


# PolyTube Recommender Systems



Andrew Thenedi

# Introduction

- PolyTube is an online media platform which provides online teleplay services.
  - For teleplays already in service, the platform will record user ratings and give an average rating for each teleplay on the web page.
  - For recently published teleplays without ratings, how to predict them is essential for the investment policy of PolyTube.
  - On the other hand, to attract users watching more teleplays on the platform, providing accurate and personalized recommendation services is also meaningful for revenue increase.
  - This project will help PolyTube improve their system.
- 

# Prerequisites

You need to have installed following softwares and libraries in your machine before running this project.

1. Python 3

2. Anaconda

- It will install ipython notebook and most of the libraries which are needed like NumPy, Pandas, Matplotlib, Seaborn, scikit-learn, and Tensorflow.



# Getting Started

Start by downloading the project and run "COMP4434\_IPROJ\_Final\_Submission.ipynb" file in ipython-notebook.



# Data Cleaning



## Get and Preprocess the Training Dataset ('Teleplay.csv' and 'Rating.csv') as 'df\_teleplay' and 'df\_rating' data frames.

- Remove 'rating' or 'members' column values less or equal to 0, in the 'df\_teleplay' data frame.
- Remove 'rating' column values less or equal to 0, in the 'df\_rating' data frame.
- Find the mean ratings for each 'teleplay\_id' column values in the 'df\_rating' data frame.
  - Merge such values in the 'df\_teleplay' data frame on the 'rating\_2' column.
- Check empty values existence in 'df\_teleplay' data frame, per column.
  - In this case, there exists an empty value in the 'genre' and 'rating\_2' column values.
- Drop the empty values of the 'genre' column in the 'df\_teleplay' data frame.
  - Check empty values existence in the 'df\_teleplay' data frame, per column.
    - There is an empty value in the 'rating\_2' column values in this case.

# Get and Preprocess the Testing Dataset 'New\_Teleplay.csv' as 'df\_new\_teleplay' data frame.

- Merge the mean ratings for each 'teleplay\_id' column values in 'df\_rating' dataframe values in 'df\_new\_teleplay' dataframe on 'rating\_2' column.
- Check empty values existence in 'df\_new\_teleplay' dataframe, per column:
  - In this case, there exists an empty value in 'genre' , 'type', 'rating' , and 'rating\_2' column values.
- Replace empty values in 'df\_new\_teleplay' dataframe with "nan" string data type, except 'rating' and 'rating\_2' column values.
  - Check empty values existence in 'df\_new\_teleplay' dataframe, per column.
    - In this case, there exists an empty value in the 'rating' and 'rating\_2' column values.

Predict the empty values of 'rating\_2' column from 'df\_teleplay' and  
'df\_new\_teleplay' dataframes





- Split 'df\_teleplay' dataframe into 'df\_teleplay\_train\_predict\_rating\_2' dataframe as training dataset and 'df\_teleplay\_test\_predict\_rating\_2' dataframe as test dataset.
  - Reset index of 'df\_teleplay\_train\_predict\_rating\_2' and 'df\_teleplay\_test\_predict\_rating\_2' dataframes.
- Split 'df\_new\_teleplay' dataframe into 'df\_new\_teleplay\_train\_predict\_rating\_2' dataframe as training dataset and 'df\_new\_teleplay\_test\_predict\_rating\_2' dataframe as test dataset:
  - Reset index of 'df\_new\_teleplay\_train\_predict\_rating\_2' and 'df\_new\_teleplay\_test\_predict\_rating\_2' dataframes.



- Encode and Normalize the values of 'type', 'genre', and 'members' columns:
  - Concatenate 'df\_teleplay\_train\_predict\_rating\_2', 'df\_teleplay\_test\_predict\_rating\_2', 'df\_new\_teleplay\_train\_predict\_rating\_2' and 'df\_new\_teleplay\_test\_predict\_rating\_2' data frames as 'df\_merge\_encode\_normalize' dataframe, for 'rating', 'rating\_2', 'members', 'type', and 'genre' column values.
    - Reset index of 'df\_merge\_encode\_normalize' dataframe.
    - Normalize 'rating', 'rating\_2', and 'members' columns values in 'df\_merge\_encode\_normalize' dataframe, which automatically skip the empty values.
    - Encode 'type' and 'genre' columns values in 'df\_merge\_encode\_normalize' dataframe.
  - Merge normalized 'members', 'rating' & 'rating\_2' and encoded 'type' and 'genre' columns values of 'df\_merge\_encode\_normalize' into 'arr\_merge\_encode\_normalize' nparray.
    - Split 'arr\_merge\_encode\_normalize' np array into 'x\_noNA\_teleplay', 'x\_NA\_teleplay', 'y\_noNA\_teleplay', 'y\_NA\_teleplay', 'x\_noNA\_new\_teleplay', 'x\_NA\_new\_teleplay', 'y\_noNA\_new\_teleplay', and 'y\_NA\_new\_teleplay' np arrays.

- Predict the empty values of the 'rating\_2' column from the 'df\_teleplay' data frame.
  - Evaluate the model performance using the non-empty values of the 'rating\_2' column from the 'df\_teleplay' data frame:
    - Split 'x\_noNA\_teleplay' and 'y\_noNA\_teleplay' np arrays into 'x\_train', 'x\_test', 'y\_train', and 'y\_test' np arrays.
    - Perform Artificial Neural Network (ANN) Model in the training dataset.
    - Calculate RMSE of using ANN model in the training dataset.
      - In this case, the value of RMSE is 0.7051800875594707.
  - Perform ANN model, with 'x\_NA\_teleplay' np array as the given input values, to predict the empty values, in 'y\_NA\_teleplay', of 'rating\_2' column from 'df\_teleplay' data frame:
    - Replace the empty values with the predicted values of 'y\_NA\_teleplay' np array to 'df\_teleplay\_test\_predict\_rating\_2' dataframe of 'rating\_2' column.
    - Concatenate 'df\_teleplay\_train\_predict\_rating\_2' and 'df\_teleplay\_test\_predict\_rating\_2' data frames as 'df\_teleplay' dataframe.

- Predict the empty values of 'rating\_2' column from 'df\_new\_teleplay' dataframe:
  - Evaluate the model performance using the non-empty values of the 'rating\_2' column from the 'df\_new\_teleplay' data frame.
    - Combine 'x\_noNA\_new\_teleplay' with 'x\_noNA\_teleplay' (excluding 'rating' column values).
    - Split 'x\_noNA\_new\_teleplay' and 'y\_noNA\_new\_teleplay' np arrays into 'x\_train', 'x\_test', 'y\_train', and 'y\_test' np arrays.
    - Perform Artificial Neural Network (ANN) Model in the training dataset.
    - Calculate RMSE of using ANN model in the training dataset.
      - In this case, the value of RMSE is 1.0958634072744535.
  - Perform ANN model, with 'x\_NA\_new\_teleplay' np array as the given input values, to predict the empty values, in 'y\_NA\_new\_teleplay', of 'rating\_2' column from 'df\_new\_teleplay' dataframe:
    - Replace the empty values with the predicted values of 'y\_NA\_new\_teleplay' np array to 'df\_new\_teleplay\_test\_predict\_rating\_2' dataframe of 'rating\_2' column.
    - Concatenate 'df\_new\_teleplay\_train\_predict\_rating\_2' and 'df\_new\_teleplay\_test\_predict\_rating\_2' dataframes as 'df\_new\_teleplay' dataframe

## Check empty values existence in 'df\_teleplay' and 'df\_new\_teleplay' data frames, per column

- In 'df\_teleplay' data frame, there does not exist an empty value.
- In 'df\_new\_teleplay' data frame, There is an empty value in this case in the 'rating' column values.



# Exploratory Data Analysis (EDA)



- Describe the values of the 'df\_teleplay' data frame on 'members' and 'rating' columns
  - Check the distribution of the 'df\_teleplay' data frame 'members' values.
    - It is shown that the minimum values are the teleplays that have at least 17 'members' value.
  - Check the distribution of the 'df\_teleplay' data frame 'rating' values.
    - It is shown that the minimum values are the teleplays that have at least 1.670 'rating' value.
- Visualize the trend and correlations among data points of the 'df\_teleplay' data frame.
  - 'members' column:
    - It is shown that the number of 'members' is dominated by a small group of teleplays that are blockbusters. The distribution plummets for other teleplays subsequently.
  - 'rating' column:
    - It is shown that the average 'rating' value is around 6.5, as the bell-curve shape permits.
  - 'members' and 'rating' columns:
    - It is shown that the 'rating' values are increasing with the number of 'members'.
  - 'type' column:
    - It is shown that medium-'type' teleplays are the most sought-after, while music-type teleplays have the least demand:
  - 'genre' column:
    - It is shown that the comedy-'genre' has the most demand, followed closely by the action-'genre'.
- After considering all possible relevant relationships among columns in the 'df\_teleplay' data frame, 'members', 'type', and 'genre' column values are being chosen as the input values for task 1, in addition to the 'rating\_2' column values.

Task 1: Predict the empty values of the 'rating' column from the  
'df\_new\_teleplay' data frame





- Concatenate 'df\_teleplay' and 'df\_new\_teleplay' dataframes as 'df\_merge\_encode\_normalize' dataframe, for 'rating\_2' , 'members' , 'type', and 'genre' column values.
  - Reset index of 'df\_merge\_encode\_normalize' dataframe.
  - Save indexes information details about 'df\_merge\_encode\_normalize'.
  - Normalize 'rating\_2' , and 'members' columns values.
  - Encode 'type' and 'genre' columns values.
  - Merge normalized 'rating\_2', and 'members' and encoded 'type' and 'genre' columns values into 'arr\_merge\_encode\_normalize' np array.
    - Split 'arr\_merge\_encode\_normalize' np array into 'x\_teleplay', 'y\_teleplay', 'x\_new\_teleplay', and 'y\_new\_teleplay' np arrays.
- Evaluate the model performance from the 'df\_teleplay' data frame:
  - Split 'x\_teleplay' and 'y\_teleplay' np arrays into 'x\_train', 'x\_test', 'y\_train', and 'y\_test' np arrays.
  - Perform Artificial Neural Network (ANN) Model in the training dataset.
  - Calculate RMSE of using ANN model in the training dataset.
    - In this case, the value of RMSE is 0.49096476235066167.
- Perform ANN model, with 'x\_new\_teleplay' np array as the given input values, to predict the values of 'rating' column from 'df\_new\_teleplay' data frame.
  - Replace the values with the predicted values of 'y\_new\_teleplay' np array to the 'df\_new\_teleplay' data frame of the 'rating' column.
  - Submission.
  - The resulted values are stored in the '18079969D\_task1.csv' file.

Task 2: Predict user 53698's personalized rating of all teleplays



- Get the average rating of each teleplay, based on the rating of each user, to each teleplay in 'df\_rating'.
- Normalize the column values of 'rating\_per\_user\_id' and perform aggregation to get the input values for collaborative filtering.
  - Store the resulted aggregation values as 'ratings' np array.
  - Split 'ratings' np array to 'x\_train', 'x\_test', 'y\_train', and 'y\_test' np arrays.
- Perform Collaborative Filtering Model in the training dataset.
- Plot the RMSE value of using the Collaborative Filtering model in the training dataset with the number of epochs.
- Perform a prediction for user 53698's personalized rating of all teleplays.
  - Submission:
  - The resulted values are stored in the '18079969D\_task2.csv' file.

THANK  
YOU