each word once

✓ Total Time Complexity: $O(n + m)$

- In practical terms: Linear, i.e., $O(n)$
(since $m \leq n$ for most realistic inputs)

✓ Summary

Step	What It Does	T
<code>nlp(text)</code>	Tokenizes input using spaCy	0

List of <code>token.text.lower()</code>	Converts to lowercase words	0
Counter()	Creates frequency dictionary	0
Total	Word frequency via smart tokens	0

•

◆ What is spaCy?

spaCy ek high-performance NLP (Natural Language Processing) library hai Python me, jo specially banaya gaya hai:

- Speed + accuracy ke liye
- Industry-level NLP applications me use karne ke liye

🔧 Use Cases:

- Tokenization
- POS tagging (part-of-speech)
- Named Entity Recognition (NER)

- Lemmatization
 - Dependency Parsing
 - Sentence segmentation
 - Similarity detection
 - Word vectors
 - Text classification
-

◆ What is `nlp = spacy.load("en_core_web_sm")`?

Ye line:

python

CopyEdit

```
nlp = spacy.load("en_core_web_sm")
```

ka matlab:

Load karo English ka pre-trained small language model
Ye ek AI model hota hai jo already training le chuka hai
thousands of English sentences par.

`nlp` = spaCy NLP pipeline object
Ye object input text ko process karke usme:

- Words

- Punctuation
 - Named entities
 - POS tags
 - Lemmas
- sab assign karta hai.
-

◆ What is nlp(text)?

Jab tu likhta hai:

```
python
```

```
CopyEdit
```

```
doc = nlp(text)
```

To spaCy:

1. Tokenizes the text (i.e., breaks into words, punctuations, etc.)
2. Analyzes linguistically – jaise ki:
 - Har word ka POS tag
 - Har word ka lemma
 - Sentence boundaries
 - Named entities

3. Returns a Doc object – jisme rich info hoti hai har word ke baare me

Example:

python

CopyEdit

```
import spacy
```

```
nlp = spacy.load("en_core_web_sm")
```

```
text = "Apple is looking at buying a startup in the UK."
```

```
doc = nlp(text)
```

```
for token in doc:
```

```
    print(token.text, "|", token.pos_, "|", token.lemma_)
```

Output:

less

CopyEdit

```
Apple      | PROPN | Apple
```

```
is         | AUX   | be
```

```
looking    | VERB  | look
```

at	ADP	at
buying	VERB	buy
a	DET	a
startup	NOUN	startup
in	ADP	in
the	DET	the
UK	PROP	UK
.	PUNCT	.



Summary

Concept	Meaning
<code>spacy</code>	A powerful NLP library for analyzing and processing human language text
<code>nlp</code>	NLP pipeline object created using <code>spacy.load()</code>
<code>nlp(text)</code>	Returns a Doc object after analyzing the text
<code>Doc</code>	Sequence of tokens (words) with linguistic annotations

✓ Real World Analogy:

Imagine spaCy is a language expert, and `nlp(text)` is like giving that expert a sentence. It returns a full analysis report of that sentence:

- Where the words are
- What type of words they are
- Which words are names
- What is the meaning root (lemma) of each word, etc.

Please go through all the NLP Notebooks

Frequency Distribution

.Bag of words

Term frequency-Inverse Document Frequency (Tf-IDF)

Bag of Words(BOW)

What is Bag of Words (BoW)?

Bag of Words is a method to convert **text into numbers** so that machine learning algorithms can understand it.

✓ Core Idea:

Treat text like a **bag** that contains words — just words — and **ignore grammar, order, and context**.



Example:

Let's say you have 2 sentences (called "documents" in NLP):

Doc1: "I love NLP"

Doc2: "NLP is great and I love it"

♦ Step 1: Vocabulary creation (Unique words):

css

Copy code

```
['I', 'love', 'NLP', 'is', 'great', 'and', 'it']
```

♦ Step 2: Frequency table (how many times each word appears in each sentence):

Word	Doc1	Doc2
I	1	1
love	1	1
NLP	1	1
is	0	1
great	0	1
and	0	1
it	0	1

o, the final representation:

- Doc1 = [1, 1, 1, 0, 0, 0, 0]
- Doc2 = [1, 1, 1, 1, 1, 1, 1]

Where is Bag of Words used?

- Spam detection (emails)

- Sentiment analysis
- Text classification (movie reviews, product reviews)
- Document similarity
- Keyword extraction

⚠ Limitations of Bag of Words:

✗ Limitation

🧠 What It Means

No context	Doesn't understand meaning or sentence structure
No word order	"I love NLP" and "NLP love I" will be same
Sparse representation	Large vocab → huge matrix with lots of 0s
Doesn't handle synonyms	"awesome" and "great" are treated as different words
Ignores semantics	Can't distinguish between "not good" and "very good"

✅ When to Use Bag of Words?

Use It When...

Avoid It When...

Data is small or medium-sized	Context and word order are important
Quick, interpretable models are needed	You need semantic understanding (e.g., BERT)
Doing baseline experiments	Working on chatbots or translation tasks

🛠 Tools to Create BoW in Python:

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
corpus = ["I love NLP", "NLP is great and I love it"]
```

```
vectorizer = CountVectorizer()
```

```
X = vectorizer.fit_transform(corpus)
```

```
print(vectorizer.get_feature_names_out()) # ['and' 'great' 'is' 'it' 'love' 'nlp']
```



```
print(X.toarray())
```

Term Frequency inverse Document Frequency(TF-IDF)

What is TF-IDF?

TF-IDF stands for:

- **TF** → Term Frequency
- **IDF** → Inverse Document Frequency

Purpose:

To measure **how important a word is** in a document **relative to a collection of documents (corpus)**.

BoW just counts frequency, but TF-IDF asks:

"Is this word really meaningful **in the context of all documents**, or is it just common everywhere?"

Real-Life Analogy:

Imagine you're reading 100 resumes.

Everyone writes "hardworking", "motivated" — it appears everywhere.

But "neural networks", "docker", "NLP" — appears in few resumes.

TF-IDF gives **less weight** to common words
and **more weight** to rare-but-important words. 🔥

TF-IDF: The Formula

Let's break it down.

1. TF (Term Frequency):

How often the word appears in the document.

$TF(t,d) = \frac{\text{Number of times } t \text{ appears in document } d}{\text{Total number of terms in } d}$
 $TF(t,d) = \frac{\text{Number of times } t \text{ appears in document } d}{\text{Total number of terms in } d}$

Example:

- "I love NLP. NLP is fun"
- Term = "NLP"
- Appears 2 times in 6 words $\rightarrow TF = 2/6 = 0.33$

2. IDF (Inverse Document Frequency):

How rare the word is across **all documents**.

$IDF(t) = \log\left(\frac{N}{1 + DF(t)}\right)$ $IDF(t) = \log(1 + DF(t)N)$

- N = total number of documents
 - $DF(t)$ = number of documents containing the term t
 - Log is used to dampen the effect
- ♦ If word appears in **every doc** $\rightarrow IDF \approx 0$ (less important)
 - ♦ If word appears in **few docs only** $\rightarrow IDF \uparrow$ (more important)

3. TF-IDF Score:

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

Example

Corpus:

vbnet

Copy code

Doc1: "I love NLP"

Doc2: "NLP is powerful and fun"

Doc3: "I love machine learning"

Word = "NLP"

- TF in Doc1 = $1/3$
- Appears in 2/3 docs \rightarrow $IDF = \log(3 / (1 + 2)) = \log(1) = 0$
- So TF-IDF = 0

Word = "machine"

- TF in Doc3 = $1/4$

Appears in 1/3 docs \rightarrow $IDF = \log(3 / (1 + 1)) = \log(1.5) \approx 0.176$

So TF-IDF > 0

```
from sklearn.feature_extraction.text import TfidfVectorizer
```


```
corpus = ["I love NLP", "NLP is powerful and fun", "I love machine learning"]
```

```
vectorizer = TfidfVectorizer()
```

```
X = vectorizer.fit_transform(corpus)
```

```
print(vectorizer.get_feature_names_out())
```


```
print(X.toarray())
```

 Why TF-IDF is Powerful in NLP:

Advantage

Penalizes common words	Reduces noise from stopwords or generic words
Highlights unique, rare terms	More meaningful for classification, search, similarity tasks
Simple + interpretable	Easy to visualize and explain
Great for text classification	Spam detection, sentiment analysis, etc.

Explanation

 Limitations of TF-IDF:

Limitation

Explanation

Doesn't capture word order

Like BoW, it's a bag — "good not" ≠ "not good"

Doesn't capture context or meaning

"Java" in island vs. "Java" in programming

Large corpus → sparse matrix

Memory inefficient for big data

Static — doesn't learn semantics

No training; fixed vector per word

you calculated:

- $TF = 1/5 = 0.2$
- $IDF = \log(3/2) \approx 0.176$
- So, $TF \times IDF = 0.2 \times 0.176 = 0.0352$ ❌

But the model gave 0.534 — why?

Explanation: TF-IDF in scikit-learn is NOT raw $TF \times \log(N/DF)$

By default, **TfidfVectorizer** in **scikit-learn** uses **L2 normalization** and **sublinear TF scaling**, unless explicitly turned off.

Actual formula used by sklearn:

1. TF → by default, it uses raw count

(Unless `sublinear_tf=True`, which uses $1 + \log(tf)$)

2. IDF → uses the formula:

2. IDF → uses the formula:

$$IDF(t) = \log \left(\frac{1 + N}{1 + DF(t)} \right) + 1$$

This avoids division by zero and smooths the scores.

So for "amazing":

- Appears in 1 doc \rightarrow DF = 1
- N = 3 docs

$$\text{IDF}(\textit{amazing}) = \log\left(\frac{1+3}{1+1}\right) + 1 = \log(2) + 1 \approx 0.693 + 1 = 1.693$$

TF (raw count) = 1

TF-IDF (before normalization) = $1 \times 1.693 = 1.693$

Final Step: L2 Normalization

Each document vector is **scaled** so that the square root of sum of squares = 1.

That's why even though "amazing" has raw score 1.693, it becomes **~0.534** after normalization

What is L2 Normalization?

L2 normalization (also called **Euclidean normalization**) means:

"Scale the vector so that the sum of the squares of all values equals 1."

Doc1 (before normalization): [1.693, 1.0, 0.5, 0.2]



How L2 Normalization Works (Mathematically)

Let's say your vector is:

$$\vec{v} = [v_1, v_2, \dots, v_n]$$

Then L2 norm (length of vector):

$$\|\vec{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2}$$

Each value becomes:

$$v'_i = \frac{v_i}{\|\vec{v}\|_2}$$



Example (From Your Case):

Let's say the **raw TF-IDF vector** for **Doc1** is:

text

Copy code

```
["amazing", "fun", "love", "nlp"] → [1.693, 1.693, 0.0, 0.5]
```

Step 1: Calculate L2 norm

$$\|v\| = \sqrt{1.693^2 + 1.693^2 + 0 + 0.5^2} \approx \sqrt{2.865 + 2.865 + 0 + 0.25} \approx \sqrt{5.98} \approx 2.445$$

Step 2: Normalize each element

$$\text{amazing normalized} = \frac{1.693}{2.445} \approx 0.692$$

Same for all elements.



Why Is This Important?

Reason	Benefit
Equalizes document length	Long documents don't dominate
Useful in Cosine Similarity	Normalized vectors simplify math
Keeps relative term importance	Magnitudes scale, but meaning stays



When the Real-World Hits:

Let's say:

- You have 1 million news articles.
- Some are 20 words long, some are 3000.
- Common words: "said", "India", "company"
- Rare words: "inflation-adjusted bond yield"

If you don't normalize or smooth:

- Common words dominate
- Long docs overpower short ones
- Model overfits meaningless tokens

So we improve the formula to make it:

- ✓ More stable
- ✓ More generalizable
- ✓ Better for real-world performance

✅ Example 2: Get top 3 feature importances

python

Copy code

```
importances = np.array([0.1, 0.5, 0.3, 0.2])  
  
top_features = np.argsort(importances[::-1])[:3]  
  
print(top_features) # Output: [1 2 3]
```

💡 You now know **which features are most important** by index!

ANNOTATORS

- 1.What are annotators
- 2.types of annotators
- 3.Mathematical intuition behind NER,POS
- 4.Practical implementation
- 5.NLP task

🧠 What is Annotation in NLP?

Annotation in NLP means:

“Tagging raw text with useful labels or metadata to give it structure and meaning.”

Example (Without annotation)

Rahul went to Delhi on 25th June 2025.

It's just raw text — machines don't *understand* it.

Token	POS Tag	Entity Type
-------	---------	-------------

Rahul	NNP	PERSON
Delhi	NNP	LOCATION
25th June 2025	CD/DATE	DATE

Why Annotation Is So Important?

Because **models learn from labeled data**.

Without annotation, all your deep learning models are like:

“I see words... but what do they *mean* to humans?”

Analogy:

Imagine teaching a 2-year-old.

- Show picture of a dog → say “dog”
- Do this 1000 times → they learn to recognize dogs

Same for ML — you show:

How Do We Annotate?

Manual Annotation:

- Human experts label data using tools like:
 - Prodigy
 - Doccano
 - Label Studio

- Time-consuming but accurate

✅ Automated Annotation:

- Use pre-trained models (like SpaCy, HuggingFace) to label
- Faster but may require manual correction

🚀 Real-Life Use Cases:

Domain	NLP Task	Annotation Example
Healthcare	Entity Extraction	["fever": SYMPTOM]
E-commerce	Sentiment Classification	["This product is bad": NEGATIVE]
Finance	News Classification	["RBI increased rate": POLICY]

ANNOTATOR CREATION

🧠 1. What is NER (Named Entity Recognition)?

NER ka goal hota hai:

Identify and classify **entities** (real-world objects) in a sentence — like **person**, **organization**, **location**, **date**, **time**, **money**, etc.

📌 Example:

Text:

Rahul went to Delhi in June 2023 to join Google.

NER Output:

Word	Entity
Rahul	PERSON
Delhi	LOCATION
June 2023	DATE
Google	ORGANIZATION

Virat Kohli scored 100 in Delhi in October 2023.



Step 2: Mathematical Thinking — Why Do We Need It?

Mathematical question:

"Har word ke liye sabse sahi tag (PER, LOC, etc.) kya hai?"

Yahi kaam karte hain:

- **HMM**: purani probability se sikhta
- **CRF**: pura sentence dekh ke sikhta
- **Viterbi**: best tag ka sequence nikalta
- **Bayes**: chance nikalta kisi tag ke hone ka



Step 3: Hidden Markov Model (HMM)



Think like this:

- **Words** = dikhta hua text (e.g., "Virat")

- **Tags** = chhupi hui cheeze (e.g., "PER")

⚙️ **HMM uses 2 probabilities:**

1. Transition Probability

$$P(\text{Next Tag} \mid \text{Current Tag})$$

e.g. PERSON ke baad LOCATION aane ka chance

2. Emission Probability

$$P(\text{Word} \mid \text{Tag})$$

e.g. "Virat" aata hai PERSON tag ke saath kitni baar


 **Example:**

Sentence: **Virat plays in Delhi**

Let's say model knows:

- "Virat" = PERSON 90%
- "Delhi" = LOCATION 80%
- Transition: PERSON → LOCATION = 50%

To model sochta:

- "Virat" likely PERSON
- "Delhi" likely LOCATION
- Aur dono ka connection bhi acha hai → 

But pure sentence me *best tag ka sequence* kaise milega?
Uske liye aata hai...

Real-life Example: Weather Prediction

Imagine:

- You can't see the weather directly.
- But you can observe what someone is doing:
 - If they carry an umbrella → maybe it's raining
 - If they go for a walk → maybe it's sunny

HIDDEN STATES (Weather):

- Rainy
- Sunny

 **OBSERVATIONS (What we see):**

- Walk
- Shop
- Clean

 **HMM Needs Two Things**

1. Transition Probabilities

(Moving from one hidden state to another)

From → To	Rain	Sunny
-----------	------	-------

Rainy	0.7	0.3
-------	-----	-----

Sunny	0.4	0.6
-------	-----	-----

2. Emission Probabilities

(Probability of observing something in a hidden state)

Weather	Walk	Shop	Clean
---------	------	------	-------

Rainy	0.1	0.4	0.5
-------	-----	-----	-----

Sunny	0.6	0.3	0.1
-------	-----	-----	-----

3. Initial Probabilities

(Starting chance of a weather)

Weather	Probability
---------	-------------

Rainy	0.6
-------	-----

Sunny	0.4
-------	-----

 **Observed Sentence (3 words):**

css

Copy code

`["Time", "flies", "fast"]`

 **Hidden Tags:**

- N (Noun)
- V (Verb)

 **Initial Probabilities:**

$P(N)=0.6, P(V)=0.4$
 $P(N) = 0.6, \quad P(V) = 0.4$

 **Transition Probabilities**

From	→ To	N	V
------	------	---	---

N	0.	0.
	3	7

V	0.	0.
	8	2

🔗 Emission Probabilities:

Ta	Time	fli	fa
g		es	st

N	0.5	0.1	0.
			1

V	0.1	0.6	0.
			3

■ Step 1: First Word = "Time"

$$V_1(N) = P(N) \cdot P(\text{Time}|N) = 0.6 \cdot 0.5 = 0.3$$

$$V_1(V) = P(V) \cdot P(\text{Time}|V) = 0.4 \cdot 0.1 = 0.04$$

Step 2: Word = "flies"

low calculate for each tag using all paths from previous step:

$$\begin{aligned} \gamma_2(N) &= \max [V_1(N) \cdot P(N \rightarrow N), V_1(V) \cdot P(V \rightarrow N)] \cdot P(\text{flies}|N) \\ &= \max[0.3 \cdot 0.3, 0.04 \cdot 0.8] \cdot 0.1 = \max[0.09, 0.032] \cdot 0.1 = 0.09 \cdot 0.1 = 0. \end{aligned}$$

$$\begin{aligned} \gamma_2(V) &= \max [V_1(N) \cdot P(N \rightarrow V), V_1(V) \cdot P(V \rightarrow V)] \cdot P(\text{flies}|V) \\ &= \max[0.3 \cdot 0.7, 0.04 \cdot 0.2] \cdot 0.6 = \max[0.21, 0.008] \cdot 0.6 = 0.21 \cdot 0.6 = 0. \end{aligned}$$

Step 3: Word = "fast"

$$V_3(N) = \max [0.009 \cdot 0.3, 0.126 \cdot 0.8] \cdot P(\text{fast}|N) = \max[0.0027, 0.1008]$$

$$V_3(V) = \max [0.009 \cdot 0.7, 0.126 \cdot 0.2] \cdot P(\text{fast}|V) = \max[0.0063, 0.0252]$$

Final Step: Best Path

Now check which tag has highest final probability:

- N: 0.01008
- V: 0.00756

So best final tag = N, and we backtrack using max paths from each step to get the full sequence.

Initial Probability – Where It Comes From?

It comes from training data. You count:

For example, if out of 100 sentences:

Initial Probability — Where It Comes From?

It comes from **training data**. You count:

$$P(N) = \frac{\text{Total N at start}}{\text{Total sentences}}, \quad P(V) = \frac{\text{Total V at start}}{\text{Total sentences}}$$

60 start with a noun $\rightarrow P(N) = 0.6$

- 40 start with a verb $\rightarrow P(V) = 0.4$

This becomes our initial probability vector.

Rahul/NN is/VBZ in/IN Delhi/NNP

Assume we have 1000 such sentences in training data, where words are tagged like above.

Now count:

- How many times **VBZ** follows **NN**
- How many total times **NN** appears as previous tag

Previous Tag (tag ₁)	Next Tag (tag ₂)	Count
NN	VBZ	200
NN	NNP	100
NN	IN	300



Formula:

$$P(VBZ \mid NN) = \frac{\text{Count}(NN \rightarrow VBZ)}{\text{Total transitions from NN}} = \frac{200}{600} = 0.33$$

- Total words tagged NNP = 1000
- Total words tagged VBZ = 800



Formula:

$$P(Delhi \mid NNP) = \frac{\text{Count}(\text{Delhi}, \text{NNP})}{\text{Total words tagged as NNP}} = \frac{300}{1000} = 0.3$$



1. Count Transitions

python

Copy code

```
from collections import defaultdict

transition_counts = defaultdict(lambda: defaultdict(int))

for sentence in corpus:
    prev_tag = None
    for word, tag in sentence:
        if prev_tag is not None:
            transition_counts[prev_tag][tag] += 1
        prev_tag = tag

# Now calculate transition probabilities
transition_probs = {}
for prev_tag, next_tags in transition_counts.items():
    total = sum(next_tags.values())
```

```
transition_probs[prev_tag] = {tag: count / total for tag, count in
next_tags.items()}
```

✓ 2. Count Emissions

python

Copy code

```
emission_counts = defaultdict(lambda: defaultdict(int))
tag_counts = defaultdict(int)

for sentence in corpus:
    for word, tag in sentence:
        emission_counts[tag][word] += 1
        tag_counts[tag] += 1

# Now calculate emission probabilities
emission_probs = {}
for tag, words in emission_counts.items():
    emission_probs[tag] = {word: count / tag_counts[tag] for word,
count in words.items()}
```

✓ What is POS Tagging?

POS Tagging = Part-of-Speech Tagging

It is the process of assigning a grammatical label (tag) to each word in a sentence based on its role in the sentence.

Example:

vbnet

Copy code

Sentence: "The cat sat on the mat."

POS Tags: Det Noun Verb Prep Det Noun

Each word is tagged as:

- **Det** – Determiner (The)
 - **Noun** – Noun (cat, mat)
 - **Verb** – Verb (sat)
 - **Prep** – Preposition (on)
-

Why is POS Tagging Important?

- Helps in syntactic parsing
 - Supports text-to-speech systems
 - Improves Named Entity Recognition (NER)
 - Essential in dependency parsing, question answering, and semantic analysis
-

How Does POS Tagging Work? (Mathematical Intuition)

We treat POS tagging as a sequence labeling problem – same as NER.

We want to find the most probable sequence of POS tags for a given sequence of words.

1. Using Hidden Markov Model (HMM)

The model assumes:

- We have a sequence of observed words

- POS tags are hidden states

Goal:

Find: $\arg\max_T P(T|W)$

Where:

- T = sequence of POS tags
- W = sequence of words

Using Bayes Theorem:

$$P(T|W) = \frac{P(W|T) \cdot P(T)}{P(W)}$$

Since $P(W)$ is constant, we maximize:

$$P(W|T) \cdot P(T)$$

This leads to:

- Emission Probability: $P(\text{word} | \text{tag})$
- Transition Probability: $P(\text{tag}_i | \text{tag}_{i-1})$

Finally, we use the Viterbi Algorithm to find the most probable tag sequence efficiently.

Modern Approaches:






1. Rule-based (obsolete mostly)
2. Statistical models:
 - HMM

- CRF (Conditional Random Field)

3. Neural models:

- BiLSTM + CRF
- Transformers (BERT-based taggers)

POS Tagging vs NER – What's the Difference?

Feature	POS Tagging	Named Entity Recognition (NER)
 Purpose	Identifies grammatical role	Identifies real-world named entities
 Tags Used	Noun, Verb, Adj, Adv, etc.	PERSON, ORG, LOC, MISC, etc.
 Output Example	"Time/NN flies/VB fast/RB"	"Apple/ORG released/VB iPhone/PRODUCT"
 Approach	Sequence tagging (HMM/CRF)	Also sequence tagging (HMM/CRF/BERT)
 Use-case	Syntax parsing, translation	Info extraction, Q&A, Chatbots

Example:

Sentence: "Barack Obama was born in Hawaii."

- POS Tags:
 - Barack/NNP Obama/NNP was/VBD born/VBN in/IN Hawaii/NNP
- NER Tags:

- Barack Obama / PERSON
- Hawaii / LOCATION

✓ Summary

- POS tagging helps understand the structure of a sentence.
- It uses models like HMM, CRF, and now transformers.
- POS tagging and NER both are sequence labeling problems but serve different purposes.

```
import nltk

# Sample text
text = "Barack Obama was born in Hawaii and became the President of
the United States."

# Step 1: Tokenize the text into words
tokens = nltk.word_tokenize(text)

# Step 2: Apply POS tagging
pos_tags = nltk.pos_tag(tokens)

# Print the results
for word, tag in pos_tags:
    print(f"{word} --> {tag}")

Barack --> NNP
Obama --> NNP
was --> VBD
born --> VBN
```


■ Sample Output:

```
in --> IN
Hawaii --> NNP
and --> CC
became --> VBD
the --> DT
President --> NN
of --> IN
the --> DT
United --> NNP
States --> NNPS
. --> .
```

```
import nltk
```

```
nltk.download('punkt')           # For tokenization
nltk.download('maxent_ne_chunker') # For Named Entity Chunker
nltk.download('words')           # Wordlist for NE chunker
nltk.download('averaged_perceptron_tagger')
from nltk import word_tokenize, pos_tag, ne_chunk
```

```
text = "Barack Obama was born in Hawaii and served as the 44th
President of the United States."
```

```
# Step 1: Tokenize
```

```
tokens = word_tokenize(text)
```

```
# Step 2: POS Tagging
```

```
tagged = pos_tag(tokens)
```

```
# Step 3: Named Entity Recognition
```

```
ner_tree = ne_chunk(tagged)
```

```
# Step 4: Display the Named Entities
```

```
print(ner_tree)
```

✓ 1. POS (Part-of-Speech) Tags in NLTK (Penn Treebank Tags)

Tag	Meaning	Example
NN	Noun, singular	dog, car
NNS	Noun, plural	dogs, cars
NNP	Proper noun, singular	Obama, India
NNPS	Proper noun, plural	Americans
VB	Verb, base form	eat, run
VBD	Verb, past tense	ate, ran
VBG	Verb, gerund	eating, running
VBN	Verb, past participle	eaten, written
VBP	Verb, non-3rd person sing.	eat, run
VBZ	Verb, 3rd person singular	eats, runs
JJ	Adjective	beautiful, tall
JJR	Adjective, comparative	taller
JJS	Adjective, superlative	tallest
RB	Adverb	quickly, silently
RBR	Adverb, comparative	faster
RBS	Adverb, superlative	fastest

IN	Preposition	in, on, at
DT	Determiner	the, a, an
PRP	Personal pronoun	I, you, they
PRP \$	Possessive pronoun	my, your, their
CC	Coordinating conjunction	and, but, or
UH	Interjection	oh, wow
CD	Cardinal number	one, 2, 100
EX	Existential "there"	there is
FW	Foreign word	déjà vu
MD	Modal	can, will, should
POS	Possessive ending	's
SYM	Symbol	\$, %, +
TO	"to"	to go
WDT	Wh-determiner	which, that
WP	Wh-pronoun	who, what
WRB	Wh-adverb	where, when

```
import nltk
nltk.help.upenn_tagset('NNP')
```

✓ 2. NER Tags in NLTK (Chunk Labels)

Tag	Meaning	Example
PERSON	People's names	Barack Obama, Elon Musk
ORGANIZATION or ORGANISATION	Companies, institutions	Google, NASA
GPE	Geo-political entity	India, New York
LOCATION	Places	Himalayas, Sahara
FACILITY	Buildings, airports, highways	Golden Gate Bridge
DATE	Date expressions	10 Jan, 2023, yesterday
TIME	Time expressions	5 PM, 14:00
MONEY	Monetary values	\$5, ₹500
PERCENT	Percentage expressions	90%, 12 percent
ORGANIZATION	Institutions, Companies	UN, Apple

```
from nltk.tree import Tree
if isinstance(chunk, Tree):
    print(chunk.label()) # PERSON, GPE, etc.
```

Sentence Embedding

 What is Sentence Embedding?

A **sentence embedding** is a fixed-size vector (array of numbers) that **represents the meaning of an entire sentence**.

Just like:

- Word embedding represents a word → vector (like Word2Vec, GloVe)
- Sentence embedding represents a sentence → vector (meaning-aware)

Goal: Capture the overall **semantics (meaning)** of a sentence so that machines can understand and compare sentences numerically

Why Do We Need Sentence Embeddings?

To perform tasks like:

- Text Similarity
- Sentence Classification
- Question Answering
- Information Retrieval
- Paraphrase Detection

A model needs more than just individual word meanings — it needs to understand full sentences.

Simple Example

Suppose you have two sentences:

ini

Copy code

```
S1 = "I love playing football."
```

```
S2 = "Football is my favorite game."
```

A good sentence embedding will give vectors that are **very close** to each other because both talk about loving football.

But compare with:

ini

Copy code

```
S3 = "The earth revolves around the sun."
```

Its vector should be **far** from S1 and S2, because it's about science, not sports.

Ways of Sentence Embedding

1. Average Word Embedding

Formula= $V(S) = 1/n \cdot (\text{Summation } V(w_i))$

$V(s)$ = Sentence Embedding

$V(w_i)$ = i th word in your Sentence

n = Total number of words in the sentence

Average Word Embedding treat weight of every word equally

You already know **Word Embedding** (like Word2Vec or GloVe) represents each word as a dense vector.

Now to represent a **sentence**, we:

💡 **Take all the word vectors in a sentence → then average them.**

This gives us a **single vector** that represents the **overall meaning** of the sentence.

✅ Why Use It?

- Very fast
- Simple to implement
- Captures basic semantics
- Works well for small/medium datasets

Sentence: "I love NLP"

Let's say our word embeddings are 3-dimensional (to keep it simple):

Word	Embedding Vector
"I"	[0.1, 0.3, 0.5]
"love"	[0.7, 0.6, 0.1]
"NLP"	[0.9, 0.2, 0.4]

👉 Step-by-step:

Step 1: Add the vectors:

markdown

Copy code

```
[0.1, 0.3, 0.5]
+ [0.7, 0.6, 0.1]
+ [0.9, 0.2, 0.4]
-----
= [1.7, 1.1, 1.0]
```

Step 2: Divide by the number of words (3):

csharp

Copy code

```
[1.7 / 3, 1.1 / 3, 1.0 / 3] = [0.566, 0.366, 0.333]
```

🎯 This is the **Average Sentence Embedding** of "I love NLP".

❌ Limitations of Average Word Embedding

🔍 Issue	❌ Explanation
No context awareness	"Apple is tasty" vs. "Apple released a new iPhone" — both will treat "Apple" the same.
Ignores word order	"I love NLP" and "NLP loves I" will have the same vector .
Stopwords dilute meaning	Words like "the", "is", "a" can dominate if not removed.

Equal importance to all words

“not good” and “good” are treated similarly even though their meanings are opposite.

Not good for long/complex sentences

It doesn't capture structure, negation, sarcasm, etc.

✓ When to Use It

✓ Best For

Examples

Quick similarity search

Finding similar questions or documents

Clustering or visualization

Topic grouping, 2D plots (with PCA/TSNE)

Input for classic ML models

Logistic regression, SVM, KNN, etc.

Baseline for NLP tasks

Compare against better embeddings later

[2.TF](#)-IDF weighted sentence embeddings

It will assign higher weights to less common words



What is TF-IDF Weighted Sentence Embedding?

Instead of giving **equal weight to all words** (like average embedding),

👉 we **weight each word vector using its TF-IDF score** — this gives **important words more impact** on the final sentence vector.

$$\text{Sentence vector} = \frac{\sum(\text{TF-IDF}(\text{word}) \times \text{WordEmbedding}(\text{word}))}{\sum(\text{TF-IDF scores})}$$

✓ Manual Example (Step-by-step)

♦ Suppose Sentence:

arduino

Copy code

"I love NLP"

◆ **Word Embeddings (3-dimensional for simplicity):**

Word Embedding Vector

I [0.1, 0.3, 0.5]

love [0.7, 0.6, 0.1]

NLP [0.9, 0.2, 0.4]

◆ **TF-IDF Scores:**

Word TF-IDF Score

I 0.2

love 0.6

NLP 0.9



Multiply Word Vectors by Their TF-IDF Scores

Weighted Vectors:

Copy code

$$0.2 \times [0.1, 0.3, 0.5] = [0.02, 0.06, 0.10]$$

$$0.6 \times [0.7, 0.6, 0.1] = [0.42, 0.36, 0.06]$$

$$0.9 \times [0.9, 0.2, 0.4] = [0.81, 0.18, 0.36]$$

+ Add Weighted Vectors:

csharp

Copy code

$$[0.02 + 0.42 + 0.81, 0.06 + 0.36 + 0.18, 0.10 + 0.06 + 0.36] \\ = [1.25, 0.60, 0.52]$$

÷ Normalize by Total TF-IDF:

java

Copy code

Total TF-IDF = $0.2 + 0.6 + 0.9 = 1.7$

Final Sentence Embedding = $[1.25/1.7, 0.60/1.7, 0.52/1.7]$
 $\approx [0.735, 0.353, 0.306]$

🎯 **This vector better captures sentence meaning** than plain average!



Python Code (Using Gensim + Sklearn)

python

Copy code

```
from sklearn.feature_extraction.text import TfidfVectorizer
from gensim.models import Word2Vec
import numpy as np

# Step 1: Sample corpus and model
corpus = ["I love NLP", "NLP is great", "I enjoy machine learning"]
tokenized = [doc.lower().split() for doc in corpus]
model = Word2Vec(tokenized, vector_size=50, min_count=1)

# Step 2: TF-IDF calculation
tfidf = TfidfVectorizer()
tfidf.fit([" ".join(doc) for doc in tokenized])
word2tfidf = dict(zip(tfidf.get_feature_names_out(), tfidf.idf_))

# Step 3: Sentence embedding function
def tfidf_weighted_embedding(sentence):
    words = sentence.lower().split()
    weighted_vectors = []
    total_weight = 0
    for word in words:
        if word in model.wv and word in word2tfidf:
            weight = word2tfidf[word]
            weighted_vectors.append(model.wv[word] * weight)
            total_weight += weight
    if weighted_vectors:
        return np.sum(weighted_vectors, axis=0) / total_weight
```

```
else:
    return np.zeros(model.vector_size)

# Example
vec = tfidf_weighted_embedding("I love NLP")
print("Sentence embedding vector shape:", vec.shape)
```

❌ Limitations of TF-IDF Weighted Embedding

Limitation	Why it's a Problem
Ignores word order	Like average embeddings
Still static word vectors	"Bank" in "river bank" vs "money bank" — same
TF-IDF not deep in semantics	It's based on frequency, not meaning
Domain sensitivity	TF-IDF varies by corpus

🧠 Why PCA or SVD in Sentence Embeddings?

When you use **Average Word Embeddings**, many times:

- Common dimensions (like "the", "is", "are") **dominate** the embedding space.
- The embeddings become **less discriminative**.

🔥 Solution:

Use **PCA/SVD** to:

- Remove noise/common directions.
- Reduce dimensionality.
- Retain only the most **informative components**.

✓ What is PCA?

PCA finds the **directions (components)** in your data that explain the most **variance** and projects your data onto them.

In NLP:

- Take sentence embeddings (e.g., 300D vectors)
- Run PCA to reduce them to 100D or 50D
- Or remove the top 1 component to reduce noise

📌 Real Use: "SIF Embeddings" (Smooth Inverse Frequency)

This method:

1. Weighs word vectors using word frequency (like TF-IDF)
2. Averages them to form a sentence vector
3. Removes the **first principal component** (common words direction)

🧠 Why PCA or SVD in Sentence Embeddings?

When you use **Average Word Embeddings**, many times:

- Common dimensions (like "the", "is", "are") **dominate** the embedding space.
- The embeddings become **less discriminative**.

🔥 Solution:

Use **PCA/SVD** to:

- Remove noise/common directions.
 - Reduce dimensionality.
 - Retain only the most **informative components**.
-

What is PCA?

PCA finds the **directions (components)** in your data that explain the most **variance** and projects your data onto them.

In NLP:

- Take sentence embeddings (e.g., 300D vectors)
 - Run PCA to reduce them to 100D or 50D
 - Or remove the top 1 component to reduce noise
-

Real Use: "SIF Embeddings" (Smooth Inverse Frequency)

This method:

1. Weights word vectors using word frequency (like TF-IDF)
2. Averages them to form a sentence vector
3. Removes the **first principal component** (common words direction)

Example:

Sentence	Word	Meaning
"He went to the bank to fish"	bank	river bank (nature)
"He went to the bank to deposit cash"	bank	financial bank (institution)

- 👉 In **Word2Vec**: "bank" → same vector
- 👉 In **Contextual Embedding**: "bank" → different vectors for each sentence

⚙️ How it Works (Background)

Contextual embeddings come from **deep language models** trained on **massive text corpora**. Popular models:

- ELMo (2018)
- BERT (2019)
- RoBERTa, GPT, etc.

📖 Architecture Behind It

Contextual embeddings use:

- **RNNs (ELMo)** or
- **Transformers (BERT)**

In BERT, each word is represented as:

- A **token** (after subword tokenization)
- Passed through **multiple transformer layers**
- At each layer, attention allows it to **look at surrounding words** (context)
- Final layer output gives the **contextual vector**

✅ Advantages

Strength	Why it's powerful
Understands polysemy	Different meanings → different vectors
Considers full sentence	Not just fixed window

Boosts downstream accuracy

Improves NER, QA, classification

Transfer learning ready

Pretrained on huge data (e.g., BERT)

✗ Limitations

Limitation	Why it's a problem
Heavy computation	Models like BERT are large (110M+ params)
Needs GPU for speed	Especially for large-scale tasks
Not interpretable	Hard to explain why a vector looks like that
Latency	Slow for real-time applications
Context length limit (512)	Cannot handle very long documents at once

Stemming And Lemmatization

🔍 What is Stemming?

Stemming is the process of reducing a word to its **root/base form** by chopping off prefixes or suffixes.

♦ Example:

- "playing" → "play"
- "played" → "play"
- "flies" → "fli" (not always perfect!)

📌 Stemming doesn't always give a **real word**, but it reduces words to a **common base** to help models treat them as **same conceptually**.



Why Use Stemming?

In NLP tasks like:

- Sentiment analysis
- Text classification
- Search engines

You want "played", "playing", "plays" all treated as the **same word** — "play".



1. Porter Stemmer (1979) – *Most Popular*



Idea:

- Uses **rule-based suffix stripping**
- Applies **5 steps** of rules in **sequence**, each with many conditions



Working (Mathematical Logic):

- Each word is checked against a list of suffix patterns (like **ing**, **ed**, **es**, **ly**, etc.)
- Before stripping, it **calculates a "measure" m** of the word:

$m = \text{number of vowel-consonant sequences in the stem}$

Example:

- Word: "consign"
- Vowel-consonant form: c(o)n(s)(i)g(n) → **VCVC** → $m = 2$

Only if $m > 0$, certain suffixes are removed.

Example Flow:

Word: **"relational"**

1. Rule: **"ational" → "ate"**
2. **relational → relate**

Word: **"ponies"**

- **"ies" → "i" → poni**

Word: **"caresses"**

- **"sses" → "ss" → caress**

Strength:

- Well-tested, robust
- Language-specific tuning

Weakness:

- Not aggressive
- Output not always a real word (e.g., **"relational" → "relate"**, **"ponies" → "poni"**)
- ```
from nltk.stem import PorterStemmer
```
- ```
ps = PorterStemmer()
```
- ```
words = ["playing", "played", "plays", "player", "fly", "flies"]
```
- ```
stems = [ps.stem(word) for word in words]
```
- ```
print(stems)
```
- # Output: ['play', 'play', 'play', 'player', 'fli', 'fli']

## 2. Lancaster Stemmer – *Very Aggressive*

### Idea:

- Shortens words **brutally** — even at the cost of real meaning
- Uses a **recursive algorithm**:

Keep applying rules until nothing more can be removed.

### Working:

- Removes common suffixes like **ing**, **ed**, **ly**, **tion**, etc.
- Doesn't use **measure** like Porter
- Built-in **rule list** is shorter and harsher

---

### Example:

Word: "maximum" → "max"

Word: "crying" → "cry"

Word: "happiness" → "happy" → "hap" (too much!)

python

Copy code

```
from nltk.stem import LancasterStemmer
stemmer = LancasterStemmer()
print(stemmer.stem("happiness")) # → hap
```

---

### Strength:

- Fast and very compact words
- Good for **IR (Information Retrieval)**

## ❌ Weakness:

- **Over-stemming:** "university" → "univers" → "univer"
  - **Bad for ML tasks** needing context
- 

## ❄️ 3. Snowball Stemmer – *Improved Porter*

### 📌 Idea:

- Developed by **Porter himself** as an improvement
- More languages supported
- Better structured and easier to understand

### 🧠 Logic:

- Same as Porter, but rules are **organized better**
- Handles edge cases like "lying" → "lie" instead of "ly"

It still uses **m-measure** and checks **vowel/consonant sequences**, but with:

- Cleaner rule handling
  - Better suffix priority management
- 

### 🔄 Example:

Word: "generously" → "generous"

Word: "happiness" → "happi"

Word: "playing" → "play"

python

Copy code

```
from nltk.stem import SnowballStemmer
```

```
stemmer = SnowballStemmer("english")
print(stemmer.stem("generously")) # → generous
```

---

### Quick Comparison Table:

| Feature                   | Porter       | Lancaster  | Snowball    |
|---------------------------|--------------|------------|-------------|
| Aggressiveness            | Medium       | High       | Medium      |
| Real Words Returned       | Sometimes    | Rarely     | Often       |
| Recursive                 | ✗            | ✓          | ✗           |
| Multilingual              | ✗            | ✗          | ✓           |
| Based on <i>m</i> measure | ✓            | ✗          | ✓           |
| Best for                  | Academic NLP | IR / Speed | General NLP |

---

### Summary of Mathematical Intuition:

| Concept                 | Used By          | Purpose                             |
|-------------------------|------------------|-------------------------------------|
| VC (vowel-consonant)    | Porter, Snowball | To measure "word root strength"     |
| Rule matching           | All              | Apply suffix stripping rules        |
| Recursive application   | Lancaster        | Keep trimming until stable          |
| Language-specific rules | Snowball         | Better handling in multilingual NLP |

---

## What is Lemmatization?

**Lemmatization** is the process of **reducing a word to its base form (lemma)** using **linguistic knowledge** — including dictionary, grammar, and parts of speech.

- ◆ Unlike stemming, it **always returns a valid word**.
- ◆ For example:

| Word | Lemma |
|------|-------|
|------|-------|

|         |     |
|---------|-----|
| running | run |
|---------|-----|

|        |      |
|--------|------|
| better | good |
|--------|------|

|     |    |
|-----|----|
| was | be |
|-----|----|

|      |     |
|------|-----|
| cars | car |
|------|-----|

---

## Why Use Lemmatization?

Lemmatization:

- Understands grammar and context
  - Preserves **meaning**
  - Is ideal for **chatbots, question answering, summarization**, etc.
- 

## Lemmatization vs Stemming

| Feature       | Stemming             | Lemmatization        |
|---------------|----------------------|----------------------|
| Returns       | Not always real word | Always valid word    |
| Method        | Rule-based chopping  | Dictionary + grammar |
| Accuracy      | Low                  | High                 |
| Speed         | Fast                 | Slower               |
| Uses POS Tag? | ✗                    | ✓                    |

---

## How Does Lemmatization Work (Math/Logic)?

**Step-by-step:**

1. **Tokenize** the sentence
2. **POS tag** each word (Noun, Verb, Adjective...)

Use **morphological rules + dictionary** to map:

mathematica

Copy code

Word + POS  $\rightarrow$  Lemma

- 3.
- 



### Mathematical View:

Let  $w$  be the input word, and  $\text{pos}(w)$  is its POS tag.

Then,


perl

Copy code

$\text{Lemma}(w) = \underset{l \in L(w)}{\text{argmax}} P(l \mid w, \text{pos}(w))$

Where:

- $L(w)$  is the list of valid lemmas
- $P(l \mid w, \text{pos}(w))$  is the **probability** of lemma  $l$  given word and its POS

 In advanced systems (like spaCy, WordNetLemmatizer), this is done with **rule-based heuristics + statistical probability**.

---



### Example in Code (NLTK)

python

Copy code

```
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet
from nltk import pos_tag
```

```

from nltk.tokenize import word_tokenize

lemmatizer = WordNetLemmatizer()

Example sentence
sentence = "The leaves were falling faster than usual"
words = word_tokenize(sentence)
pos_tags = pos_tag(words)

Convert POS to WordNet format
def get_wordnet_pos(tag):
 if tag.startswith('J'):
 return wordnet.ADJ
 elif tag.startswith('V'):
 return wordnet.VERB
 elif tag.startswith('N'):
 return wordnet.NOUN
 elif tag.startswith('R'):
 return wordnet.ADV
 else:
 return wordnet.NOUN

lemmatized = [lemmatizer.lemmatize(word, get_wordnet_pos(pos)) for
word, pos in pos_tags]
print(lemmatized)

```

### Output:

python

Copy code

```
['The', 'leaf', 'be', 'fall', 'fast', 'than', 'usual']
```

---

## Types of Lemmatizers

| Lemmatizer                | Description                 | Library Used |
|---------------------------|-----------------------------|--------------|
| <b>WordNet Lemmatizer</b> | Based on WordNet lexical DB | nltk         |

|                         |                                      |                       |
|-------------------------|--------------------------------------|-----------------------|
| <b>spaCy Lemmatizer</b> | Rule + statistical + language model  | <code>spacy</code>    |
| <b>TextBlob</b>         | Simpler abstraction using WordNet    | <code>textblob</code> |
| <b>Stanford NLP</b>     | Machine learning-based lemmatization | Stanford CoreNLP      |

---

## Limitations of Lemmatization

- Requires correct **POS tagging**
  - Slower than stemming
  - May not support **domain-specific words** (e.g. medical, legal)
  - Doesn't handle **typos** or slang (for that: spelling correction + BERT)
- 

## When to Use:

Use **lemmatization over stemming** when:

- You need **correct base forms** (e.g. chatbots, summaries, answers)
- Your task depends on **word meaning**
- You're handling **grammar-sensitive tasks**