

```
#Install & Import Required Libraries
# Install required libraries
!pip install nltk spacy

# Download spaCy English model
!python -m spacy download en_core_web_sm

Requirement already satisfied: nltk in /usr/local/lib/python3.12/dist-packages (3.9.1)
Requirement already satisfied: spacy in /usr/local/lib/python3.12/dist-packages (3.8.11)
Requirement already satisfied: click in /usr/local/lib/python3.12/dist-packages (from nltk) (8.3.1)
Requirement already satisfied: joblib in /usr/local/lib/python3.12/dist-packages (from nltk) (1.5.3)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.12/dist-packages (from nltk) (2025.11.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.12/dist-packages (from nltk) (4.67.1)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.12/dist-packages (from spacy) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (1.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (1.0.15)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.12/dist-packages (from spacy) (2.0.13)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.12/dist-packages (from spacy) (3.0.12)
Requirement already satisfied: thinc<8.4.0,>=8.3.4 in /usr/local/lib/python3.12/dist-packages (from spacy) (8.3.10)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.12/dist-packages (from spacy) (1.1.3)
Requirement already satisfied: srslr<3.0.0,>=2.4.3 in /usr/local/lib/python3.12/dist-packages (from spacy) (2.5.2)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.12/dist-packages (from spacy) (2.0.10)
Requirement already satisfied: weasel<0.5.0,>=0.4.2 in /usr/local/lib/python3.12/dist-packages (from spacy) (0.4.3)
Requirement already satisfied: typer-slim<1.0.0,>=0.3.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (0.20.0)
Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (2.0.2)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (2.32.4)
Requirement already satisfied: pydantic!=1.8,!=1.8.1,<3.0.0,>=1.7.4 in /usr/local/lib/python3.12/dist-packages (from spacy)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.12/dist-packages (from spacy) (3.1.6)
Requirement already satisfied: setuptools in /usr/local/lib/python3.12/dist-packages (from spacy) (75.2.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.12/dist-packages (from spacy) (25.0)
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.12/dist-packages (from pydantic!=1.8,!=1.8.1)
Requirement already satisfied: pydantic-core==2.41.4 in /usr/local/lib/python3.12/dist-packages (from pydantic!=1.8,!=1.8.1)
Requirement already satisfied: typing-extensions>=4.14.1 in /usr/local/lib/python3.12/dist-packages (from pydantic!=1.8,!=1.8.1)
Requirement already satisfied: typing-inspection>=0.4.2 in /usr/local/lib/python3.12/dist-packages (from pydantic!=1.8,!=1.8.1)
Requirement already satisfied: charset_normalizer<4,>=2 in /usr/local/lib/python3.12/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.12/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.12/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.12/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
Requirement already satisfied: blis<1.4.0,>=1.3.0 in /usr/local/lib/python3.12/dist-packages (from thinc<8.4.0,>=8.3.4->spacy)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.12/dist-packages (from thinc<8.4.0,>=8.3.4->spacy)
Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in /usr/local/lib/python3.12/dist-packages (from weasel<0.5.0,>=0.4.2)
Requirement already satisfied: smart-open<8.0.0,>=5.2.1 in /usr/local/lib/python3.12/dist-packages (from weasel<0.5.0,>=0.4.2)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.12/dist-packages (from jinja2->spacy) (3.0.3)
Requirement already satisfied: wrapt in /usr/local/lib/python3.12/dist-packages (from smart-open<8.0.0,>=5.2.1->weasel<0.5.0)
Collecting en-core-web-sm==3.8.0
  Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.8.0/en_core_web_sm-3.8.0-py3-none-any.whl (12.8/12.8 MB 79.8 MB/s eta 0:00:00)
✓ Download and installation successful
You can now load the package via spacy.load('en_core_web_sm')
⚠ Restart to reload dependencies
If you are in a Jupyter or Colab notebook, you may need to restart Python in
order to load all the package's dependencies. You can do this by selecting the
'Restart kernel' or 'Restart runtime' option.
```

```
#Import Libraries
import nltk
import spacy
from nltk.tokenize import sent_tokenize, word_tokenize
from nltk.stem import PorterStemmer
```

```
medical_text = """
Patients with diabetes mellitus often experience increased blood glucose levels.
The treatment includes insulin therapy and lifestyle modifications.
Hypertension and cardiovascular diseases are common comorbidities in diabetic patients.
"""
```

```
#Sentence & Word Tokenization (NLTK)
import nltk
nltk.download('punkt_tab')

# Sentence Tokenization
sentences = sent_tokenize(medical_text)
print("Sentences:")
for s in sentences:
    print("-", s)

Sentences:
-
Patients with diabetes mellitus often experience increased blood glucose levels.
```

```
- The treatment includes insulin therapy and lifestyle modifications.
- Hypertension and cardiovascular diseases are common comorbidities in diabetic patients.
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt_tab.zip.
```

```
# Word Tokenization
words = word_tokenize(medical_text)
print("\nWords:")
print(words)
```

Words:
['Patients', 'with', 'diabetes', 'mellitus', 'often', 'experience', 'increased', 'blood', 'glucose', 'levels', '.', 'The', '.

```
#Tokenization using spaCy
# Load spaCy model
nlp = spacy.load("en_core_web_sm")

doc = nlp(medical_text)

print("spaCy Tokens:")
for token in doc:
    print(token.text)
```

spaCy Tokens:

Patients
with
diabetes
mellitus
often
experience
increased
blood
glucose
levels
.

The
treatment
includes
insulin
therapy
and
lifestyle
modifications
.

Hypertension
and
cardiovascular
diseases
are
common
comorbidities
in
diabetic
patients
.

```
#Stemming
stemmer = PorterStemmer()

stemmed_words = [stemmer.stem(word) for word in words if word.isalpha()]

print("Stemmed Words:")
print(stemmed_words)
```

Stemmed Words:
['patient', 'with', 'diabet', 'mellitu', 'often', 'experi', 'increas', 'blood', 'glucos', 'level', 'the', 'treatment', 'incl']

```
#Lemmatization
lemmatized_words = [token.lemma_ for token in doc if token.is_alpha]

print("Lemmatized Words:")
```

```
print(lemmatized_words)
```

Lemmatized Words:
['patient', 'with', 'diabete', 'mellitus', 'often', 'experience', 'increase', 'blood', 'glucose', 'level', 'the', 'treatment']

```
#Comparison: Stemming vs Lemmatization
comparison = list(zip(stemmed_words, lemmatized_words[:len(stemmed_words)]))
```

```
print("Stemming vs Lemmatization:")
for s, l in comparison:
    print(f"{s} --> {l}")
```

```
Stemming vs Lemmatization:
patient --> patient
with --> with
diabet --> diabete
mellitu --> mellitus
often --> often
experi --> experience
increas --> increase
blood --> blood
glucos --> glucose
level --> level
the --> the
treatment --> treatment
includ --> include
insulin --> insulin
therapi --> therapy
and --> and
lifestyl --> lifestyle
modif --> modification
hypertens --> Hypertension
and --> and
cardiovascular --> cardiovascular
diseas --> disease
are --> be
common --> common
comorbid --> comorbiditie
in --> in
diabet --> diabetic
patient --> patient
```

PDF Question

```
import nltk
nltk.download('punkt')
nltk.download('wordnet')

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
True
```

```
import nltk
nltk.download('punkt_tab')
nltk.download('punkt')
nltk.download('wordnet')

from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer, WordNetLemmatizer

sentence = "NLP models are transforming the world rapidly!"

tokens = word_tokenize(sentence)
stemmed = [PorterStemmer().stem(w) for w in tokens]
lemmatized = [WordNetLemmatizer().lemmatize(w) for w in tokens]

print(tokens)
print(stemmed)
print(lemmatized)

[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt_tab.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]  Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
['NLP', 'models', 'are', 'transforming', 'the', 'world', 'rapidly', '!']
['nlp', 'model', 'are', 'transform', 'the', 'world', 'rapidli', '!']
['NLP', 'model', 'are', 'transforming', 'the', 'world', 'rapidly', '!']
```

CLASS ASSIGNMENT - SR UNIVERSITY**TOKENIZING**

```
SRUniversity="""
The SR University campus is located in Ananthasagar village of Hasanparthy Mandal in Warangal, Telangana, India.
It is in 150 acres, with both separate hostel facilities for boys and girls.
There is a huge central library along with Indias largest Technology Business Incubator (TBI) in tier 2 cities."""
```

```
import nltk
nltk.download('punkt')
from nltk.tokenize import word_tokenize
word_tokenize(SRUniversity)

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
['The',
 'SR',
 'University',
 'campus',
 'is',
 'located',
 'in',
 'Ananthasagar',
 'village',
 'of',
 'Hasanparthy',
 'Mandal',
 'in',
 'Warangal',
 ',',
 'Telangana',
 ',',
 'India',
 '',
 'It',
 'is',
 'in',
 '150',
 'acres',
 ',',
 'with',
 'both',
 'separate',
 'hostel',
 'facilities',
 'for',
 'boys',
 'and',
 'girls',
 '',
 'There',
 'is',
 'a',
 'huge',
 'central',
 'library',
 'along',
 'with',
 'Indias',
 'largest',
 'Technology',
 'Business',
 'Incubator',
 '(',
 'TBI',
 ')',
 'in',
 'tier',
 '2',
 'cities',
 '.']
```

SENTENCE TOKENIZING

```
from nltk.tokenize import sent_tokenize
sent_tokenize(SRUniversity)
```

```
['The SR University campus is located in Ananthasagar village of Hasanparthy Mandal in Warangal, Telangana, India.',
 'It is in 150 acres, with both separate hostel facilities for boys and girls.',
 'There is a huge central library along with Indias largest Technology Business Incubator (TBI) in tier 2 cities.]
```

Filtering Stop Words

```
nltk.download("stopwords")
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]  Unzipping corpora/stopwords.zip.
```

```
words_in_quote = word_tokenize(SRUniversity)
words_in_quote
```

```
['The',
 'SR',
 'University',
 'campus',
 'is',
 'located',
 'in',
 'Ananthasagar',
 'village',
 'of',
 'Hasanparthy',
 'Mandal',
 'in',
 'Warangal',
 ',',
 'Telangana',
 ',',
 'India',
 '.',
 'It',
 'is',
 'in',
 '150',
 'acres',
 ',',
 'with',
 'both',
 'separate',
 'hostel',
 'facilities',
 'for',
 'boys',
 'and',
 'girls',
 '',
 'There',
 'is',
 'a',
 'huge',
 'central',
 'library',
 'along',
 'with',
 'Indias',
 'largest',
 'Technology',
 'Business',
 'Incubator',
 '(',
 ')',
 'TBI',
 ')',
 'in',
 'tier',
 '2',
 'cities',
 '..']
```

```
stop_words = set(stopwords.words("english"))
filtered_list = []
for word in words_in_quote:
    if word.casfold() not in stop_words:
        filtered_list.append(word)
filtered_list
```

```
['SR',
 'University',
 'campus',
 'located',
 'Ananthasagar',
 'village',
 'Hasanparthy',
 'Mandal',
 'Warangal',
 '',
 'Telangana',
 ',']
```

```
'India',
'',
'150',
'acres',
'',',
'separate',
'hostel',
'facilities',
'boys',
'girls',
'.',
'huge',
'central',
'library',
'along',
'Indias',
'largest',
'Technology',
'Business',
'Incubator',
'(',
'TBI',
')',
'tier',
'',
'cities',
'..']
```

Stemming

```
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
stemmer = PorterStemmer()
words = word_tokenize(SRUniversity)
stemmed_words = [stemmer.stem(word) for word in words]
stemmed_words
```

```
['the',
'sr',
'univers',
'campu',
'is',
'locat',
'in',
'ananthasagar',
'vellag',
'of',
'hasanparthi',
'mandal',
'in',
'warang',
',',
'telangana',
',',
'india',
'..',
'it',
'is',
'in',
'150',
'acr',
',',
'with',
'both',
'sepan',
'hostel',
'facil',
'for',
'boy',
'and',
'girl',
'.',
'there',
'is',
'a',
'huge',
'central',
'librari',
'along',
'with',
'india',
'largest',
'technolog',
'busi',
'incub',
'(',
'tbi',
')',
```

```
'in',
'tier',
'2',
'citi',
'..']
```

Lemmatization

```
nltk.download('omw-1.4')
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
words = word_tokenize(SRUUniversity)
for word in words:
    print(word,"--->", lemmatizer.lemmatize(word))

[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
The ---> The
SR ---> SR
University ---> University
campus ---> campus
is ---> is
located ---> located
in ---> in
Ananthasagar ---> Ananthasagar
village ---> village
of ---> of
Hasanparthy ---> Hasanparthy
Mandal ---> Mandal
in ---> in
Warangal ---> Warangal
, ---> ,
Telangana ---> Telangana
, ---> ,
India ---> India
. ---> .
It ---> It
is ---> is
in ---> in
150 ---> 150
acres ---> acre
, ---> ,
with ---> with
both ---> both
separate ---> separate
hostel ---> hostel
facilities ---> facility
for ---> for
boys ---> boy
and ---> and
girls ---> girl
. ---> .
There ---> There
is ---> is
a ---> a
huge ---> huge
central ---> central
library ---> library
along ---> along
with ---> with
Indias ---> Indias
largest ---> largest
Technology ---> Technology
Business ---> Business
Incubator ---> Incubator
( ---> (
TBI ---> TBI
) ---> )
in ---> in
tier ---> tier
2 ---> 2
cities ---> city
. ---> .
```

```
lemmatizer.lemmatize("worst")
```

```
'worst'
```

```
lemmatizer.lemmatize("worst", pos="a")
```

```
'bad'
```

COMPARISON:Stemming VS Lemmatization

```
from nltk.stem import PorterStemmer, SnowballStemmer, LancasterStemmer, RegexpStemmer, WordNetLemmatizer
```

```
porter = PorterStemmer()
lancaster = LancasterStemmer()
snowball = SnowballStemmer(language='english')
regexp = RegexpStemmer('ing|e', min=4) # Corrected: Regex pattern on a single line
lemmatizer = WordNetLemmatizer()

word_list = ["friend", "friendship", "friends", "friendships"]
print("{0:20}{1:20}{2:20}{3:30}{4:40}{5:50}".format("Word","Porter Stemmer","Snowball Stemmer","Lancaster Stemmer",'Regexp Stemmer'))
for word in word_list:
    print("{0:20}{1:20}{2:20}{3:30}{4:40}{5:50}".format(word,porter.stem(word),snowball.stem(word),lancaster.stem(word),rege
```

Word	Porter Stemmer	Snowball Stemmer	Lancaster Stemmer	Regexp Stemmer
friend	friend	friend	friend	frind
friendship	friendship	friendship	friend	frindship
friends	friend	friend	friend	frinds
friendships	friendship	friendship	friend	frindships