

```
In [2]: import json
import pandas as pd
import numpy as np

# Load JSON file
with open("legal_advice_india_all.json", "r", encoding="utf-8") as f:
    data = json.load(f)

# Convert JSON to DataFrame
df = pd.DataFrame(data)

# Show the first few rows
print(df.head())
```

	title	author \
0	How to legally keep my father away from my wor...	divc99
1	If the person is about to get married and secu...	notms16
2	FNF not paid by startup even after legal notice	shadowslay97
3	Property Ownership Between Joint Owners: Right...	Glum_Success8717
4	Account funds put on hold Need Help	scshiv29

	url	score	created_utc \
0	https://www.reddit.com/r/LegalAdviceIndia/comm...	2	1.738060e+09
1	https://www.reddit.com/r/LegalAdviceIndia/comm...	1	1.738059e+09
2	https://www.reddit.com/r/LegalAdviceIndia/comm...	3	1.738059e+09
3	https://www.reddit.com/r/LegalAdviceIndia/comm...	2	1.738059e+09
4	https://www.reddit.com/r/LegalAdviceIndia/comm...	1	1.738058e+09

	created_date	num_comments \
0	2025-01-28 10:24:09	1
1	2025-01-28 10:10:48	1
2	2025-01-28 10:06:03	2
3	2025-01-28 10:04:55	0
4	2025-01-28 09:57:25	1

	selftext	id \
0	My father is an alcoholic and very controlling...	1ibyjsx2
1	If the person is about to get married and secu...	1bydlb
2	I worked in a very early stage startup for ove...	1bybfp
3	Note: This query is not related to me, I am as...	1byaw0
4	Some of my money (45195) in the account was pu...	1by789

	subreddit	location
0	LegalAdviceIndia	Unknown
1	LegalAdviceIndia	Unknown
2	LegalAdviceIndia	Unknown
3	LegalAdviceIndia	Unknown
4	LegalAdviceIndia	Gurgaon

```
In [3]: df["created_date"] = pd.to_datetime(df["created_utc"], unit="s")
```

```
In [4]: df["location"] = df["location"].replace("Unknown", np.nan)
```

```
In [5]: df.drop_duplicates(subset=["id"], keep="first", inplace=True)
```

```
In [6]: labour_keywords = ["labour", "employee", "work", "salary", "job", "wages", "term"]
df_labour = df[df["title"].str.contains("|".join(labour_keywords), case=False), n
```

```
In [7]: df_cleaned = df[["title", "author", "created_date", "score", "num_comments", "se
```

```
In [8]: df_cleaned = df[["title", "author", "created_date", "score", "num_comments", "se
```

```
In [9]: import re
import nltk
import pandas as pd
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer

# Download necessary data for NLTK
nltk.download("punkt")
nltk.download("stopwords")
nltk.download("wordnet")
nltk.download('all')

# Initialize Lemmatizer
lemmatizer = WordNetLemmatizer()
stop_words = set(stopwords.words("english"))
```

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading collection 'all'
[nltk_data] |
[nltk_data] | Downloading package abc to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package abc is already up-to-date!
[nltk_data] | Downloading package alpino to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package alpino is already up-to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger is already up-
[nltk_data] | to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger_eng to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger_eng is already
[nltk_data] | up-to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger_ru to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger_ru is already
[nltk_data] | up-to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger_rus to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger_rus is already
[nltk_data] | up-to-date!
[nltk_data] | Downloading package basque_grammars to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package basque_grammars is already up-to-date!
[nltk_data] | Downloading package bcp47 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package bcp47 is already up-to-date!
[nltk_data] | Downloading package biocreative_ppi to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package biocreative_ppi is already up-to-date!
[nltk_data] | Downloading package bllip_wsj_no_aux to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package bllip_wsj_no_aux is already up-to-date!
[nltk_data] | Downloading package book_grammars to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package book_grammars is already up-to-date!
[nltk_data] | Downloading package brown to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package brown is already up-to-date!
[nltk_data] | Downloading package brown_tei to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package brown_tei is already up-to-date!
[nltk_data] | Downloading package cess_cat to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package cess_cat is already up-to-date!
[nltk_data] | Downloading package cess_esp to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package cess_esp is already up-to-date!
```

```
[nltk_data] | Downloading package chat80 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package chat80 is already up-to-date!
[nltk_data] | Downloading package city_database to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package city_database is already up-to-date!
[nltk_data] | Downloading package cmudict to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package cmudict is already up-to-date!
[nltk_data] | Downloading package comparative_sentences to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package comparative_sentences is already up-to-
[nltk_data] | date!
[nltk_data] | Downloading package comtrans to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package comtrans is already up-to-date!
[nltk_data] | Downloading package conll2000 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2000 is already up-to-date!
[nltk_data] | Downloading package conll2002 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2002 is already up-to-date!
[nltk_data] | Downloading package conll2007 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2007 is already up-to-date!
[nltk_data] | Downloading package crubadan to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package crubadan is already up-to-date!
[nltk_data] | Downloading package dependency_treebank to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package dependency_treebank is already up-to-date!
[nltk_data] | Downloading package dolch to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package dolch is already up-to-date!
[nltk_data] | Downloading package europarl_raw to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package europarl_raw is already up-to-date!
[nltk_data] | Downloading package extended_omw to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package extended_omw is already up-to-date!
[nltk_data] | Downloading package floresta to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package floresta is already up-to-date!
[nltk_data] | Downloading package framenet_v15 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package framenet_v15 is already up-to-date!
[nltk_data] | Downloading package framenet_v17 to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package framenet_v17 is already up-to-date!
[nltk_data] | Downloading package gazetteers to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package gazetteers is already up-to-date!
[nltk_data] | Downloading package genesis to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package genesis is already up-to-date!
[nltk_data] | Downloading package gutenberg to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package gutenberg is already up-to-date!
[nltk_data] | Downloading package ieer to
[nltk_data] |   C:\Users\Komal\AppData\Roaming\nltk_data...
```

```

[nltk_data] | Package ieer is already up-to-date!
[nltk_data] | Downloading package inaugural to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package inaugural is already up-to-date!
[nltk_data] | Downloading package indian to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package indian is already up-to-date!
[nltk_data] | Downloading package jeita to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package jeita is already up-to-date!
[nltk_data] | Downloading package kimmo to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package kimmo is already up-to-date!
[nltk_data] | Downloading package knbc to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package knbc is already up-to-date!
[nltk_data] | Downloading package large_grammars to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package large_grammars is already up-to-date!
[nltk_data] | Downloading package lin_thesaurus to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package lin_thesaurus is already up-to-date!
[nltk_data] | Downloading package mac_morpho to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package mac_morpho is already up-to-date!
[nltk_data] | Downloading package machado to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package machado is already up-to-date!
[nltk_data] | Downloading package masc_tagged to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package masc_tagged is already up-to-date!
[nltk_data] | Downloading package maxent_ne_chunker to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_ne_chunker is already up-to-date!
[nltk_data] | Downloading package maxent_ne_chunker_tab to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_ne_chunker_tab is already up-to-
[nltk_data] | date!
[nltk_data] | Downloading package maxent_treebank_pos_tagger to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_treebank_pos_tagger is already up-
[nltk_data] | to-date!
[nltk_data] | Downloading package maxent_treebank_pos_tagger_tab to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_treebank_pos_tagger_tab is already
[nltk_data] | up-to-date!
[nltk_data] | Downloading package moses_sample to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package moses_sample is already up-to-date!
[nltk_data] | Downloading package movie_reviews to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package movie_reviews is already up-to-date!
[nltk_data] | Downloading package mte_teip5 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package mte_teip5 is already up-to-date!
[nltk_data] | Downloading package mwa_ppdb to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package mwa_ppdb is already up-to-date!
[nltk_data] | Downloading package names to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...

```

```

[nltk_data] | Package names is already up-to-date!
[nltk_data] | Downloading package nombank.1.0 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package nombank.1.0 is already up-to-date!
[nltk_data] | Downloading package nonbreaking_prefixes to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package nonbreaking_prefixes is already up-to-date!
[nltk_data] | Downloading package nps_chat to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package nps_chat is already up-to-date!
[nltk_data] | Downloading package omw to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package omw is already up-to-date!
[nltk_data] | Downloading package omw-1.4 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package omw-1.4 is already up-to-date!
[nltk_data] | Downloading package opinion_lexicon to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package opinion_lexicon is already up-to-date!
[nltk_data] | Downloading package panlex_swadesh to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package panlex_swadesh is already up-to-date!
[nltk_data] | Downloading package paradigms to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package paradigms is already up-to-date!
[nltk_data] | Downloading package pe08 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package pe08 is already up-to-date!
[nltk_data] | Downloading package perluniprops to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package perluniprops is already up-to-date!
[nltk_data] | Downloading package pil to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package pil is already up-to-date!
[nltk_data] | Downloading package pl196x to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package pl196x is already up-to-date!
[nltk_data] | Downloading package porter_test to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package porter_test is already up-to-date!
[nltk_data] | Downloading package ppattach to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package ppattach is already up-to-date!
[nltk_data] | Downloading package problem_reports to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package problem_reports is already up-to-date!
[nltk_data] | Downloading package product_reviews_1 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package product_reviews_1 is already up-to-date!
[nltk_data] | Downloading package product_reviews_2 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package product_reviews_2 is already up-to-date!
[nltk_data] | Downloading package propbank to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package propbank is already up-to-date!
[nltk_data] | Downloading package pros_cons to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package pros_cons is already up-to-date!
[nltk_data] | Downloading package ptb to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...

```

```
[nltk_data] | Package ptb is already up-to-date!
[nltk_data] | Downloading package punkt to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package punkt is already up-to-date!
[nltk_data] | Downloading package punkt_tab to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package punkt_tab is already up-to-date!
[nltk_data] | Downloading package qc to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package qc is already up-to-date!
[nltk_data] | Downloading package reuters to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package reuters is already up-to-date!
[nltk_data] | Downloading package rslp to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package rslp is already up-to-date!
[nltk_data] | Downloading package rte to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package rte is already up-to-date!
[nltk_data] | Downloading package sample_grammars to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package sample_grammars is already up-to-date!
[nltk_data] | Downloading package semcor to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package semcor is already up-to-date!
[nltk_data] | Downloading package senseval to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package senseval is already up-to-date!
[nltk_data] | Downloading package sentence_polarity to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package sentence_polarity is already up-to-date!
[nltk_data] | Downloading package sentiwordnet to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package sentiwordnet is already up-to-date!
[nltk_data] | Downloading package shakespeare to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package shakespeare is already up-to-date!
[nltk_data] | Downloading package sinica_treebank to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package sinica_treebank is already up-to-date!
[nltk_data] | Downloading package smultron to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package smultron is already up-to-date!
[nltk_data] | Downloading package snowball_data to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package snowball_data is already up-to-date!
[nltk_data] | Downloading package spanish_grammars to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package spanish_grammars is already up-to-date!
[nltk_data] | Downloading package state_union to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package state_union is already up-to-date!
[nltk_data] | Downloading package stopwords to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package stopwords is already up-to-date!
[nltk_data] | Downloading package subjectivity to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package subjectivity is already up-to-date!
[nltk_data] | Downloading package swadesh to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
```



```
[nltk_data] | Package swadesh is already up-to-date!
[nltk_data] | Downloading package switchboard to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package switchboard is already up-to-date!
[nltk_data] | Downloading package tagsets to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package tagsets is already up-to-date!
[nltk_data] | Downloading package tagsets_json to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package tagsets_json is already up-to-date!
[nltk_data] | Downloading package timit to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package timit is already up-to-date!
[nltk_data] | Downloading package toolbox to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package toolbox is already up-to-date!
[nltk_data] | Downloading package treebank to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package treebank is already up-to-date!
[nltk_data] | Downloading package twitter_samples to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package twitter_samples is already up-to-date!
[nltk_data] | Downloading package udhr to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package udhr is already up-to-date!
[nltk_data] | Downloading package udhr2 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package udhr2 is already up-to-date!
[nltk_data] | Downloading package unicode_samples to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package unicode_samples is already up-to-date!
[nltk_data] | Downloading package universal_tagset to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package universal_tagset is already up-to-date!
[nltk_data] | Downloading package universal_treebanks_v20 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package universal_treebanks_v20 is already up-to-
[nltk_data] | date!
[nltk_data] | Downloading package vader_lexicon to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package vader_lexicon is already up-to-date!
[nltk_data] | Downloading package verbnet to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package verbnet is already up-to-date!
[nltk_data] | Downloading package verbnet3 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package verbnet3 is already up-to-date!
[nltk_data] | Downloading package webtext to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package webtext is already up-to-date!
[nltk_data] | Downloading package wmt15_eval to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wmt15_eval is already up-to-date!
[nltk_data] | Downloading package word2vec_sample to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package word2vec_sample is already up-to-date!
[nltk_data] | Downloading package wordnet to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet is already up-to-date!
[nltk_data] | Downloading package wordnet2021 to
```



```

[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet2021 is already up-to-date!
[nltk_data] | Downloading package wordnet2022 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet2022 is already up-to-date!
[nltk_data] | Downloading package wordnet31 to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet31 is already up-to-date!
[nltk_data] | Downloading package wordnet_ic to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet_ic is already up-to-date!
[nltk_data] | Downloading package words to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package words is already up-to-date!
[nltk_data] | Downloading package ycoe to
[nltk_data] | C:\Users\Komal\AppData\Roaming\nltk_data...
[nltk_data] | Package ycoe is already up-to-date!
[nltk_data] |
[nltk_data] Done downloading collection all

```

```

In [10]: def preprocess_text(text):
          if pd.isna(text): # Handle missing values
              return ""

          text = text.lower() # Convert to Lowercase
          text = re.sub(r"http\S+|www\S+", "", text) # Remove URLs
          text = re.sub(r"^[a-zA-Z\s]", "", text) # Remove special characters & punct
          words = word_tokenize(text) # Tokenization
          words = [word for word in words if word not in stop_words] # Remove stopwords
          words = [lemmatizer.lemmatize(word) for word in words] # Lemmatization

          return " ".join(words)

```

```

In [11]: df_cleaned["title"] = df_cleaned["title"].apply(preprocess_text)
          df_cleaned["selftext"] = df_cleaned["selftext"].apply(preprocess_text)

```

C:\Users\Komal\AppData\Local\Temp\ipykernel_6644\1990440792.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df_cleaned["title"] = df_cleaned["title"].apply(preprocess_text)
```

C:\Users\Komal\AppData\Local\Temp\ipykernel_6644\1990440792.py:2: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df_cleaned["selftext"] = df_cleaned["selftext"].apply(preprocess_text)
```

```

In [15]: df_cleaned.to_json("preprocessed_legal_advice.json", orient="records", indent=4)
          print("Preprocessed data saved successfully!")

```

Preprocessed data saved successfully!

```

In [16]: 5

```

Out[16]: 5

```
In [13]: import mysql.connector

df_cleaned = df_cleaned.replace({np.nan: None})

# Connect to MySQL database
conn = mysql.connector.connect(
    host="localhost",
    user="root",
    password="",
    database="sma_reddit"
)

cursor = conn.cursor()

# Create table if it doesn't exist
cursor.execute('''
CREATE TABLE IF NOT EXISTS legal_data (
    id INT AUTO_INCREMENT PRIMARY KEY,
    title VARCHAR(255),
    author VARCHAR(255),
    url VARCHAR(255),
    score INT,
    created_date DATETIME,
    num_comments INT,
    selftext TEXT,
    location VARCHAR(255)
)
''')

for index, row in df_cleaned.iterrows():
    cursor.execute('''
        INSERT INTO legal_data (title, author, url, score, created_date, num_com
        VALUES (%s, %s, %s, %s, %s, %s, %s, %s)
        ''', (
            row['title'],
            row['author'],
            row['url'],
            row['score'],
            row['created_date'],
            row['num_comments'],
            row['selftext'],
            row['location']
        ))

conn.commit()
conn.close()
```

In []: