# 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. Thus, this is a very basic project depicting the recommendation system used in the OTT platforms.
- The models used are NLP-based Tfidf-Vectorizer to convert the words or text to vectorized format and then Nearest Neighbors algorithm is used to find the interaction between the content using cosine similarity.

## Advantages of TFIDF Vectorizer

Content based recommendation deals with the content i.e., the description/reviews corresponding to the items in the dataset.

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))

- In the given text not, all words are important in terms of vectorizing.
- Since different words have different importance, it is necessary for us to give less important words less importance.
- So, this vectorizer gives unequal importance to different words of the text. This helps in forming efficient vectors from the given text.
- 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.

## Advantages of Nearest Neighbors Algorithm (NNA)

- Unsupervised algorithm to find the cosine similarity between the item vectors. Cosine similarity gives the nearest or closest items pertaining to the current item.
- 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.
- Cosine similarity uses the cosine law between vectors. Less is the cosine similarity more is the cosine distance.
- Thus, NNA helps us in finding the best k items pertaining to the current item all based upon cosine similarity between the items.

# Using the repository (use the commands given in the dataset section of Readme File)

## 1] Fetching the dataset from Kaggle:

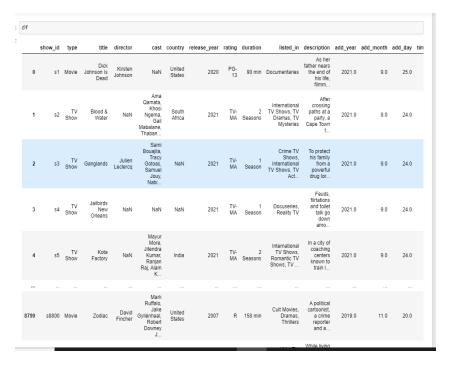Active Kaggle account is required. Datasets used in the repo are

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.



## 2] Training of the models:

- There are three datasets and so we need to pre-process and clean them individually. After they are in proper format then we have to merge the movies and ratings datasets on the Netflix dataset to obtain a complete dataset corresponding to the content.

- Then vectorize the required column(s) using the Tfidf Vectorizer and then use the vectors to find the cosine similarity between them. Lastly apply the Nearest Neighbour Algorithm to find the cosine similarity between the vectors.

- As we can see here the results because of content i.e., taking only description into consideration are poor. If we were to take other features into consideration like directors, cast, etc. we would have to use Count Vectorizer by giving space between those features and not between text of a single feature.



3] For running of the API

- First, we build the API using fastapi module leveraging the python async.io web framework.
- Then we run this API with the uvicorn ASGI web server.
- Also, we perform logging using loguru module.

Refer to the video:
https://drive.google.com/file/d/1V24rJx_ih04pNahCqhbxJcoR5XR--YP_/view?usp=sharing

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