

```
In [1]: import numpy as np

import pandas as pd
```

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

credits = pd.read_csv('tmdb_5000_credits.csv')
```

```
In [3]: movies = movies.merge(credits, on='title')
```

```
In [4]: movies.head(1)
```

Out[4]:

	budget	genres	homepage	id	keywords	original_language	original_title	overview	popularity	production_companies
0	237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": 1464, "name": "culture clash"}]	en	Avatar	In the 22nd century, a paraplegic Marine is di...	150.437577	[{"name": "Ingenious Film Partners", "id": 289..

1 rows × 23 columns



```
In [5]: # removing not needed columns from the data set
# say, for a content based movie recommender system 'budget' attribute won't have any affect
# removing only such columns for more ex: Like; homepage, original_language, original_title (title is already there)..
```

```
In [6]: movies = movies[['movie_id', 'title', 'overview', 'genres', 'keywords', 'cast', 'crew']]
```

In [7]: `movies.head(5)`

Out[7]:

	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...
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...	[{"id": 28, "name": "Action"}, {"id": 80, "nam...	[{"id": 849, "name": "dc comics"}, {"id": 853,...	[{"cast_id": 2, "character": "Bruce Wayne / Ba...	[{"credit_id": "52fe4781c3a36847f81398c3", "de...
4	49529	John Carter	John Carter is a war-weary, former military ca...	[{"id": 28, "name": "Action"}, {"id": 12, "nam...	[{"id": 818, "name": "based on novel"}, {"id":...	[{"cast_id": 5, "character": "John Carter", "c...	[{"credit_id": "52fe479ac3a36847f813eaa3", "de...

In [8]: `movies.isnull().sum()`

Out[8]:

```

movie_id    0
title       0
overview    3
genres      0
keywords    0
cast        0
crew        0
dtype: int64

```

In [9]: `movies.dropna(inplace = True)`

In [10]: `movies.duplicated().sum()`

Out[10]: 0

In [11]: *# preprocessing genres*

```
movies.iloc[0].genres
```

Out[11]: '[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'

In [12]: *# found stringified list of dictionaries*

firstly have to convert the string into a pure list using 'ast' module

```
import ast
```

now will use ast.literal_eval(string) where ever needed.

In [13]: **def** convert(string):

```
    L = []
```

```
    for x in ast.literal_eval(string):
```

```
        L.append(x['name'])
```

```
    return L
```

In [14]: `movies['genres'] = movies['genres'].apply(convert)`

In [15]: *# same can be applied on the keywords column*

```
movies['keywords'] = movies['keywords'].apply(convert)
```

```
In [16]: def convertcast(string):
        L = []
        counter = 0
        for x in ast.literal_eval(string):
            if counter != 3:
                L.append(x['name'])
                counter += 1
            else:
                break

        return L
```

```
In [17]: movies['cast'] = movies['cast'].apply(convertcast)
```

```
In [18]: movies.head(5)
```

Out[18]:

	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 Barbossa, 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...
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]	[{"credit_id": "52fe4781c3a36847f81398c3", "de...
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]	[{"credit_id": "52fe479ac3a36847f81398c3", "de...

```
In [19]: def fetchdirector(string):
    L = []
    for x in ast.literal_eval(string):
        if x['job'] == 'Director':
            L.append(x['name'])
            break

    return L
```

```
In [20]: movies['crew'] = movies['crew'].apply(fetchdirector)
```

```
In [21]: movies.head(5)
```

Out[21]:

	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 [22]: # converting overview from 'string' to 'list'

movies['overview'] = movies['overview'].apply(lambda x:x.split())
```

In [23]: `movies.head(5)`

Out[23]:

	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]
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...	[Action, Adventure, Science Fiction]	[based on novel, mars, medallion, space travel...	[Taylor Kitsch, Lynn Collins, Samantha Morton]	[Andrew Stanton]

In [24]: *# removing spaces from between words like 'Sam Worthington' & 'Sam Mendes',
if not removed may result in alternate responses*

```
movies['genres'] = movies['genres'].apply(lambda x:[i.replace(" ", "") for i in x])
movies['cast'] = movies['cast'].apply(lambda x:[i.replace(" ", "") for i in x])
movies['keywords'] = movies['keywords'].apply(lambda x:[i.replace(" ", "") for i in x])
movies['crew'] = movies['crew'].apply(lambda x:[i.replace(" ", "") for i in x])
```

In [25]: `movies.head()`

Out[25]:

	movie_id	title	overview	genres	keywords	cast	crew
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...	[Action, Adventure, Fantasy, ScienceFiction]	[cultureclash, future, spacewar, spacecolony, ...	[SamWorthington, ZoeSaldana, SigourneyWeaver]	[JamesCameron]
1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d...	[Adventure, Fantasy, Action]	[ocean, drugabuse, exoticisland, eastindiatrad...	[JohnnyDepp, OrlandoBloom, KeiraKnightley]	[GoreVerbinski]
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...	[Action, Adventure, Crime]	[spy, basedonnovel, secretagent, sequel, mi6, ...	[DanielCraig, ChristophWaltz, LéaSeydoux]	[SamMendes]
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...	[Action, Crime, Drama, Thriller]	[dccomics, crimefighter, terrorist, secretiden...	[ChristianBale, MichaelCaine, GaryOldman]	[ChristopherNolan]
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...	[Action, Adventure, ScienceFiction]	[basedonnovel, mars, medallion, spacetravel, p...	[TaylorKitsch, LynnCollins, SamanthaMorton]	[AndrewStanton]

In [26]: `# combining all columns into a single column named 'tags'`

```
movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movies['cast'] + movies['crew']
```

```
In [27]: movies.head()
```

```
Out[27]:
```

	movie_id	title	overview	genres	keywords	cast	crew	tags
0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin...	[Action, Adventure, Fantasy, ScienceFiction]	[culturedclash, future, spacewar, spacecolony, ...	[SamWorthington, ZoeSaldana, SigourneyWeaver]	[JamesCameron]	[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, drugabuse, exoticisland, eastindiatrad...	[JohnnyDepp, OrlandoBloom, KeiraKnightley]	[GoreVerbinski]	[Captain, Barbossa,, long, believed, to, be, d...
2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send...	[Action, Adventure, Crime]	[spy, basedonnovel, secretagent, sequel, mi6, ...	[DanielCraig, ChristophWaltz, LéaSeydoux]	[SamMendes]	[A, cryptic, message, from, Bond's, past, send...
3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney...	[Action, Crime, Drama, Thriller]	[dccomics, crimefighter, terrorist, secretiden...	[ChristianBale, MichaelCaine, GaryOldman]	[ChristopherNolan]	[Following, the, death, of, District, Attorney...
4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili...	[Action, Adventure, ScienceFiction]	[basedonnovel, mars, medallion, spacetravel, p...	[TaylorKitsch, LynnCollins, SamanthaMorton]	[AndrewStanton]	[John, Carter, is, a, war-weary,, former, mili...

```
In [28]: new_df = movies[['movie_id', 'title', 'tags']]
```

```
In [29]: # converting the (list)tags into (string)tags
```

```
new_df['tags'] = new_df['tags'].apply(lambda x: " ".join(x))
```

C:\Users\Vishwas\AppData\Local\Temp\ipykernel_3504\160216261.py:3: 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 (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
new_df['tags'] = new_df['tags'].apply(lambda x: " ".join(x))
```



```
In [30]: new_df.head()
```

```
Out[30]:
```

	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...

```
In [31]: new_df['tags'][0]
```

```
Out[31]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn be
tween following orders and protecting an alien civilization. Action Adventure Fantasy ScienceFiction cultureclash fu
ture spacewar spacecolony society spacetravel futuristic romance space alien tribe alienplanet cgi marine soldier ba
ttle loveaffair antiwar powerrelations mindandsoul 3d SamWorthington ZoeSaldana SigourneyWeaver JamesCameron'
```

```
In [32]: new_df['tags'] = new_df['tags'].apply(lambda x: x.lower())
```

C:\Users\Vishwas\AppData\Local\Temp\ipykernel_3504\1380776331.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 (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
new_df['tags'] = new_df['tags'].apply(lambda x: x.lower())
```

In [33]: `new_df.head()`

Out[33]:

	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...

In [34]: `import nltk
from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()`

In [35]: `def stem(text):
 y = []
 for i in text.split():
 y.append(ps.stem(i))
 return " ".join(y)`

In [36]: `new_df['tags'] = new_df['tags'].apply(stem)`

C:\Users\Vishwas\AppData\Local\Temp\ipykernel_3504\3213734980.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 (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

`new_df['tags'] = new_df['tags'].apply(stem)`

```
In [37]: # performing vectorization: text --> vector
# using 'bag of words';

from sklearn.feature_extraction.text import CountVectorizer

cv = CountVectorizer(max_features = 5000, stop_words = 'english')
```

```
In [38]: vectors = cv.fit_transform(new_df['tags']).toarray()
```

```
In [39]: vectors
```

```
Out[39]: 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 [40]: vectors[0]
```

```
Out[40]: array([0, 0, 0, ..., 0, 0, 0], dtype=int64)
```

```
In [41]: cv.get_feature_names_out()
```

```
Out[41]: array(['000', '007', '10', ..., 'zone', 'zoo', 'zoeydeschanel'],
               dtype=object)
```

```
In [42]: new_df['tags'][0]
```

```
Out[42]: 'in the 22nd century, a parapleg marin is dispatch to the moon pandora on a uniqu mission, but becom torn between fo
llow order and protect an alien civilization. action adventur fantasi sciencefict cultureclash futur spacewar spacec
oloni societi spacetravel futurist romanc space alien tribe alienplanet cgi marin soldier battl loveaffair antiwar p
owerrel mindandsoul 3d samworthington zoesaldana sigourneyweav jamescameron'
```

```
In [43]: # checking for similarities
```

```
from sklearn.metrics.pairwise import cosine_similarity
```

```
In [44]: similarity = cosine_similarity(vectors)
```

```
In [45]: similarity[1]
```

```
Out[45]: array([0.08346223, 1.          , 0.06063391, ..., 0.02378257, 0.          ,  
               0.02615329])
```

```
In [46]: # recommending Logic
```

```
def recommend(movie):  
    # finding movie index in data  
    movie_index = new_df[new_df['title'] == movie].index[0]  
    # distances with other movies  
    distances = similarity[movie_index]  
    # getting top 5 similar movies  
    movies_list = sorted(list(enumerate(distances)), reverse = True, key=lambda x:x[1])[1:6]  
  
    for i in movies_list:  
        print(new_df.iloc[i[0]].title)
```

```
In [47]: recommend('Batman Begins')
```

```
The Dark Knight  
Batman  
Batman  
The Dark Knight Rises  
10th & Wolf
```

```
In [48]: import pickle
```

```
In [49]: pickle.dump(new_df, open('movies.pkl', 'wb'))
```

```
In [50]: pickle.dump(new_df.to_dict(), open('movie_dict.pkl', 'wb'))
```

```
In [51]: pickle.dump(similarity, open('similarity.pkl', 'wb'))
```

```
In [ ]:
```