# WEEK3
# LINEAR REGRESSION

6438169421 Pattaradanai Lakkananithiphan

# CONTENT

## 5043 Rows

Loaded from the CSV file
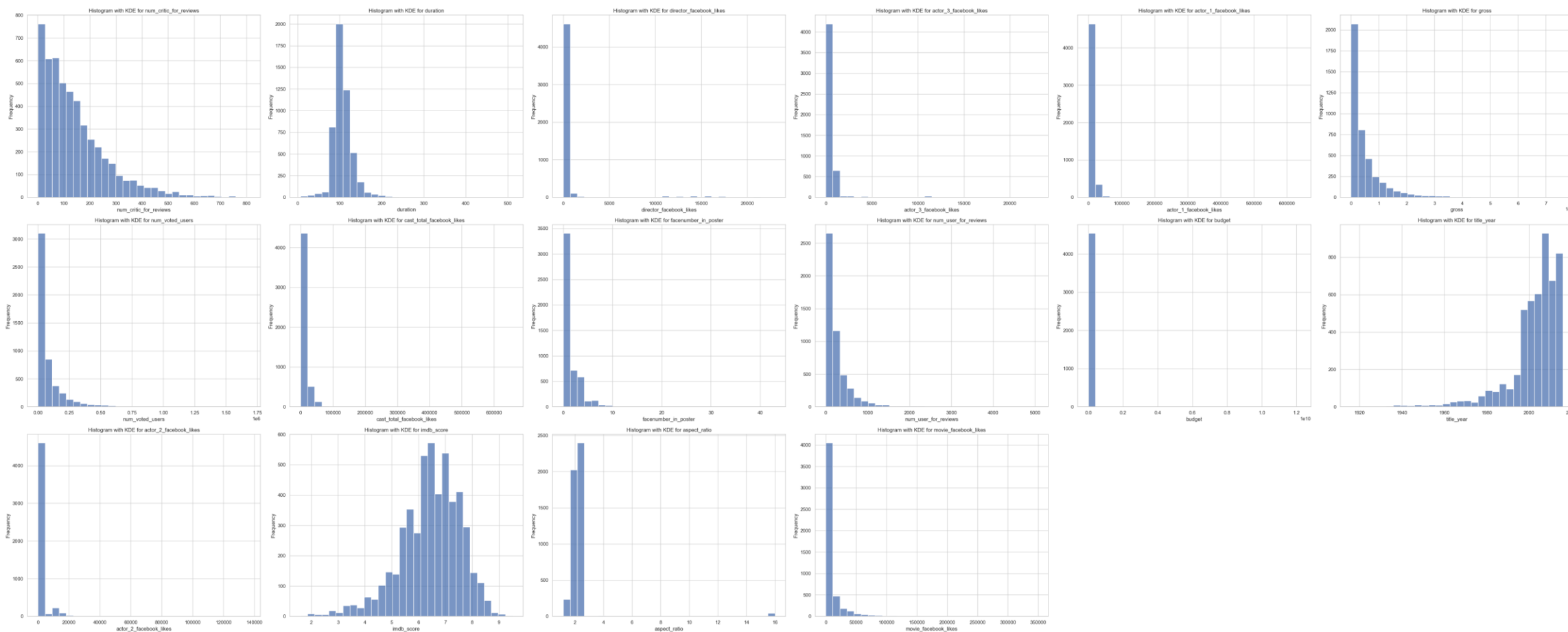
## 27 Columns

16 numerical + 11 categorical

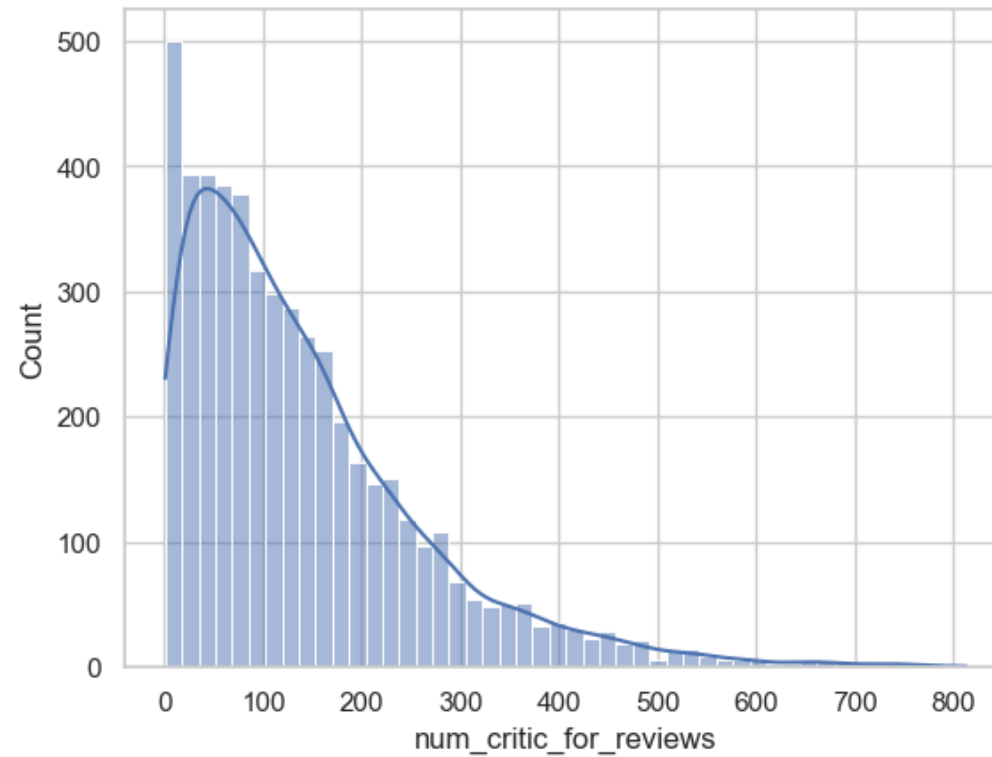## >884 rows have missing values

50 rows have missing label

## Numerical Label

The label is a numerical value of type float

# I: LOADING

# THE DISTRIBUTION OF THE LABEL



The data is right-skewed

# CATEGORY COUNTS

| CATEGORICAL COLUMN NAMES | NUMBER OF CATEGORIES |
|---|---|
| director_name | 2398 |
| actor_2_name | 3032 |
| genres | 914 |
| actor_1_name | 2097 |
| movie_title | 4917 |
| actor_3_name | 3521 |
| plot_keywords | 4760 |
| movie_imdb_link | 4919 |
| language | 47 |
| country | 65 |
| content_rating | 18 |

## II: PREPROCESS
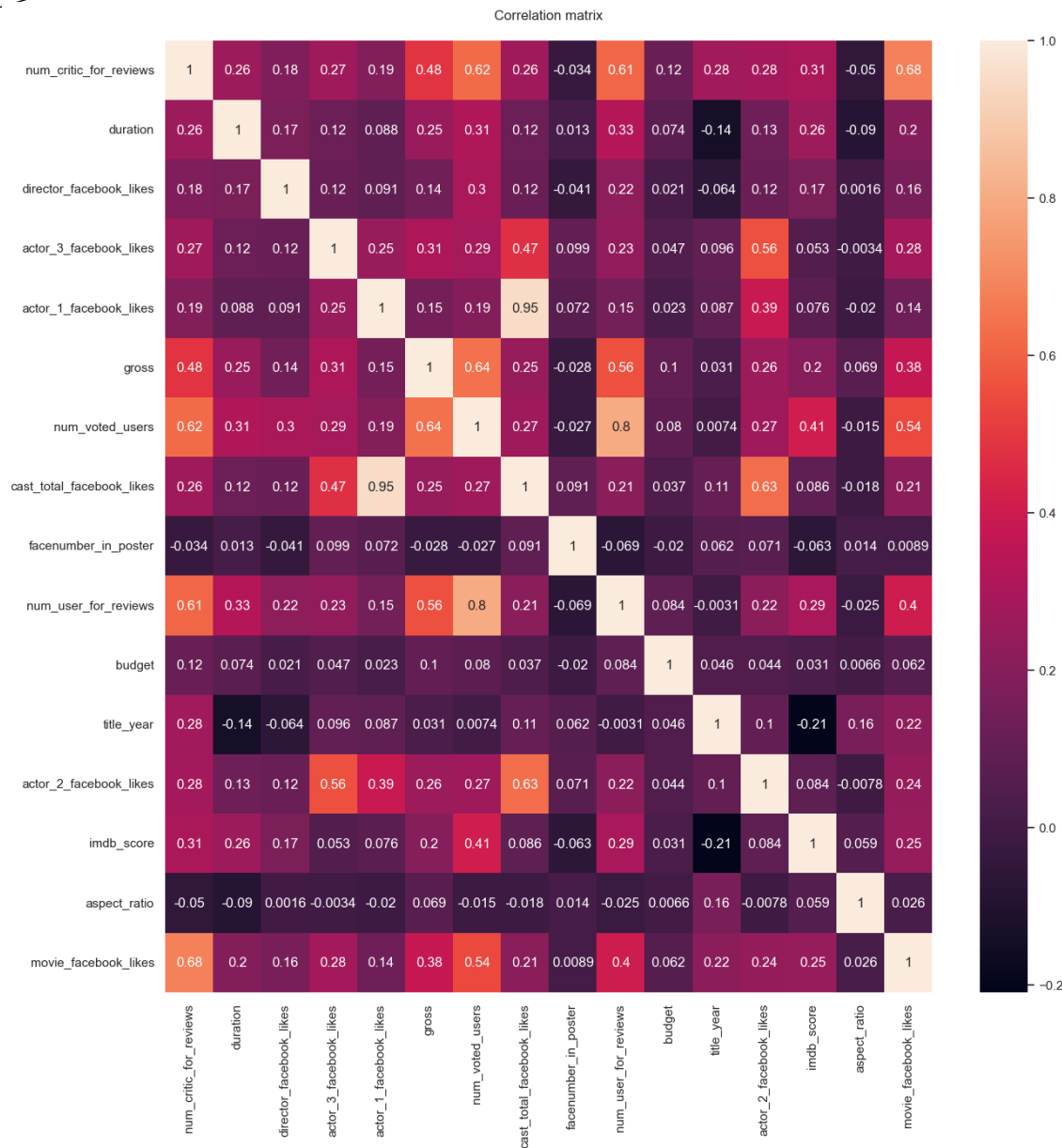
Dropping columns

Dropping/Filling rows

Scaling

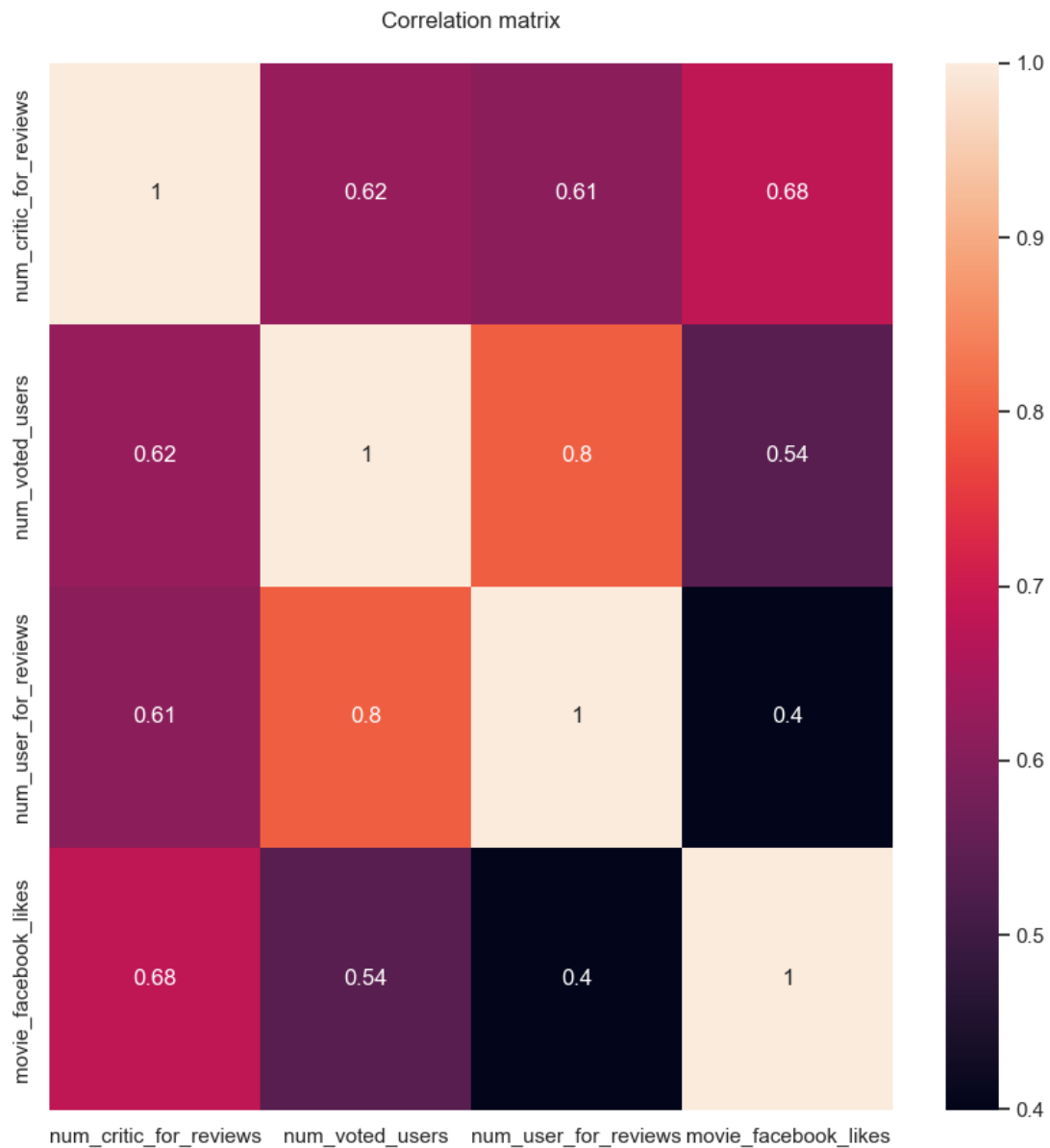Splitting

# DROPPING COLUMNS: CATEGORICAL

- Most if not all categorical columns have a lot more categories than appropriate

- Many are highly unrelated to the label we want to predict

- Some are very biased toward certain categories

- Many have a lot of NaN rows

- As guided by the solution slide

- **CONCLUSION: DROP THE CATEGORICAL COLUMNS**

Correlation matrix

# DROPPING COLUMNS: NUMERICAL

## <- CORRELATION MATRIX

DROP ALL COLUMNS OF ABSOLUTE CORRELATION LESS THAN **0.5** AS GUIDED BY THE SLIDE

Correlation matrix

# DROPPING COLUMNS: NUMERICAL

## <- CORRELATION MATRIX

## THE RESULT

**NOTE\*\*:** THE **NUM_USER_FOR_REVIEWS** AND **NUM_VOTED_USERS** COLUMNS ARE HIGHLY CORRELATED AND WE COULD DROP THE LOWER CORRELATED COLUMN IE **NUM_USER_FOR_REVIEWS**
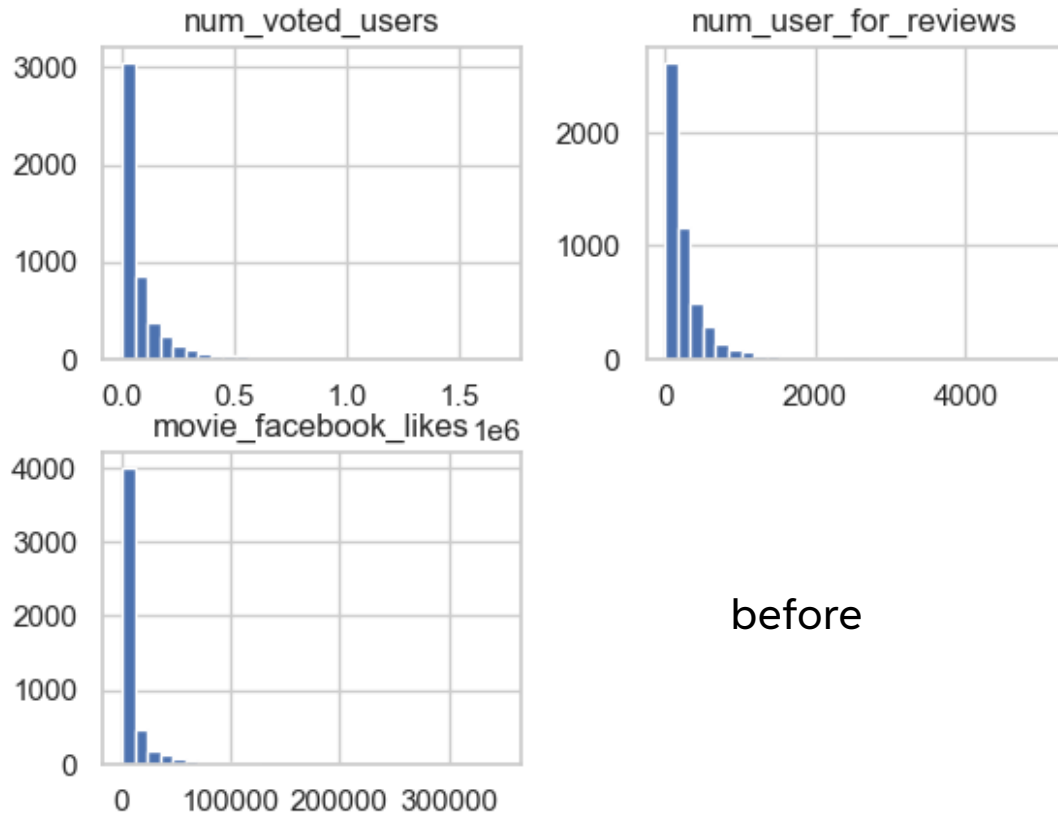
# DROP/FILL THE EMPTY ROWS

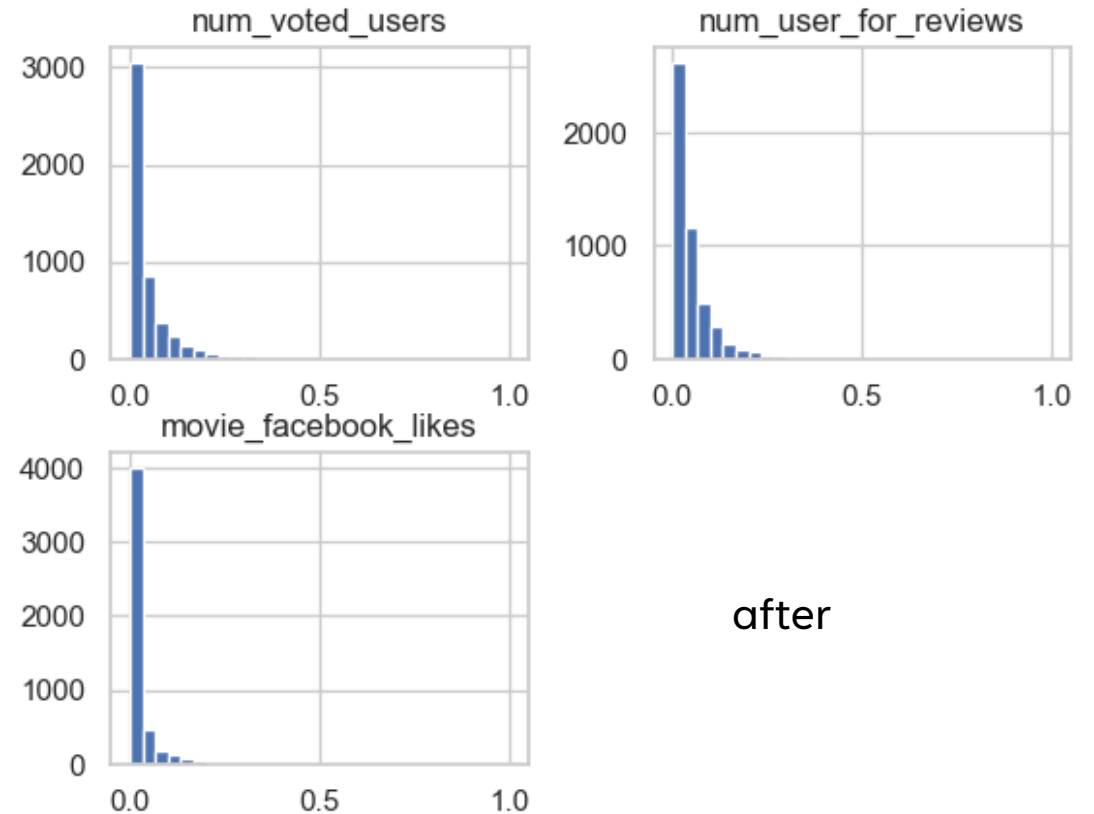| COLUMN NAME | NUMBER OF EMPTY CELLS |
|---:|---|
| num_critic_for_reviews | 50 |
| num_voted_users | 0 |
| num_user_for_reviews | 21 |
| movie_facebook_likes | 0 |

SINCE THERE ARE ONLY A FEW ROWS WITH PROBLEMS
I CHOSE TO DROP THEM

**#ROWS: 5042 -> 4993 (0.97% LOST)**

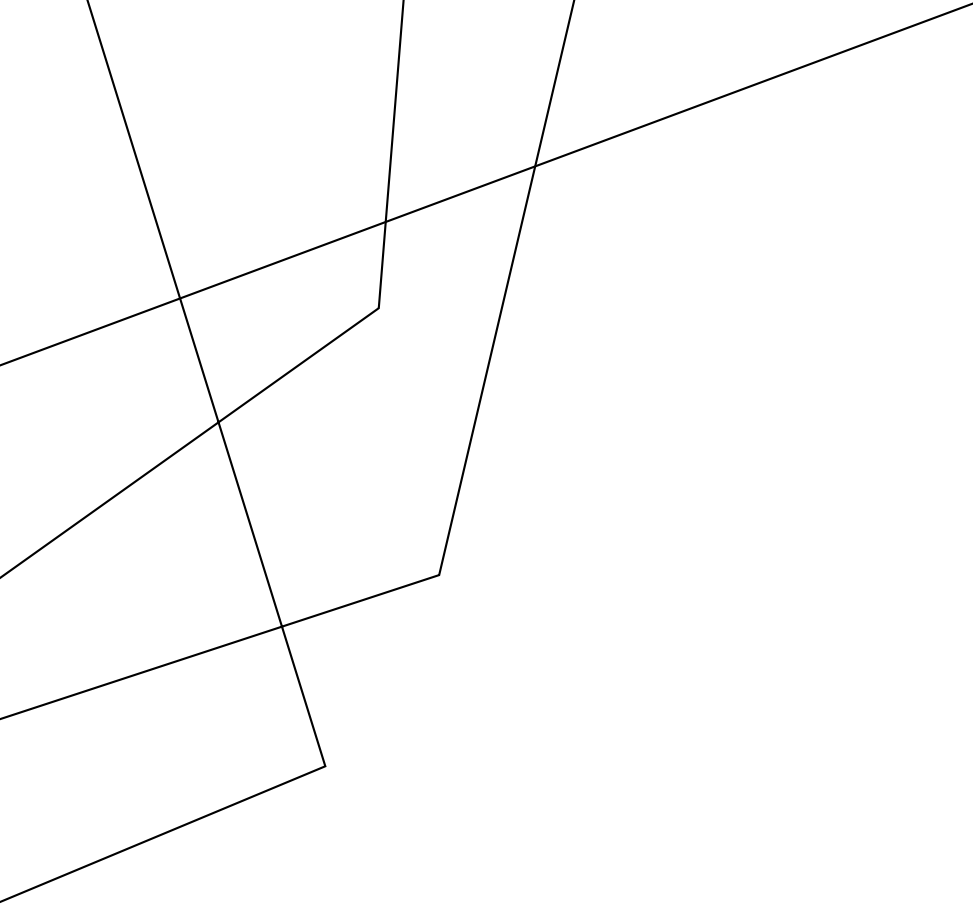# SCALE THE FEATURES: MINMAX SCALER (0~1)



before

after

# SPLIT THE DATA INTO TRAINING AN D TESTING SET

## 70% TRAINING

## 30% TESTING

## SHUFFLE & RANDOMIZED

SGDR model

Linear regression model
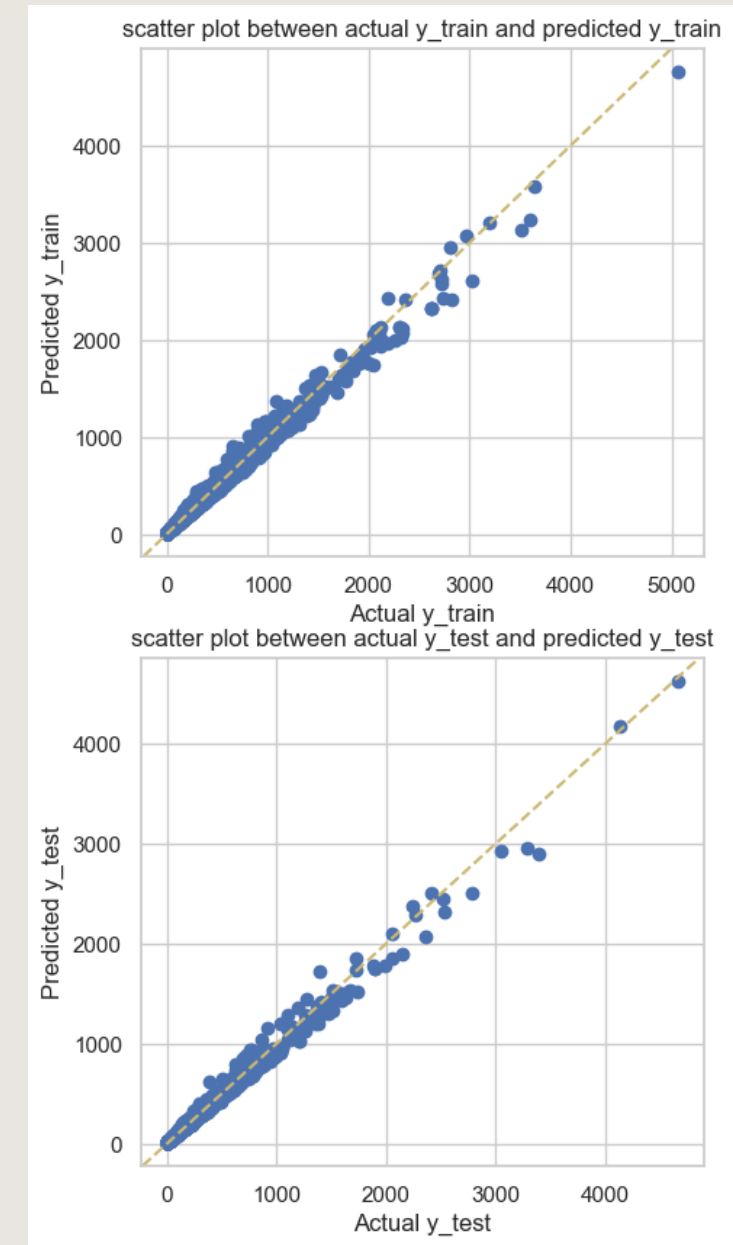
SGDR model**

**only two feature columns

# III: MODELS

# SGDREGRESSOR MODEL

- MAE : 21.66

- MSE : 1842.16

- RMSE : 42.92

- MAPE : 0.2946

- R2 : 0.9883

- Adjusted R2 : 0.9882

| features | coefficents |
|---|---|
| num_voted_users | 781 |
| num_user_for_reviews | 4183 |
| movie_facebook_likes | -122 |
| INTERCEPT | 12 |

scatter plot between actual y_train and predicted y_train
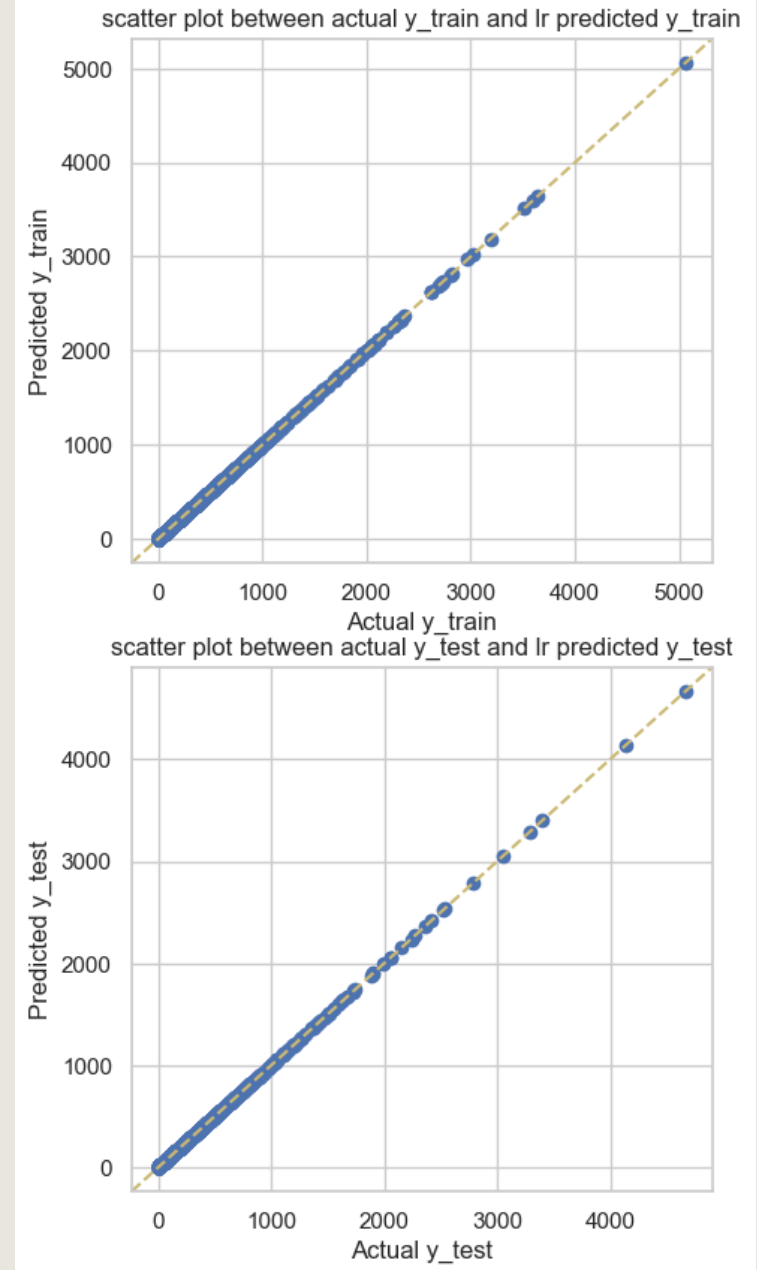
scatter plot between actual y_test and predicted y_test

# LINEAR REGRESSION MODEL

- MAE : 0

- MSE : 0

- RMSE : 0

- MAPE : 0

- R2 : 1

- Adjusted R2 : 1

Overfitting !

| features | coefficents |
|---|---|
| num_voted_users | 0 |
| num_user_for_reviews | 5059 |
| movie_facebook_likes | 0 |
| INTERCEPT | 1 |



scatter plot between actual y_train and lr predicted y_train



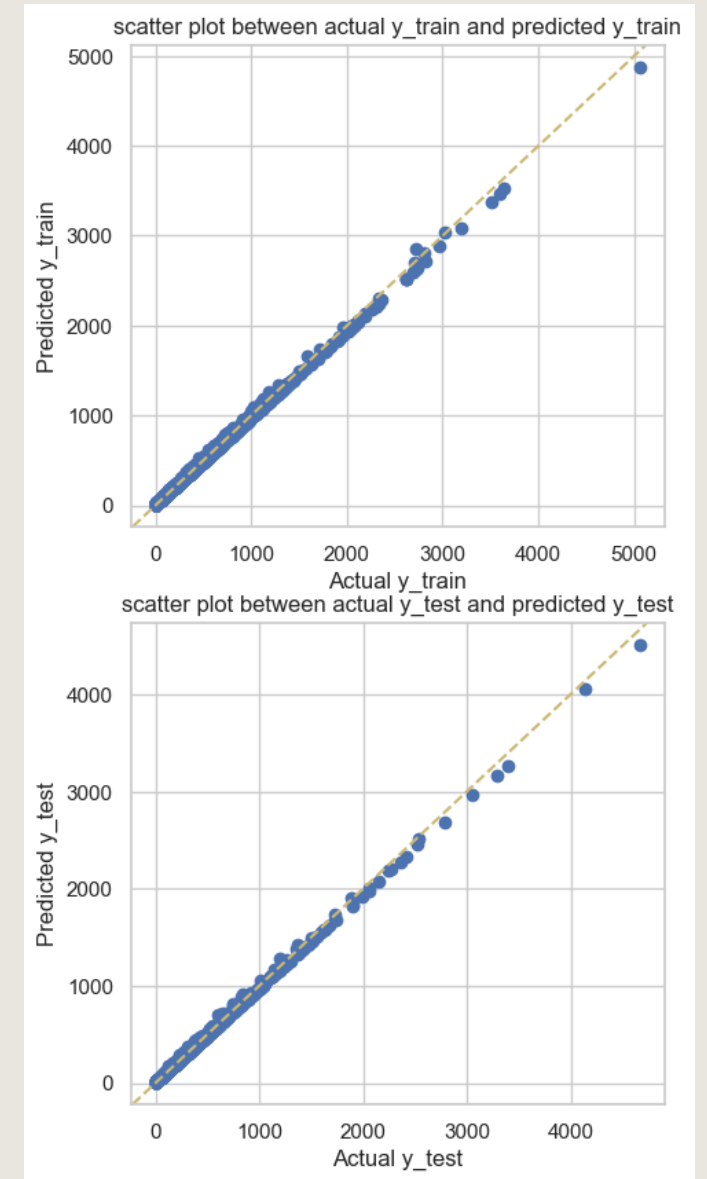scatter plot between actual y_test and lr predicted y_test

# SGDREGRESSOR MODEL
# (ONLY 2 COLUMNS**)

- MAE : 8.926

- MSE : 265.33

- RMSE : 16.29

- MAPE : 0.1633

- R2 : 0.9983

- Adjusted R2 : 0.9983

Better results
But
Possible Overfitting



| features | coefficents |
|---|---|
| num_user_for_reviews | 4856 |
| movie_facebook_likes | 227 |
| INTERCEPT | 7 |

# INTERPRETATION

- The 1$^{st}$ and 3$^{rd}$ models are candidates for practical application
  - The 1$^{st}$ model is less overfitted
  - The 3$^{rd}$ model performs better
- For every model, the features' importance is ranked:
  1. Number of users for reviews
  2. Movie Facebook likes
  3. Number of voted users