

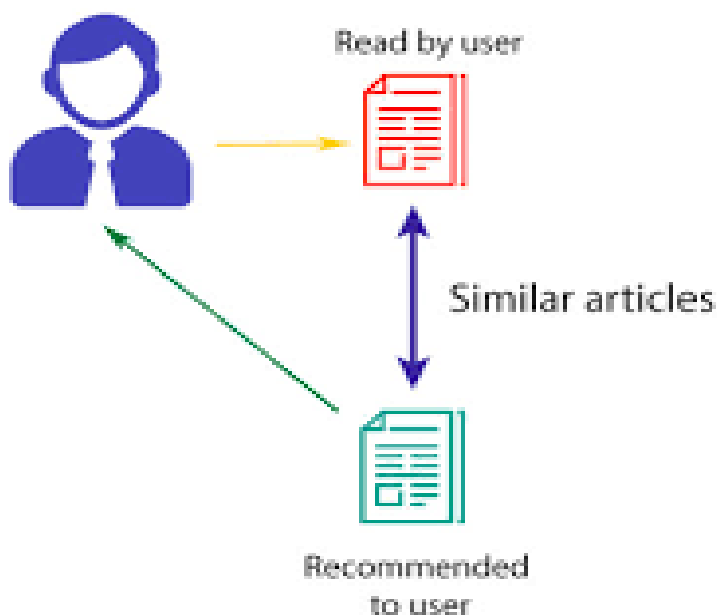
# CONTENT BASED RECOMMENDER SYSTEM

## Description

- This project as the title suggests focuses on recommending based upon content present in the dataset i.e., the item-item interaction.
- This dataset uses the Netflix dataset and so the items here are the movies and the TV shows. Item-item interaction essentially refers matching users based upon their item history.
- Let's say a user performs an action on items A, B, C and another user performs actions on items A, C. Then automatically the user will be recommended item B by the algorithm.

## Content Based Filtering

### CONTENT-BASED FILTERING



- Not based upon user ratings and preferences. Instead, solely based upon the context of the movie or the TV show.
- Here we have used the text preprocessing techniques of punctuation marks removal, stripping and removal of stop-words along with making the words in the text lowercase.
- After this use the Tfidf vectorizer as it gives importance to the words based upon their value in the text to convert the text into vectorized format.
- TFIDF Vectorizer stands for term frequency inverse document frequency vectorizer. It is given by the formula:

WORD (W) = (Number of times the word repeats in a sentence/Total number of words in that sentence) \* (log (Total number of sentences/Total number of sentences containing that word))

- Then each movie/ TV show will have its context vectors corresponding to the content present in their description.
- Since different words have different importance, it is necessary for us to give less important words less importance.
- Present in scikit learn and no separate downloading required. So, a matrix is formed with columns as the words and rows as the items with row of the item being its corresponding vector.
- Finally use the Nearest Neighbors unsupervised to find the k nearest movies or TV shows based on the cosine similarity.
- Cosine similarity between item vector v1 and v2 is  $(v1.v2)/(\text{mod}(v1) * \text{mod}(v2))$ . Thus, the items closer to each other have more cosine similarity and so less cosine distance which is 1-cosine similarity.
- Item-item interaction is all about comparing between the items and then recommending based upon the item. Content based recommendation uses comparing and finding the relation between items based upon vectors of the description or the reviews.

## Using the repository (Use the commands given in the Readme file)

### 1] For fetching of the dataset from Kaggle:

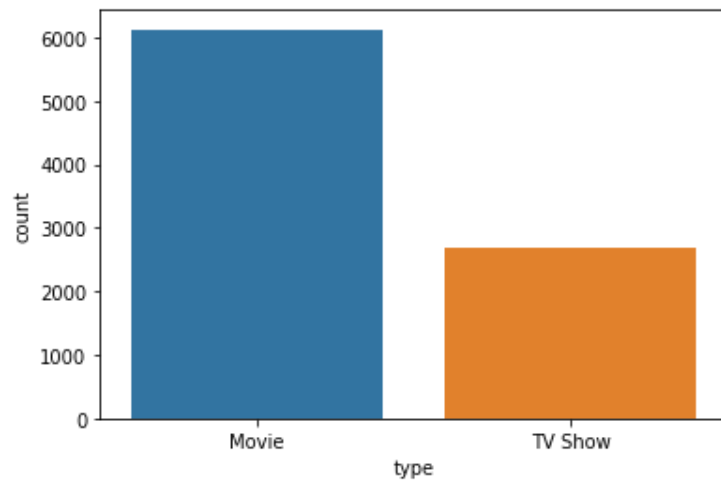
Active Kaggle account is required. Datasets used in the repo are

Netflix Movies and Shows: <https://www.kaggle.com/shivamb/netflix-shows>

IMDB Movies and Shows: <https://www.kaggle.com/stefanoleone992/imdb-extensive-dataset>

IMDB Ratings: <https://www.kaggle.com/stefanoleone992/imdb-extensive-dataset>

- Netflix is one of the most popular media and video streaming platforms. They have over 8000 movies or tv shows available on their platform.
- IMDb is the most popular movie website and it combines movie plot description, Megastore ratings, critic and user ratings and reviews, release dates, and many more aspects.
- Run the tox command mentioned and your datasets will be ready to use. For live data just input the name of the Movie or TV show in string format in the text box.
- After merging the datasets, we can infer those movies are in excess compared to tv shows.



- Final dataset after merging will look like this

show_id	type	title	director	cast	country	release_year	rating	duration	listed_in	description	add_year	add_month	add_day	tm
0	s1	Movie	Dick Johnson Is Dead	Kristen Johnson	NaN	United States	2020	PG-13	90 min	Documentaries	As her father nears the end of his life, from...	2021.0	9.0	25.0
1	s2	TV Show	Blood & Water	NaN	Amma Osumba, Khosi Ngema, Qali Mabasa, Thabani...	South Africa	2021	TV-MA	Seasons	International TV Shows, TV Dramas, TV Mysteries	After crossing paths at a party, a Cape Town l...	2021.0	9.0	24.0
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajja, Tracy Coombs, Samuel Jaffi, Nabil...	NaN	2021	TV-MA	Season	Crime TV Shows, International TV Shows, TV Ad...	To protect his family from a powerful drug lor...	2021.0	9.0	24.0
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	2021	TV-MA	Season	Documentaries, Reality TV	Faith, frictions and toilet talk go down amo...	2021.0	9.0	24.0
4	s5	TV Show	Kala Factory	NaN	Mayur Moh, Jitendra Kumar, Ranjan Raj, Alank...	India	2021	TV-MA	Seasons	International TV Shows, Romantic TV Shows, TV ...	In a city of coaching centers known to train l...	2021.0	9.0	24.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
8799	s3800	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	2007	R	150 min	Cult Movies, Dramas, Thrillers	A political cartoonist, a crime reporter and a...	2019.0	11.0	20.0

- The above dataset uses the Netflix dataset which has important features like the type of content, director, cast, description (extremely important), country of origin, etc.
- On this we have merged the IMDB dataset for ratings and genre of the Netflix content. For live data just input the name of the Movie or TV show in string format in the text box.

```

kshir@DESKTOP-NASUUL9 MINGW64 ~
$ cd C:/Users/kshir/OneDrive/Desktop/DS_PROJECT_2/ds_project/package
kshir@DESKTOP-NASUUL9 MINGW64 ~OneDrive/Desktop/DS_PROJECT_2/ds_project/package (project)
$ cd ..
kshir@DESKTOP-NASUUL9 MINGW64 ~OneDrive/Desktop/DS_PROJECT_2/ds_project (project)
$ tox -e fetch_data
Fetch_data create: C:/Users/kshir/OneDrive/Desktop/DS_PROJECT_2/ds_project/tox/train_test_package
Fetch_data installdeps: --requirements test_requirements.txt
Fetch_data install: anyio==3.3.4, appdirs==1.4.0, atomicwrites==1.4.0, attrs==21.2.0, bash==0.6, black==20.8b1, certifi==2021.10.8, chardet==3.0.4, click==7.1.2, colorama==0.4.4, fastapi==0.70.0, feature-engine==1.0.2, flake8==3.9.2, hll==0.9.0, idna==2.10, iniconfig==1.1.1, isort==5.8.0, Jinja2==3.0.2, joblib==1.0.1, kaggle==1.5.2, loguru==0.5.3, MarkupSafe==2.0.1, mccabe==0.6.1, mypy==0.812, mypy-extensions==0.4.3, numpy==1.20.3, packaging==21.2, pandas==1.2.5, pathspec==0.9.0, patsy==0.5.2, pluggy==1.0.0, py==1.11.0, pycodestyle==2.7.0, pydantic==1.8.2, pyflakes==2.3.1, pygments==2.4.7, pytest==6.2.5, python-dateutil==2.8.2, python-json-logger==0.1.11, python-multipart==0.5.5, python-slugify==5.0.2, pytz==2021.3, PyYAML==5.4.1, requests==2.23.0, ruamel.yaml==0.16.12, ruamel.yaml.clib==0.2.6, scikit-learn==0.24.2, scipy==1.7.2, six==1.16.0, sniffio==1.2.0, starlette==0.16.0, statsmodels==0.13.0, strictyaml==1.3.2, text-unidecode==1.3, threadpoolctl==3.0.0, tomli==0.10.2, toml==0.10.2, tqdm==4.62.3, typed-ast==1.4.3, typing-extensions==3.7.4.3, urllib3==1.22, uvicorn==0.11.8, websocket==8.1, win32-setctime==1.0.3
Fetch_data run-test-pre: PYTHONHASHSEED='0'
Fetch_data run-test: commands[0] | kaggle datasets download -d shivamb/netflix-shows -p ./package/recommender_model/datasets
Downloading netflix-shows.zip to ./package/recommender_model/datasets
100%##### 1.34M/1.34M [00:00<00:00, 3.93MB/s]
Fetch_data run-test: commands[1] | unzip package/recommender_model/datasets/netflix-shows.zip -d package/recommender_model/datasets
WARNING: test command found but not installed in testenv
cmd: C:/Program Files/Git/usr/bin/unzip.exe
env: C:/Users/kshir/OneDrive/Desktop/DS_PROJECT_2/ds_project/tox/train_test_package
Maybe you forgot to specify a dependency? See also the allowlist_external envconfig setting.
DEPRECATION WARNING: this will be an error in tox 4 and above!
Archive: package/recommender_model/datasets/netflix-shows.zip
Inflating: package/recommender_model/datasets/netflix_titles.csv
Fetch_data run-test: commands[2] | kaggle datasets download -d stefano1eone992/imdb-extensive-dataset -p ./package/recommender_model/datasets
Downloading imdb-extensive-dataset.zip to ./package/recommender_model/datasets
100%##### 82.3K/82.3K [00:18<00:00, 4.66MB/s]
Fetch_data run-test: commands[3] | unzip package/recommender_model/datasets/imdb-extensive-dataset.zip -d package/recommender_model/datasets
WARNING: test command found but not installed in testenv
cmd: C:/Program Files/Git/usr/bin/unzip.exe
env: C:/Users/kshir/OneDrive/Desktop/DS_PROJECT_2/ds_project/tox/train_test_package
Maybe you forgot to specify a dependency? See also the allowlist_external envconfig setting.
DEPRECATION WARNING: this will be an error in tox 4 and above!
Archive: package/recommender_model/datasets/imdb-extensive-dataset.zip
Inflating: package/recommender_model/datasets/IMDb movies.csv
Inflating: package/recommender_model/datasets/IMDb names.csv
Inflating: package/recommender_model/datasets/IMDb ratings.csv
Inflating: package/recommender_model/datasets/IMDb title_principals.csv
Fetch_data run-test: commands[4] | mv package/recommender_model/datasets/IMDb title_principals.csv'

```

## 2] For training the MODEL:

- First of all, we need to pre-process and clean the data in all the three datasets. Then we need to use the left join and join the ratings and movies dataset on the Netflix dataset using left join.
- After the merging we need to drop the duplicates and the extract the important columns like 'description' from the dataset. The final dataset has 8804 entries with 69.6% of them being movies while the rest being TV shows.
- Then TfIdf vectorizer is used with each content vector being of (1,19183) dimensions as the corresponding to the context of the items using the above-mentioned formula for the text separated by spacing. Thus, the data frame comprises of words in columns and items in the index.
- With these content vectors we can then find out the cosine similarity. More is the cosine similarity more likely is the movie or the TV show i.e., the content to be recommended.
- Algorithm used is the Nearest Neighbours which is an unsupervised algorithm with the content being the index of the data frame and the rows being the content vectors. The nearest neighbours taken here are 10 i.e., top 10 recommendations are to be made.

```

kshir@DESKTOP-NASUUL9 MINGW64 ~OneDrive/Desktop/DS_PROJECT_2/ds_project (project)
$ tox -e train_test_package
train_test_package installed: anyio==3.3.4, appdirs==1.4.0, atomicwrites==1.4.0, attrs==21.2.0, bash==0.6, black==20.8b1, certifi==2021.10.8, chardet==3.0.4, click==7.1.2, colorama==0.4.4, fastapi==0.70.0, feature-engine==1.0.2, flake8==3.9.2, hll==0.9.0, idna==2.10, iniconfig==1.1.1, isort==5.8.0, Jinja2==3.0.2, joblib==1.0.1, kaggle==1.5.2, loguru==0.5.3, MarkupSafe==2.0.1, mccabe==0.6.1, mypy==0.812, mypy-extensions==0.4.3, numpy==1.20.3, packaging==21.2, pandas==1.2.5, pathspec==0.9.0, patsy==0.5.2, pluggy==1.0.0, py==1.11.0, pycodestyle==2.7.0, pydantic==1.8.2, pyflakes==2.3.1, pygments==2.4.7, pytest==6.2.5, python-dateutil==2.8.2, python-json-logger==0.1.11, python-multipart==0.5.5, python-slugify==5.0.2, pytz==2021.3, PyYAML==5.4.1, requests==2.23.0, ruamel.yaml==0.16.12, ruamel.yaml.clib==0.2.6, scikit-learn==0.24.2, scipy==1.7.2, six==1.16.0, sniffio==1.2.0, starlette==0.16.0, statsmodels==0.13.0, strictyaml==1.3.2, text-unidecode==1.3, threadpoolctl==3.0.0, tomli==0.10.2, toml==0.10.2, tqdm==4.62.3, typed-ast==1.4.3, typing-extensions==3.7.4.3, urllib3==1.22, uvicorn==0.11.8, websockets==8.1, win32-setctime==1.0.3
train_test_package run-test-pre: PYTHONHASHSEED='0'
train_test_package run-test: commands[0] | python package/recommender_model/train_pipeline.py
C:/Users/kshir/OneDrive/Desktop/DS_PROJECT_2/ds_project/package/recommender_model/train_pipeline.py:16: SyntaxWarning: "is" with a literal. Did you mean "=="?
  # Create a list of movies and TV shows titles
{'predictions': {'Top 10 recommendations': ['The Emigrant', 'Ram Dass, Going Home', 'Jeans', 'Kaminey', 'Dil Vil Pyaar Vyaar', 'Main Hoon Na', 'My Boss's Daughter', 'Day and Night', 'Qila', 'One by Two']], 'version': '4.0.2', 'errors': None}
{'predictions': {'Classic': ['Pulp Fiction', 'Schindler's List', 'Chupke Chupke', 'Platoon', 'Rocky', 'Bawarchi', 'My Fair Lady', 'Philadelphia', 'Lolita'], 'Comedy': ['Much Ado About Nothing', 'Joker', '3 Idiots', 'Super Deluxe', 'Love Ni Bhavai', 'Taxi Driver', 'Andhadhun', 'Chupke Chupke', 'Queen'], 'International': ['City of God', 'Seven', 'Koshish', 'Ant', 'Dr. Kashinath Ghanekar', 'Much Ado About Nothing', 'Eh Janam Tumhare Lekhe', 'Oththa Seruppu Size 7', 'Mallasham', 'Joker'], 'Romance': ['Much Ado About Nothing', 'Koshish', 'Qismat', 'Love Ni Bhavai', 'Sadma', 'Chupke Chupke', 'Sairat', '2 States', 'Andaz Apna Apna']], 'version': '4.0.2', 'errors': None}
{'predictions': {'India': ['Ant', 'Dr. Kashinath Ghanekar', 'Koshish', 'Eh Janam Tumhare Lekhe', 'Oththa Seruppu Size 7', 'Manto', 'Mallasham', 'Punjab 1984', 'Merku Thodan', 'chi Malai', 'Black Friday']], 'version': '4.0.2', 'errors': None}

```

## 3] For running the API:

- First, we build the API using fastapi module leveraging the python async.io web framework.

- FastApi module is compatible with pydantic module which helps in input data validation and type checking automatically using schemas.
- Then we run this API with the uvicorn ASGI web server on <http://localhost:8001/recommendersystem>
- Also, we perform logging using loguru module to display the status messages. Live data for recommendation is only the name of the movie or the TV show.

[https://drive.google.com/file/d/1V24rJx\\_ih04pNahCqhbXJcoR5XR--YP\\_/view?usp=sharing](https://drive.google.com/file/d/1V24rJx_ih04pNahCqhbXJcoR5XR--YP_/view?usp=sharing)

Refer to the above video.

**FOR RUNNING THE REPOSITORY REFER TO THE COMMANDS IN README FILE.**