

```
In [111]: ┏ import re
      import ast
      import pickle

      import numpy as np
      import pandas as pd
      import streamlit as st

      import nltk
      from nltk.corpus import stopwords
      from nltk.stem.porter import PorterStemmer

      import sklearn
      from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
      from sklearn.metrics.pairwise import cosine_similarity

      import warnings
      warnings.filterwarnings('ignore')
```

```
In [ ]: ┏
```

```
In [2]: ┏ movies = pd.read_csv('tmdb_5000_movies.csv')
      credits = pd.read_csv('tmdb_5000_credits.csv')
```

In [3]: movies.head(3)

		budget	genres	homepage	id	keywords	original_language	original_title	overview
0	237000000		[{"id": 28, "name": "Action"}, {"id": 12, "nam...]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": ...]	en	Avatar	In the 22nd century, a paraplegic Marine is di...
1	300000000		[{"id": 12, "name": "Adventure"}, {"id": 14, "...]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "na...]	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...
2	245000000		[{"id": 28, "name": "Action"}, {"id": 12, "nam...]	http://www.sonypictures.com/movies/spectre/	206647	[{"id": 470, "name": "spy"}, {"id": 818, "name...]	en	Spectre	A cryptic message from Bond's past sends him o...

In []:

In [4]: ► credits.head()

Out[4]:	movie_id	title	cast	crew
0	19995	Avatar	[{"cast_id": 242, "character": "Jake Sully", ...}	[{"credit_id": "52fe48009251416c750aca23", "de...}
1	285	Pirates of the Caribbean: At World's End	[{"cast_id": 4, "character": "Captain Jack Spa...}	[{"credit_id": "52fe4232c3a36847f800b579", "de...}
2	206647	Spectre	[{"cast_id": 1, "character": "James Bond", "cr...}	[{"credit_id": "54805967c3a36829b5002c41", "de...}
3	49026	The Dark Knight Rises	[{"cast_id": 2, "character": "Bruce Wayne / Ba...}	[{"credit_id": "52fe4781c3a36847f81398c3", "de...}
4	49529	John Carter	[{"cast_id": 5, "character": "John Carter", "c...}	[{"credit_id": "52fe479ac3a36847f813eaa3", "de...}

In []: ➤

In [5]: ► movies.shape

Out[5]: (4803, 20)

In [6]: ► credits.shape

Out[6]: (4803, 4)

In []:

In [7]: ► movies.columns

```
Out[7]: Index(['budget', 'genres', 'homepage', 'id', 'keywords', 'original_language',
       'original_title', 'overview', 'popularity', 'production_companies',
       'production_countries', 'release_date', 'revenue', 'runtime',
       'spoken_languages', 'status', 'tagline', 'title', 'vote_average',
       'vote_count'],
      dtype='object')
```

```
In [8]: ┏ credits.columns
```

```
Out[8]: Index(['movie_id', 'title', 'cast', 'crew'], dtype='object')
```

```
In [9]: ┏ movies['title']
```

```
Out[9]: 0          Avatar
1  Pirates of the Caribbean: At World's End
2          Spectre
3          The Dark Knight Rises
4          John Carter
...
4798          El Mariachi
4799          Newlyweds
4800          Signed, Sealed, Delivered
4801          Shanghai Calling
4802          My Date with Drew
Name: title, Length: 4803, dtype: object
```

```
In [10]: ┏ credits['title']
```

```
Out[10]: 0          Avatar
1  Pirates of the Caribbean: At World's End
2          Spectre
3          The Dark Knight Rises
4          John Carter
...
4798          El Mariachi
4799          Newlyweds
4800          Signed, Sealed, Delivered
4801          Shanghai Calling
4802          My Date with Drew
Name: title, Length: 4803, dtype: object
```

```
In [ ]: ┏
```

Merge the dfs

In [11]: ⏷ movies = movies.merge(credits, on='title')

In [12]: ⏷ movies.head(2)

Out[12]:

	budget	genres	homepage	id	keywords	original_language	original_title	overview	popularity
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name..."]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": ...]	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150
1	300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "..."]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "na..."]	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	139

2 rows × 23 columns

In []: ⏷

Choosing the relevant features for movie recommendation

In [13]: movies[:2]

Out[13]:

	budget	genres	homepage	id	keywords	original_language	original_title	overview	popularity
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name..."}]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": ...}]]	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150
1	300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "..."}]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "na..."}]]	en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	139

2 rows × 23 columns

Final Features

- movie_id
- title
- overview
- genres
- keywords
- cast
- crew

In [14]: df = movies[['movie_id', 'title', 'overview', 'genres', 'keywords', 'cast', 'crew']]
df.head(3)

Out[14]:

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[{"id": 28, "name": "Action"}, {"id": 12, "nam...}	[{"id": 1463, "name": "culture clash"}, {"id": ...]	[{"cast_id": 242, "character": "Jake Sully", "...]	[{"credit_id": "52fe48009251416c750aca23", "de...]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[{"id": 12, "name": "Adventure"}, {"id": 14, "...]	[{"id": 270, "name": "ocean"}, {"id": 726, "na...]	[{"cast_id": 4, "character": "Captain Jack Spa...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[{"id": 28, "name": "Action"}, {"id": 12, "nam...}	[{"id": 470, "name": "spy"}, {"id": 818, "name...]	[{"cast_id": 1, "character": "James Bond", "cr...]	[{"credit_id": "54805967c3a36829b5002c41", "de...]

In []:

Goal -

- movie_id + title + tags

In [15]: df['overview']

Out[15]: 0 In the 22nd century, a paraplegic Marine is di...
1 Captain Barbossa, long believed to be dead, ha...
2 A cryptic message from Bond's past sends him o...
3 Following the death of District Attorney Harve...
4 John Carter is a war-weary, former military ca...
...
4804 El Mariachi just wants to play his guitar and ...
4805 A newlywed couple's honeymoon is upended by th...
4806 "Signed, Sealed, Delivered" introduces a dedic...
4807 When ambitious New York attorney Sam is sent t...
4808 Ever since the second grade when he first saw ...
Name: overview, Length: 4809, dtype: object

```
In [16]: df['overview'][0]
```

```
Out[16]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization.'
```

```
In [ ]:
```

```
In [17]: df['genres']
```

```
Out[17]: 0      [{"id": 28, "name": "Action"}, {"id": 12, "nam...
 1      [{"id": 12, "name": "Adventure"}, {"id": 14, "...
 2      [{"id": 28, "name": "Action"}, {"id": 12, "nam...
 3      [{"id": 28, "name": "Action"}, {"id": 80, "nam...
 4      [{"id": 28, "name": "Action"}, {"id": 12, "nam...
 ...
 4804     [{"id": 28, "name": "Action"}, {"id": 80, "nam...
 4805     [{"id": 35, "name": "Comedy"}, {"id": 10749, "...
 4806     [{"id": 35, "name": "Comedy"}, {"id": 18, "nam...
 4807           []
 4808           [{"id": 99, "name": "Documentary"}]
Name: genres, Length: 4809, dtype: object
```

```
In [18]: df['genres'][0]
```

```
Out[18]: '[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'
```

```
In [ ]:
```

In [19]: ► df['keywords']

```
Out[19]: 0      [{"id": 1463, "name": "culture clash"}, {"id":...  
1      [{"id": 270, "name": "ocean"}, {"id": 726, "na...  
2      [{"id": 470, "name": "spy"}, {"id": 818, "name...  
3      [{"id": 849, "name": "dc comics"}, {"id": 853,...  
4      [{"id": 818, "name": "based on novel"}, {"id":...  
       ...  
4804     [{"id": 5616, "name": "united states\u2013mexi...  
4805           []  
4806     [{"id": 248, "name": "date"}, {"id": 699, "nam...  
4807           []  
4808     [{"id": 1523, "name": "obsession"}, {"id": 224...  
Name: keywords, Length: 4809, dtype: object
```

In [20]: ► df['keywords'][0]

```
Out[20]: '[{"id": 1463, "name": "culture clash"}, {"id": 2964, "name": "future"}, {"id": 3386, "name": "space wa  
r"}, {"id": 3388, "name": "space colony"}, {"id": 3679, "name": "society"}, {"id": 3801, "name": "space  
travel"}, {"id": 9685, "name": "futuristic"}, {"id": 9840, "name": "romance"}, {"id": 9882, "name": "spa  
ce"}, {"id": 9951, "name": "alien"}, {"id": 10148, "name": "tribe"}, {"id": 10158, "name": "alien plane  
t"}, {"id": 10987, "name": "cgi"}, {"id": 11399, "name": "marine"}, {"id": 13065, "name": "soldier"},  
 {"id": 14643, "name": "battle"}, {"id": 14720, "name": "love affair"}, {"id": 165431, "name": "anti wa  
r"}, {"id": 193554, "name": "power relations"}, {"id": 206690, "name": "mind and soul"}, {"id": 209714,  
"name": "3d"}]'
```

In []: ►

```
In [21]: df['cast']
```

```
Out[21]: 0      [{"cast_id": 242, "character": "Jake Sully", "...  
1      [{"cast_id": 4, "character": "Captain Jack Spa...  
2      [{"cast_id": 1, "character": "James Bond", "cr...  
3      [{"cast_id": 2, "character": "Bruce Wayne / Ba...  
4      [{"cast_id": 5, "character": "John Carter", "c...  
       ...  
4804    [{"cast_id": 1, "character": "El Mariachi", "c...  
4805    [{"cast_id": 1, "character": "Buzzy", "credit_...  
4806    [{"cast_id": 8, "character": "Oliver O\u2019To...  
4807    [{"cast_id": 3, "character": "Sam", "credit_id...  
4808    [{"cast_id": 3, "character": "Herself", "credi...  
Name: cast, Length: 4809, dtype: object
```

In [22]: ► df['cast'][0]

In []:

In [23]: df['crew']

```
Out[23]: 0      [{"credit_id": "52fe48009251416c750aca23", "de...  
1      [{"credit_id": "52fe4232c3a36847f800b579", "de...  
2      [{"credit_id": "54805967c3a36829b5002c41", "de...  
3      [{"credit_id": "52fe4781c3a36847f81398c3", "de...  
4      [{"credit_id": "52fe479ac3a36847f813eaa3", "de...  
     ...  
4804     [{"credit_id": "52fe44eec3a36847f80b280b", "de...  
4805     [{"credit_id": "52fe487dc3a368484e0fb013", "de...  
4806     [{"credit_id": "52fe4df3c3a36847f8275ecf", "de...  
4807     [{"credit_id": "52fe4ad9c3a368484e16a36b", "de...  
4808     [{"credit_id": "58ce021b9251415a390165d9", "de...  
Name: crew, Length: 4809, dtype: object
```

In [24]: df['crew'][0]

```
Out[24]: '[{"credit_id": "52fe48009251416c750aca23", "department": "Editing", "gender": 0, "id": 1721, "job": "Editor", "name": "Stephen E. Rivkin"}, {"credit_id": "539c47ecc3a36810e3001f87", "department": "Art", "gender": 2, "id": 496, "job": "Production Design", "name": "Rick Carter"}, {"credit_id": "54491c89c3a3680fb4001cf7", "department": "Sound", "gender": 0, "id": 900, "job": "Sound Designer", "name": "Christopher Boyes"}, {"credit_id": "54491cb70e0a267480001bd0", "department": "Sound", "gender": 0, "id": 900, "job": "Supervising Sound Editor", "name": "Christopher Boyes"}, {"credit_id": "539c4a4cc3a36810c9002101", "department": "Production", "gender": 1, "id": 1262, "job": "Casting", "name": "Mali Finn"}, {"credit_id": "5544ee3b925141499f0008fc", "department": "Sound", "gender": 2, "id": 1729, "job": "Original Music Composer", "name": "James Horner"}, {"credit_id": "52fe48009251416c750ac9c3", "department": "Directing", "gender": 2, "id": 2710, "job": "Director", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750ac9d9", "department": "Writing", "gender": 2, "id": 2710, "job": "Writer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca17", "department": "Editing", "gender": 2, "id": 2710, "job": "Editor", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca29", "department": "Production", "gender": 2, "id": 2710, "job": "Producer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca3f", "department": "Writing", "gender": 2, "id": 2710, "job": "Screenplay", "name": "James Cameron"}, {"credit_id": "539c4987c3a36810ba0021a4", "department": "Art", "gender": 2, "id": 7236, "job": "Art Direction", "name": "Andrew Menzies"}, {"credit_id": "549598c3c3a3686ae9004383", "department": "Visual Effects", "gender": 0, "id": 6690, "job": "Visual Effects Producer", "name": "Jill Brooks"}, {"credit_id": "52fe48009251416c750aca4b", "department": "Production", "gender": 1, "id": 6347, "job": "Casting", "name": "Mangan, Simkin"}, {"credit_id": "570bcf4107f1117d700775e", "department": "Sound"}]
```

In []:

In []: ┌

Cleaning the genres features

In [25]: ┌ df[:2]

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[{"id": 28, "name": "Action"}, {"id": 12, "name": ...]	[{"id": 1463, "name": "culture clash"}, {"id": ...]	[{"cast_id": 242, "character": "Jake Sully", "...]	[{"credit_id": "52fe48009251416c750aca23", "de...]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[{"id": 12, "name": "Adventure"}, {"id": 14, "name": ...]	[{"id": 270, "name": "ocean"}, {"id": 726, "na...]	[{"cast_id": 4, "character": "Captain Jack Spa...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]

In []: ┌

In [26]: ┌ df['genres']

```
Out[26]: 0      [{"id": 28, "name": "Action"}, {"id": 12, "nam...}
 1      [{"id": 12, "name": "Adventure"}, {"id": 14, "...
 2      [{"id": 28, "name": "Action"}, {"id": 12, "nam...
 3      [{"id": 28, "name": "Action"}, {"id": 80, "nam...
 4      [{"id": 28, "name": "Action"}, {"id": 12, "nam...
 ...
 4804     [{"id": 28, "name": "Action"}, {"id": 80, "nam...
 4805     [{"id": 35, "name": "Comedy"}, {"id": 10749, "...
 4806     [{"id": 35, "name": "Comedy"}, {"id": 18, "nam...
 4807          []
 4808          [{"id": 99, "name": "Documentary"}]
Name: genres, Length: 4809, dtype: object
```

In [27]: ┌ df['genres'][0]

```
Out[27]: '[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'
```

```
In [28]: ► type(df['genres'][0])
```

```
Out[28]: str
```

```
In [29]: ► ast.literal_eval(df['genres'][0])
```

```
Out[29]: [{"id": 28, "name": "Action"},  
          {"id": 12, "name": "Adventure"},  
          {"id": 14, "name": "Fantasy"},  
          {"id": 878, "name": "Science Fiction"}]
```

```
In [ ]: ►
```

```
In [30]: ► for i in ast.literal_eval(df['genres'][0]):  
           print(i['name'])
```

```
Action  
Adventure  
Fantasy  
Science Fiction
```

```
In [ ]: ►
```

```
In [31]: ► def fetch_genres(text):
```

```
    l = []  
  
    for i in ast.literal_eval(text):  
        l.append(i['name'])  
  
    return l
```

```
In [32]: ► fetch_genres(df['genres'][0])
```

```
Out[32]: ['Action', 'Adventure', 'Fantasy', 'Science Fiction']
```

In []:

In [33]:

```
Out[33]: 0      [Action, Adventure, Fantasy, Science Fiction]
1          [Adventure, Fantasy, Action]
2          [Action, Adventure, Crime]
3          [Action, Crime, Drama, Thriller]
4          [Action, Adventure, Science Fiction]
...
4804      [Action, Crime, Thriller]
4805      [Comedy, Romance]
4806      [Comedy, Drama, Romance, TV Movie]
4807      []
4808      [Documentary]
Name: genres, Length: 4809, dtype: object
```

In [34]:

df['genres'] = df['genres'].apply(fetch_genres)

In [35]:

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[Action, Adventure, Fantasy, Science Fiction]	[{"id": 1463, "name": "culture clash"}, {"id": ...]	[{"cast_id": 242, "character": "Jake Sully", "...]	[{"credit_id": "52fe48009251416c750aca23", "de...]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[Adventure, Fantasy, Action]	[{"id": 270, "name": "ocean"}, {"id": 726, "na...]	[{"cast_id": 4, "character": "Captain Jack Spa...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[Action, Adventure, Crime]	[{"id": 470, "name": "spy"}, {"id": 818, "name...]	[{"cast_id": 1, "character": "James Bond", "cr...]	[{"credit_id": "54805967c3a36829b5002c41", "de...]

In [36]: df[:3]

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[Action, Adventure, Fantasy, Science Fiction]	[{"id": 1463, "name": "culture clash"}, {"id": ...]	[{"cast_id": 242, "character": "Jake Sully", "na...]	[{"credit_id": "52fe48009251416c750aca23", "de...]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[Adventure, Fantasy, Action]	[{"id": 270, "name": "ocean"}, {"id": 726, "na...]	[{"cast_id": 4, "character": "Captain Jack Spa...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[Action, Adventure, Crime]	[{"id": 470, "name": "spy"}, {"id": 818, "name...]	[{"cast_id": 1, "character": "James Bond", "cr...]	[{"credit_id": "54805967c3a36829b5002c41", "de...]

In []:

Keywords

In [37]: df['keywords']

Out[37]:

```

0      [{"id": 1463, "name": "culture clash"}, {"id": ...}
1      [{"id": 270, "name": "ocean"}, {"id": 726, "na...}
2      [{"id": 470, "name": "spy"}, {"id": 818, "name...}
3      [{"id": 849, "name": "dc comics"}, {"id": 853, ...
4      [{"id": 818, "name": "based on novel"}, {"id": ...
...
4804     [{"id": 5616, "name": "united states\u2013mexi...
4805           []
4806     [{"id": 248, "name": "date"}, {"id": 699, "nam...
4807           []
4808     [{"id": 1523, "name": "obsession"}, {"id": 224...
Name: keywords, Length: 4809, dtype: object

```

```
In [38]: df['keywords'][0]
```

```
Out[38]: '[{"id": 1463, "name": "culture clash"}, {"id": 2964, "name": "future"}, {"id": 3386, "name": "space war"}, {"id": 3388, "name": "space colony"}, {"id": 3679, "name": "society"}, {"id": 3801, "name": "space travel"}, {"id": 9685, "name": "futuristic"}, {"id": 9840, "name": "romance"}, {"id": 9882, "name": "space"}, {"id": 9951, "name": "alien"}, {"id": 10148, "name": "tribe"}, {"id": 10158, "name": "alien planet"}, {"id": 10987, "name": "cgi"}, {"id": 11399, "name": "marine"}, {"id": 13065, "name": "soldier"}, {"id": 14643, "name": "battle"}, {"id": 14720, "name": "love affair"}, {"id": 165431, "name": "anti war"}, {"id": 193554, "name": "power relations"}, {"id": 206690, "name": "mind and soul"}, {"id": 209714, "name": "3d"}]'
```

```
In [39]: ast.literal_eval(df['keywords'][0])
```

```
Out[39]: [{"id": 1463, "name": 'culture clash'}, {"id": 2964, "name": 'future'}, {"id": 3386, "name": 'space war'}, {"id": 3388, "name": 'space colony'}, {"id": 3679, "name": 'society'}, {"id": 3801, "name": 'space travel'}, {"id": 9685, "name": 'futuristic'}, {"id": 9840, "name": 'romance'}, {"id": 9882, "name": 'space'}, {"id": 9951, "name": 'alien'}, {"id": 10148, "name": 'tribe'}, {"id": 10158, "name": 'alien planet'}, {"id": 10987, "name": 'cgi'}, {"id": 11399, "name": 'marine'}, {"id": 13065, "name": 'soldier'}, {"id": 14643, "name": 'battle'}, {"id": 14720, "name": 'love affair'}, {"id": 165431, "name": 'anti war'}, {"id": 193554, "name": 'power relations'}, {"id": 206690, "name": 'mind and soul'}]
```

```
In [ ]:
```

In [40]: ► def fetch_keywords(text):

```
    l = []
    for i in ast.literal_eval(text):
        l.append(i['name'])
    return l
```

In [41]: ► fetch_keywords(df['keywords'][0])

Out[41]: ['culture clash',
 'future',
 'space war',
 'space colony',
 'society',
 'space travel',
 'futuristic',
 'romance',
 'space',
 'alien',
 'tribe',
 'alien planet',
 'cgi',
 'marine',
 'soldier',
 'battle',
 'love affair',
 'anti war',
 'power relations',
 ' ',
 ' ']

In []: ►

In [42]: df['keywords'].apply(fetch_keywords)

```
Out[42]: 0      [culture clash, future, space war, space colon...
 1      [ocean, drug abuse, exotic island, east india ...
 2      [spy, based on novel, secret agent, sequel, mi...
 3      [dc comics, crime fighter, terrorist, secret i...
 4      [based on novel, mars, medallion, space travel...
 ...
 4804    [united states-mexico barrier, legs, arms, pap...
 4805    []
 4806    [date, love at first sight, narration, investi...
 4807    []
 4808    [obsession, camcorder, crush, dream girl]
Name: keywords, Length: 4809, dtype: object
```

In [43]: df['keywords'] = df['keywords'].apply(fetch_keywords)

In [44]: df[:3]

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...	[{"cast_id": 242, "character": "Jake Sully", ...]	[{"credit_id": "52fe48009251416c750aca23", "de...]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...]	[{"cast_id": 4, "character": "Captain Jack Spa...]	[{"credit_id": "52fe4232c3a36847f800b579", "de...]
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...]	[{"cast_id": 1, "character": "James Bond", "cr...]	[{"credit_id": "54805967c3a36829b5002c41", "de...]

In []:

In [45]: ► df.isna().sum()

```
Out[45]: movie_id    0
          title      0
         overview   3
        genres     0
      keywords    0
        cast      0
       crew      0
      dtype: int64
```

In [46]: ► df.dropna(inplace=True)

In [47]: ► df.isna().sum()

```
Out[47]: movie_id    0
          title      0
         overview   0
        genres     0
      keywords    0
        cast      0
       crew      0
      dtype: int64
```

In []: ►

Cast

In [48]: ► df['cast']

```
Out[48]: 0      [{"cast_id": 242, "character": "Jake Sully", "...  
1      [{"cast_id": 4, "character": "Captain Jack Spa...  
2      [{"cast_id": 1, "character": "James Bond", "cr...  
3      [{"cast_id": 2, "character": "Bruce Wayne / Ba...  
4      [{"cast_id": 5, "character": "John Carter", "c...  
       ...  
4804  [{"cast_id": 1, "character": "El Mariachi", "c...  
4805  [{"cast_id": 1, "character": "Buzzy", "credit_...  
4806  [{"cast_id": 8, "character": "Oliver O\u2019To...  
4807  [{"cast_id": 3, "character": "Sam", "credit_id...  
4808  [{"cast_id": 3, "character": "Herself", "credi...  
Name: cast, Length: 4806, dtype: object
```

In [49]: ► df['cast'][0]

In []:

```
In [50]: ┏ def fetch_cast(text):
```

```
    l = []
    counter = 0

    for i in ast.literal_eval(text):
        if counter != 3:
            l.append(i['name'])
            counter += 1
        else:
            break

    return l
```

```
In [51]: ┏ fetch_cast(df['cast'][0])
```

```
Out[51]: ['Sam Worthington', 'Zoe Saldana', 'Sigourney Weaver']
```

```
In [ ]: ┏
```

```
In [52]: ┏ df['cast'].apply(fetch_cast)
```

```
Out[52]: 0      [Sam Worthington, Zoe Saldana, Sigourney Weaver]
1      [Johnny Depp, Orlando Bloom, Keira Knightley]
2      [Daniel Craig, Christoph Waltz, Léa Seydoux]
3      [Christian Bale, Michael Caine, Gary Oldman]
4      [Taylor Kitsch, Lynn Collins, Samantha Morton]
...
4804     [Carlos Gallardo, Jaime de Hoyos, Peter Marqua...
4805     [Edward Burns, Kerry Bishé, Marsha Dietlein]
4806     [Eric Mabius, Kristin Booth, Crystal Lowe]
4807     [Daniel Henney, Eliza Coupe, Bill Paxton]
4808     [Drew Barrymore, Brian Herzlinger, Corey Feldman]
Name: cast, Length: 4806, dtype: object
```

```
In [53]: ┏ df['cast'] = df['cast'].apply(fetch_cast)
```

In [54]: df[:3]

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...]	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[{"credit_id": "52fe48009251416c750aca23", "de...
1	285	Pirates of the Caribbean: At World's End	Captain Barbosa, long believed to be dead, ha...	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...]	[Johnny Depp, Orlando Bloom, Keira Knightley]	[{"credit_id": "52fe4232c3a36847f800b579", "de...
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...]	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[{"credit_id": "54805967c3a36829b5002c41", "de...

In []:

In [55]: # Crew

df['crew']

Out[55]:

```

0    [{"credit_id": "52fe48009251416c750aca23", "de...
1    [{"credit_id": "52fe4232c3a36847f800b579", "de...
2    [{"credit_id": "54805967c3a36829b5002c41", "de...
3    [{"credit_id": "52fe4781c3a36847f81398c3", "de...
4    [{"credit_id": "52fe479ac3a36847f813eaa3", "de...
...
4804   [{"credit_id": "52fe44eec3a36847f80b280b", "de...
4805   [{"credit_id": "52fe487dc3a368484e0fb013", "de...
4806   [{"credit_id": "52fe4df3c3a36847f8275ecf", "de...
4807   [{"credit_id": "52fe4ad9c3a368484e16a36b", "de...
4808   [{"credit_id": "58ce021b9251415a390165d9", "de...
Name: crew, Length: 4806, dtype: object

```

In [56]: df['crew'][0]

```
Out[56]: '[{"credit_id": "52fe48009251416c750aca23", "department": "Editing", "gender": 0, "id": 1721, "job": "Editor", "name": "Stephen E. Rivkin"}, {"credit_id": "539c47ecc3a36810e3001f87", "department": "Art", "gender": 2, "id": 496, "job": "Production Design", "name": "Rick Carter"}, {"credit_id": "54491c89c3a3680fb4001cf7", "department": "Sound", "gender": 0, "id": 900, "job": "Sound Designer", "name": "Christopher Boyes"}, {"credit_id": "54491cb70e0a267480001bd0", "department": "Sound", "gender": 0, "id": 900, "job": "Supervising Sound Editor", "name": "Christopher Boyes"}, {"credit_id": "539c4a4cc3a36810c9002101", "department": "Production", "gender": 1, "id": 1262, "job": "Casting", "name": "Mali Finn"}, {"credit_id": "5544ee3b925141499f0008fc", "department": "Sound", "gender": 2, "id": 1729, "job": "Original Music Composer", "name": "James Horner"}, {"credit_id": "52fe48009251416c750ac9c3", "department": "Directing", "gender": 2, "id": 2710, "job": "Director", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750ac9d9", "department": "Writing", "gender": 2, "id": 2710, "job": "Writer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca17", "department": "Editing", "gender": 2, "id": 2710, "job": "Editor", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca29", "department": "Production", "gender": 2, "id": 2710, "job": "Producer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca3f", "department": "Writing", "gender": 2, "id": 2710, "job": "Screenplay", "name": "James Cameron"}, {"credit_id": "539c4987c3a36810ba0021a4", "department": "Art", "gender": 2, "id": 7236, "job": "Art Direction", "name": "Andrew Menzies"}, {"credit_id": "549598c3c3a3686ae9004383", "department": "Visual Effects", "gender": 0, "id": 6690, "job": "Visual Effects Producer", "name": "Jill Brooks"}, {"credit_id": "52fe48009251416c750aca4b", "department": "Production", "gender": 1, "id": 6347, "job": "Production Design", "name": "Natalie Radford"}, {"credit_id": "52fe48009251416c750aca5f", "department": "Cinematography", "gender": 1, "id": 6347, "job": "Cinematographer", "name": "Robert McLachlan"}, {"credit_id": "52fe48009251416c750aca6b", "department": "Production", "gender": 1, "id": 6347, "job": "Production Design", "name": "Natalie Radford"}, {"credit_id": "52fe48009251416c750aca77", "department": "Cinematography", "gender": 1, "id": 6347, "job": "Cinematographer", "name": "Robert McLachlan"}, {"credit_id": "52fe48009251416c750aca85", "department": "Production", "gender": 1, "id": 6347, "job": "Production Design", "name": "Natalie Radford"}, {"credit_id": "52fe48009251416c750aca9b", "department": "Cinematography", "gender": 1, "id": 6347, "job": "Cinematographer", "name": "Robert McLachlan"}, {"credit_id": "52fe48009251416c750acaab", "department": "Production", "gender": 1, "id": 6347, "job": "Production Design", "name": "Natalie Radford"}, {"credit_id": "52fe48009251416c750acaeb", "department": "Cinematography", "gender": 1, "id": 6347, "job": "Cinematographer", "name": "Robert McLachlan"}, {"credit_id": "52fe48009251416c750acaef", "department": "Production", "gender": 1, "id": 6347, "job": "Production Design", "name": "Natalie Radford"}, {"credit_id": "52fe48009251416c750acaef", "department": "Cinematography", "gender": 1, "id": 6347, "job": "Cinematographer", "name": "Robert McLachlan"}]
```

In []:

In [57]: def fetch_director(text):

```
l = []

for i in ast.literal_eval(text):
    if i['job'] == 'Director':
        l.append(i['name'])

return l
```

In [58]: fetch_director(df['crew'][0])

```
Out[58]: ['James Cameron']
```

```
In [ ]: █
```

```
In [59]: █ df['crew'].apply(fetch_director)
```

```
Out[59]: 0                [James Cameron]
          1                [Gore Verbinski]
          2                [Sam Mendes]
          3                [Christopher Nolan]
          4                [Andrew Stanton]
          ...
          4804              [Robert Rodriguez]
          4805              [Edward Burns]
          4806              [Scott Smith]
          4807              [Daniel Hsia]
          4808      [Brian Herzlinger, Jon Gunn, Brett Winn]
Name: crew, Length: 4806, dtype: object
```

```
In [60]: █ df['crew'] = df['crew'].apply(fetch_director)
```

```
In [ ]: █
```

In [61]: df.head()

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...]	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[James Cameron]
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...]	[Johnny Depp, Orlando Bloom, Keira Knightley]	[Gore Verbinski]
2	206647	Spectre	A cryptic message from Bond's past sends him o...	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...]	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[Sam Mendes]
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i...]	[Christian Bale, Michael Caine, Gary Oldman]	[Christopher Nolan]
4	49529	John Carter	John Carter is a war-weary, former military ca...	[Action, Adventure, Science Fiction]	[based on novel, mars, medallion, space travel...]	[Taylor Kitsch, Lynn Collins, Samantha Morton]	[Andrew Stanton]

In []:

In [62]: # overview

```
df['overview']
```

Out[62]: 0 In the 22nd century, a paraplegic Marine is di...
 1 Captain Barbossa, long believed to be dead, ha...
 2 A cryptic message from Bond's past sends him o...
 3 Following the death of District Attorney Harve...
 4 John Carter is a war-weary, former military ca...
 ...
 4804 El Mariachi just wants to play his guitar and ...
 4805 A newlywed couple's honeymoon is upended by th...
 4806 "Signed, Sealed, Delivered" introduces a dedic...
 4807 When ambitious New York attorney Sam is sent t...
 4808 Ever since the second grade when he first saw ...
Name: overview, Length: 4806, dtype: object

```
In [63]: df['overview'][0]
```

```
Out[63]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization.'
```

```
In [64]: 'i am a boy'
```

```
Out[64]: 'i am a boy'
```

```
In [65]: 'i am a boy'.split()
```

```
Out[65]: ['i', 'am', 'a', 'boy']
```

```
In [ ]:
```

```
In [66]: df['overview'].apply(lambda x: x.split())
```

```
Out[66]: 0      [In, the, 22nd, century,, a, paraplegic, Marin...
 1      [Captain, Barbossa,, long, believed, to, be, d...
 2      [A, cryptic, message, from, Bond's, past, send...
 3      [Following, the, death, of, District, Attorney...
 4      [John, Carter, is, a, war-weary,, former, mili...
 ...
 4804    [El, Mariachi, just, wants, to, play, his, gui...
 4805    [A, newlywed, couple's, honeymoon, is, upended...
 4806    ["Signed,, Sealed,, Delivered", introduces, a, ...
 4807    [When, ambitious, New, York, attorney, Sam, is...
 4808    [Ever, since, the, second, grade, when, he, fi...
Name: overview, Length: 4806, dtype: object
```

```
In [ ]:
```

```
In [67]: df['overview'] = df['overview'].apply(lambda x: x.split())
```

In [68]: df[:4]

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...]	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...]	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[James Cameron]
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...]	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...]	[Johnny Depp, Orlando Bloom, Keira Knightley]	[Gore Verbinski]
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...]	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...]	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[Sam Mendes]
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...]	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i...]	[Christian Bale, Michael Caine, Gary Oldman]	[Christopher Nolan]

In []:

In [69]: [4] + [78] + ['a']

Out[69]: [4, 78, 'a']

In []:

In [70]: # (df['overview'] + df['genres'] + df['keywords'] + df['cast'] + df['crew'])[0]

In [71]: df['tags'] = df['overview'] + df['genres'] + df['keywords'] + df['cast'] + df['crew']

In [72]: df[:3]

Out[72]:

	movie_id	title	overview	genres	keywords	cast	crew	tags
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...]	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon...]	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[James Cameron]	[In, the, 22nd, century,, a, paraplegic, Marin...]
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...]	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india ...]	[Johnny Depp, Orlando Bloom, Keira Knightley]	[Gore Verbinski]	[Captain, Barbossa,, long, believed, to, be, d...]
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...]	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi...]	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[Sam Mendes]	[A, cryptic, message, from, Bond's, past, send...]

In []:

In []:

Final Dataframe

In [73]:

```
data = df[['movie_id', 'title', 'tags']]
data
```

Out[73]:

	movie_id	title	tags
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...
...
4804	9367	EI Mariachi	[EI, Mariachi, just, wants, to, play, his, gui...
4805	72766	Newlyweds	[A, newlywed, couple's, honeymoon, is, upended...
4806	231617	Signed, Sealed, Delivered	["Signed., Sealed., Delivered", introduces, a,...
4807	126186	Shanghai Calling	[When, ambitious, New, York, attorney, Sam, is...
4808	25975	My Date with Drew	[Ever, since, the, second, grade, when, he, fi...

4806 rows × 3 columns

In []:

In [74]:

```
print(data['tags'][0])
```

```
['In', 'the', '22nd', 'century,', 'a', 'paraplegic', 'Marine', 'is', 'dispatched', 'to', 'the', 'moon',
'Pandora', 'on', 'a', 'unique', 'mission,', 'but', 'becomes', 'torn', 'between', 'following', 'orders',
'and', 'protecting', 'an', 'alien', 'civilization.', 'Action', 'Adventure', 'Fantasy', 'Science Fiction',
'culture clash', 'future', 'space war', 'space colony', 'society', 'space travel', 'futuristic', 'romance',
'space', 'alien', 'tribe', 'alien planet', 'cgi', 'marine', 'soldier', 'battle', 'love affair',
'anti war', 'power relations', 'mind and soul', '3d', 'Sam Worthington', 'Zoe Saldana', 'Sigourney Weaver',
'James Cameron']
```

In [95]:

```
'i am a boy'
```

Out[95]:

```
'i am a boy'
```

```
In [98]: ┏ 'i am a boy'.replace(" ", '')
```

```
Out[98]: 'iamaboy'
```

```
In [76]: ┏ data['tags'].apply(lambda x: [i.replace(' ', '') for i in x])
```

```
Out[76]: 0      [In, the, 22nd, century,, a, paraplegic, Marin...
 1      [Captain, Barbossa,, long, believed, to, be, d...
 2      [A, cryptic, message, from, Bond's, past, send...
 3      [Following, the, death, of, District, Attorney...
 4      [John, Carter, is, a, war-weary,, former, mili...
 ...
 4804    [El, Mariachi, just, wants, to, play, his, gui...
 4805    [A, newlywed, couple's, honeymoon, is, upended...
 4806    ["Signed,, Sealed,, Delivered", introduces, a, ...
 4807    [When, ambitious, New, York, attorney, Sam, is...
 4808    [Ever, since, the, second, grade, when, he, fi...
Name: tags, Length: 4806, dtype: object
```

```
In [77]: ┏ data['tags'] = data['tags'].apply(lambda x: [i.replace(' ', '') for i in x])
```

```
In [78]: ┏ data[:5]
```

	movie_id	title	tags
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...

```
In [80]: ┏ print(data['tags'][0])
```

```
['In', 'the', '22nd', 'century,', 'a', 'paraplegic', 'Marine', 'is', 'dispatched', 'to', 'the', 'moon',  
'Pandora', 'on', 'a', 'unique', 'mission,', 'but', 'becomes', 'torn', 'between', 'following', 'orders',  
'and', 'protecting', 'an', 'alien', 'civilization.', 'Action', 'Adventure', 'Fantasy', 'ScienceFiction',  
'cultureclash', 'future', 'spacewar', 'spacecolony', 'society', 'spacettravel', 'futuristic', 'romance',  
'space', 'alien', 'tribe', 'alienplanet', 'cgi', 'marine', 'soldier', 'battle', 'loveaffair', 'antiwar',  
'powerrelations', 'mindandsoul', '3d', 'SamWorthington', 'ZoeSaldana', 'SigourneyWeaver', 'JamesCamer  
n']
```

```
In [ ]: ┏
```

```
In [ ]: ┏
```

```
In [81]: ┏ 'i am a boy'.split()
```

```
Out[81]: ['i', 'am', 'a', 'boy']
```

```
In [82]: ┏ " ".join(['i', 'am', 'a', 'boy'])
```

```
Out[82]: 'i am a boy'
```

```
In [ ]: ┏
```

```
In [83]: ┏ data['tags'].apply(lambda x: " ".join(x))
```

```
Out[83]: 0      In the 22nd century, a paraplegic Marine is di...  
1      Captain Barbosa, long believed to be dead, ha...  
2      A cryptic message from Bond's past sends him o...  
3      Following the death of District Attorney Harve...  
4      John Carter is a war-weary, former military ca...  
      ...  
4804     El Mariachi just wants to play his guitar and ...  
4805     A newlywed couple's honeymoon is upended by th...  
4806     "Signed, Sealed, Delivered" introduces a dedic...  
4807     When ambitious New York attorney Sam is sent t...  
4808     Ever since the second grade when he first saw ...  
Name: tags, Length: 4806, dtype: object
```

In [84]: ┌ data['tags'] = data['tags'].apply(lambda x: " ".join(x))

In [85]: ┌ data

Out[85]:

	movie_id	title	tags
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...
2	206647	Spectre	A cryptic message from Bond's past sends him o...
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...
4	49529	John Carter	John Carter is a war-weary, former military ca...
...
4804	9367	El Mariachi	El Mariachi just wants to play his guitar and ...
4805	72766	Newlyweds	A newlywed couple's honeymoon is upended by th...
4806	231617	Signed, Sealed, Delivered	"Signed, Sealed, Delivered" introduces a dedic...
4807	126186	Shanghai Calling	When ambitious New York attorney Sam is sent t...
4808	25975	My Date with Drew	Ever since the second grade when he first saw ...

4806 rows × 3 columns

In []: ┌

In [86]: ┌ data['tags'][0]

Out[86]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization. Action Adventure Fantasy Scienc eFiction cultureclash future spacewar spacecolony society spacetravel futuristic romance space alien tri be alienplanet cgi marine soldier battle loveaffair antiwar powerrelations mindandsoul 3d SamWorthington ZoeSaldana SigourneyWeaver JamesCameron'

In []: ┌

In []:

NLP Concepts to preprocess textual data

- Lower Case
- Tokenization
- Stemming
- Stopwords Removal

In [89]:

data[:5]

Out[89]:

	movie_id	title	tags
0	19995	Avatar	In the 22nd century, a paraplegic Marine is di...
1	285	Pirates of the Caribbean: At World's End	Captain Barbosa, long believed to be dead, ha...
2	206647	Spectre	A cryptic message from Bond's past sends him o...
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...
4	49529	John Carter	John Carter is a war-weary, former military ca...

In [91]:

data['tags'][0]

Out[91]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization. Action Adventure Fantasy Scienc eFiction cultureclash future spacewar spacecolony society spacetravel futuristic romance space alien tri be alienplanet cgi marine soldier battle loveaffair antiwar powerrelations mindandsoul 3d SamWorthington ZoeSaldana SigourneyWeaver JamesCameron'

In [93]:

'I am a Boy'.lower()

Out[93]: 'i am a boy'

```
In [ ]: █
```

```
In [112]: █ 'act', 'action', 'actions', 'actionable'
```

```
Out[112]: ('act', 'action', 'actions', 'actionable')
```

```
In [ ]: █
```

```
In [115]: █ ps = PorterStemmer()
```

```
def preprocess_text(text):  
  
    updated_text = []  
  
    for i in text.split():  
        lower = i.lower()  
        updated_text.append(ps.stem(lower))  
  
    return " ".join(updated_text)
```

```
In [116]: █ preprocess_text(data['tags'][0])
```

```
Out[116]: 'in the 22nd century, a parapleg marin is dispatch to the moon pandora on a uniqu mission, but becom tor  
n between follow order and protect an alien civilization. action adventur fantasi sciencefict culturecla  
sh futur spacewar spacecoloni societi spacetravel futurist romanc space alien tribe alienplanet cgi mari  
n soldier battl loveaffair antiwar powerrel mindandsoul 3d samworthington zoesaldana sigourneyweav james  
cameron'
```

```
In [ ]: █
```

```
In [117]: ► data['tags'].apply(preprocess_text)
```

```
Out[117]: 0      in the 22nd century, a parapleg marin is disp...
 1      captain barbossa, long believ to be dead, ha c...
 2      a cryptic messag from bond' past send him on a...
 3      follow the death of district attorney harvey d...
 4      john carter is a war-weary, former militari ca...
 ...
 4804    el mariachi just want to play hi guitar and ca...
 4805    a newlyw couple' honeymoon is upend by the arr...
 4806    "signed, sealed, delivered" introduc a dedic q...
 4807    when ambiti new york attorney sam is sent to s...
 4808    ever sinc the second grade when he first saw h...
Name: tags, Length: 4806, dtype: object
```

```
In [ ]: ►
```

```
In [118]: ► data['tags'] = data['tags'].apply(preprocess_text)
```

```
In [ ]: ►
```

In [119]: ⏷ data

Out[119]:

	movie_id	title	tags
0	19995	Avatar	in the 22nd century, a parapleg marin is disp...
1	285	Pirates of the Caribbean: At World's End	captain barbossa, long believ to be dead, ha c...
2	206647	Spectre	a cryptic messag from bond' past send him on a...
3	49026	The Dark Knight Rises	follow the death of district attorney harvey d...
4	49529	John Carter	john carter is a war-weary, former militari ca...
...
4804	9367	EI Mariachi	el mariachi just want to play hi guitar and ca...
4805	72766	Newlyweds	a newlyw couple' honeymoon is upend by the arr...
4806	231617	Signed, Sealed, Delivered	"signed, sealed, delivered" introduc a dedic q...
4807	126186	Shanghai Calling	when ambiti new york attorney sam is sent to s...
4808	25975	My Date with Drew	ever sinc the second grade when he first saw h...

4806 rows × 3 columns

In []: ⏷

Bag of Words (BOW) Method

In [120]: ⏷ cv = CountVectorizer(max_features=5000, stop_words='english')

```
In [121]: ► data['tags']
```

```
Out[121]: 0      in the 22nd century, a parapleg marin is disp...
1      captain barbossa, long believ to be dead, ha c...
2      a cryptic messag from bond' past send him on a...
3      follow the death of district attorney harvey d...
4      john carter is a war-weary, former militari ca...
...
4804    el mariachi just want to play hi guitar and ca...
4805    a newlyw couple' honeymoon is upend by the arr...
4806    "signed, sealed, delivered" introduc a dedic q...
4807    when ambiti new york attorney sam is sent to s...
4808    ever sinc the second grade when he first saw h...
Name: tags, Length: 4806, dtype: object
```

```
In [123]: ► vectors = cv.fit_transform(data['tags']).toarray()
vectors
```

```
Out[123]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]], dtype=int64)
```

```
In [124]: ► vectors.shape
```

```
Out[124]: (4806, 5000)
```

```
In [128]: ► len(vectors)
```

```
Out[128]: 4806
```

```
In [125]: ► cv.get_feature_names_out()
```

```
Out[125]: array(['000', '007', '10', ..., 'zone', 'zoo', 'zooeydeschanel'],
 dtype=object)
```

```
In [127]: # for i in cv.get_feature_names_out():
    print(i)
```

```
000
007
10
100
11
12
13
14
15
16
17
17th
18
18th
18thcenturi
19
1910
1920
1930
1940
```

```
In [ ]: #
```

```
In [ ]: #
```

Calculate Cosine Similarity between vectors

```
In [129]: # vectors
```

```
Out[129]: array([[0, 0, 0, ..., 0, 0, 0],
                  [0, 0, 0, ..., 0, 0, 0],
                  [0, 0, 0, ..., 0, 0, 0],
                  ...,
                  [0, 0, 0, ..., 0, 0, 0],
                  [0, 0, 0, ..., 0, 0, 0],
                  [0, 0, 0, ..., 0, 0, 0]], dtype=int64)
```

```
In [130]: ┏ similarity = cosine_similarity(vectors)
similarity
```

```
Out[130]: array([[1.          , 0.08346223, 0.0860309 , ... , 0.04499213, 0.          ,
                   0.          ],
                  [0.08346223, 1.          , 0.06063391, ... , 0.02378257, 0.          ,
                   0.02615329],
                  [0.0860309 , 0.06063391, 1.          , ... , 0.02451452, 0.          ,
                   0.          ],
                  ...,
                  [0.04499213, 0.02378257, 0.02451452, ... , 1.          , 0.03962144,
                   0.04229549],
                  [0.          , 0.          , 0.          , ... , 0.03962144, 1.          ,
                   0.08714204],
                  [0.          , 0.02615329, 0.          , ... , 0.04229549, 0.08714204,
                   1.          ]])
```

```
In [131]: ┏ similarity[0]
```

```
Out[131]: array([1.          , 0.08346223, 0.0860309 , ... , 0.04499213, 0.          ,
                   0.          ])
```

```
In [132]: ┏ similarity[1]
```

```
Out[132]: array([0.08346223, 1.          , 0.06063391, ... , 0.02378257, 0.          ,
                   0.02615329])
```

```
In [ ]: ┏
```

```
In [133]: ┏ similarity[0].shape
```

```
Out[133]: (4806,)
```

```
In [134]: ► sorted(similarity[0])
```

In []: ➤

```
In [135]: ┆ sorted(similarity[0])[-10: -1]
```

```
Out[135]: [0.23174488732966075,  
          0.23179316248638276,  
          0.24455799402225922,  
          0.24511108480187255,  
          0.25038669783359574,  
          0.255608593705383,  
          0.2605130246476754,  
          0.26901379342448517,  
          0.28676966733820225]
```

In []:

In [136]: ► data[:3]

Out[136]:

	movie_id	title	tags
0	19995	Avatar	in the 22nd century, a parapleg marin is dispa...
1	285	Pirates of the Caribbean: At World's End	captain barbossa, long believe to be dead, ha c...
2	206647	Spectre	a cryptic messag from bond' past send him on a...

In []: ►

In [140]: ► sorted(similarity[0], reverse=True)

Out[140]: [1.0000000000000002,
 0.28676966733820225,
 0.26901379342448517,
 0.2605130246476754,
 0.255608593705383,
 0.25038669783359574,
 0.24511108480187255,
 0.24455799402225922,
 0.23179316248638276,
 0.23174488732966075,
 0.2278389747471728,
 0.2252817784447915,
 0.22269966704152225,
 0.21853668936906193,
 0.21239769762143662,
 0.2108663315950723,
 0.2105263157894737,
 0.20443988269091456,
 0.20437977982832192,
 0.20205070126182276

In []: ►

```
In [142]: ┌─ list(enumerate(similarity[0]))
```

```
Out[142]: [(0, 1.0000000000000002),  
 (1, 0.08346223261119858),  
 (2, 0.08603090020146065),  
 (3, 0.0734718358370645),  
 (4, 0.1892994097121204),  
 (5, 0.10838874619051501),  
 (6, 0.04024218182927669),  
 (7, 0.14673479641335554),  
 (8, 0.05923488777590923),  
 (9, 0.0967301666813349),  
 (10, 0.10259783520851541),  
 (11, 0.09464970485606021),  
 (12, 0.09037128496931669),  
 (13, 0.04499212706658476),  
 (14, 0.12824729401064427),  
 (15, 0.06282808624375433),  
 (16, 0.07894736842105264),  
 (17, 0.13977653617040256),  
 (18, 0.09493290614465533),  
 (19, 0.00000100017045201)]
```

```
In [ ]: ┌─
```

```
In [144]: ┆ sorted(list(enumerate(similarity[0])), reverse=True)
```

```
Out[144]: [(4805, 0.0),  
 (4804, 0.0),  
 (4803, 0.04499212706658476),  
 (4802, 0.046829290579084706),  
 (4801, 0.019252140716412975),  
 (4800, 0.0),  
 (4799, 0.052631578947368425),  
 (4798, 0.04223886030955117),  
 (4797, 0.0),  
 (4796, 0.0),  
 (4795, 0.0),  
 (4794, 0.0),  
 (4793, 0.05407380704358751),  
 (4792, 0.0),  
 (4791, 0.0),  
 (4790, 0.0582716546748065),  
 (4789, 0.060833032924035954),  
 (4788, 0.0),  
 (4787, 0.019117977822546817),  
 (4786, 0.0)]
```

```
In [146]: ┏▶ sorted(list(enumerate(similarity[0])), reverse=True, key=lambda x: x[1])
```

```
Out[146]: [(0, 1.0000000000000002),  
 (1216, 0.28676966733820225),  
 (2409, 0.26901379342448517),  
 (3730, 0.2605130246476754),  
 (507, 0.255608593705383),  
 (539, 0.25038669783359574),  
 (582, 0.24511108480187255),  
 (1204, 0.24455799402225922),  
 (1194, 0.23179316248638276),  
 (778, 0.23174488732966075),  
 (4048, 0.2278389747471728),  
 (1920, 0.2252817784447915),  
 (61, 0.22269966704152225),  
 (2786, 0.21853668936906193),  
 (172, 0.21239769762143662),  
 (972, 0.2108663315950723),  
 (322, 0.2105263157894737),  
 (2333, 0.20443988269091456),  
 (3608, 0.20437977982832192),  
 (266, 0.20395070126182276)]
```

```
In [148]: ┏▶ sorted(list(enumerate(similarity[0])), reverse=True, key=lambda x: x[1])[1:6]
```

```
Out[148]: [(1216, 0.28676966733820225),  
 (2409, 0.26901379342448517),  
 (3730, 0.2605130246476754),  
 (507, 0.255608593705383),  
 (539, 0.25038669783359574)]
```

```
In [ ]: ┏▶
```

```
In [150]: ┏▶ data.iloc[0]
```

```
Out[150]: movie_id          19995  
title                Avatar  
tags      in the 22nd century, a parapleg marin is dispa...  
Name: 0, dtype: object
```

```
In [152]: ► data['tags'][0]
```

```
Out[152]: 'in the 22nd century, a parapleg marin is dispatch to the moon pandora on a uniqu mission, but becom tor  
n between follow order and protect an alien civilization. action adventur fantasi sciencefict culturecla  
sh futur spacewar spacecoloni societi spacetravel futurist romanc space alien tribe alienplanet cgi mari  
n soldier battl loveaffair antiwar powerrel mindandsoul 3d samworthington zoesaldana sigourneyweav james  
cameron'
```

```
In [151]: ► data.iloc[1216]
```

```
Out[151]: movie_id          440  
title           Aliens vs Predator: Requiem  
tags      a sequel to 2004' alien vs. predator, the icon...  
Name: 1216, dtype: object
```

```
In [153]: ► data.iloc[2409]
```

```
Out[153]: movie_id          679  
title           Aliens  
tags      when ripley' lifepod is found by a salvag crew...  
Name: 2409, dtype: object
```

```
In [ ]: ►
```

Final Recommendation Function

```
In [182]: ► def recommend(movie):
```

```
    movie_index = data[data['title'] == movie].index[0]
    distance = similarity[movie_index]
    movie_list = sorted(list(enumerate(distance)), reverse=True, key=lambda x: x[1])[1:6]

    for i in movie_list:
        #     print(i[0])
        print(data.iloc[i[0]].title)
```

In [183]: ► recommend('Avatar')

Aliens vs Predator: Requiem
Aliens
Falcon Rising
Independence Day
Titan A.E.

In [184]: ► recommend('Iron Man')

Iron Man 3
Iron Man 2
Avengers: Age of Ultron
The Avengers
Captain America: Civil War

In [185]: ► recommend('Spider-Man')

Spider-Man 3
Spider-Man 2
The Amazing Spider-Man 2
Arachnophobia
Kick-Ass

In [186]: ► recommend('Batman')

Batman
Batman & Robin
Batman Begins
Batman Returns
The R.M.

In [190]: ► recommend('''Pirates of the Caribbean: At World's End''')

Pirates of the Caribbean: Dead Man's Chest
Pirates of the Caribbean: The Curse of the Black Pearl
Pirates of the Caribbean: On Stranger Tides
Life of Pi
20,000 Leagues Under the Sea

```
In [ ]: █
```

```
In [ ]: █
```

```
In [ ]: █
```

```
In [188]: █ data
```

Out[188]:

	movie_id	title	tags
0	19995	Avatar	in the 22nd century, a parapleg marin is disp...
1	285	Pirates of the Caribbean: At World's End	captain barbossa, long believ to be dead, ha c...
2	206647	Spectre	a cryptic messag from bond' past send him on a...
3	49026	The Dark Knight Rises	follow the death of district attorney harvey d...
4	49529	John Carter	john carter is a war-weary, former militari ca...
...
4804	9367	EI Mariachi	el mariachi just want to play hi guitar and ca...
4805	72766	Newlyweds	a newlyw couple' honeymoon is upend by the arr...
4806	231617	Signed, Sealed, Delivered	"signed, sealed, delivered" introduc a dedic q...
4807	126186	Shanghai Calling	when ambiti new york attorney sam is sent to s...
4808	25975	My Date with Drew	ever sinc the second grade when he first saw h...

4806 rows × 3 columns

```
In [ ]: █
```

```
In [176]: █ data.iloc[1216]
```

```
Out[176]: movie_id          440
           title            'Aliens vs Predator: Requiem'
           tags             'a sequel to 2004\' alien vs. predator, the icon...
           Name: 1216, dtype: object
```

```
In [177]: █ data.iloc[1216].title
```

```
Out[177]: 'Aliens vs Predator: Requiem'
```

```
In [ ]: █
```

```
In [ ]: █
```

```
In [160]: █ data[data['title'] == 'Iron Man']
```

```
Out[160]:
       movie_id      title
       tags
68     1726 Iron Man after be held captiv in an afghan cave, billio...
```

```
In [161]: █ data[data['title'] == 'Iron Man'].index[0]
```

```
Out[161]: 68
```

```
In [ ]: █
```

```
In [ ]: █
```

```
In [ ]: █
```

In []: