

# **Linear Regression Model for Critical Reviews**

Puvit Pracharktam 6031830321

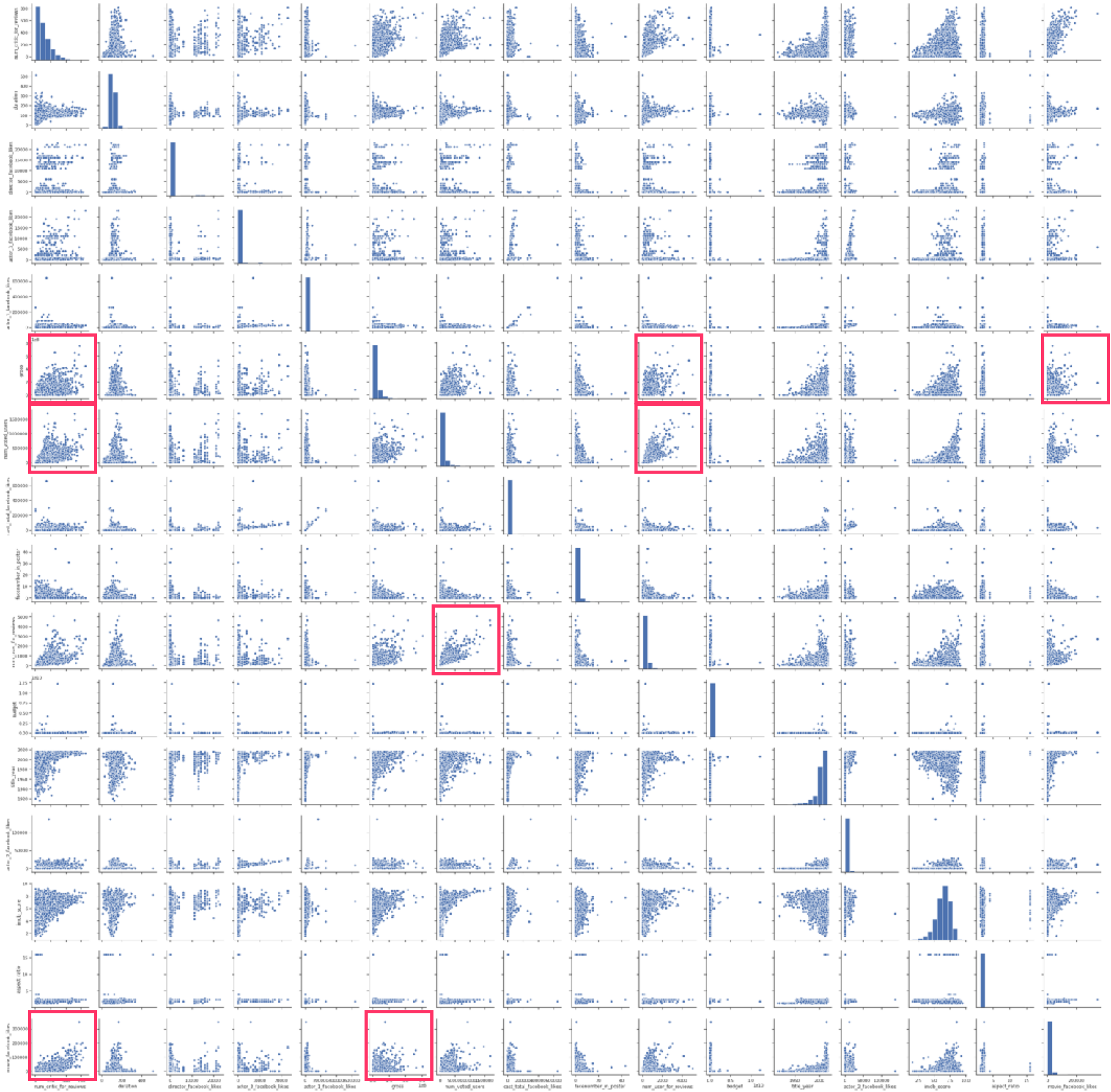
# PROCEDURE

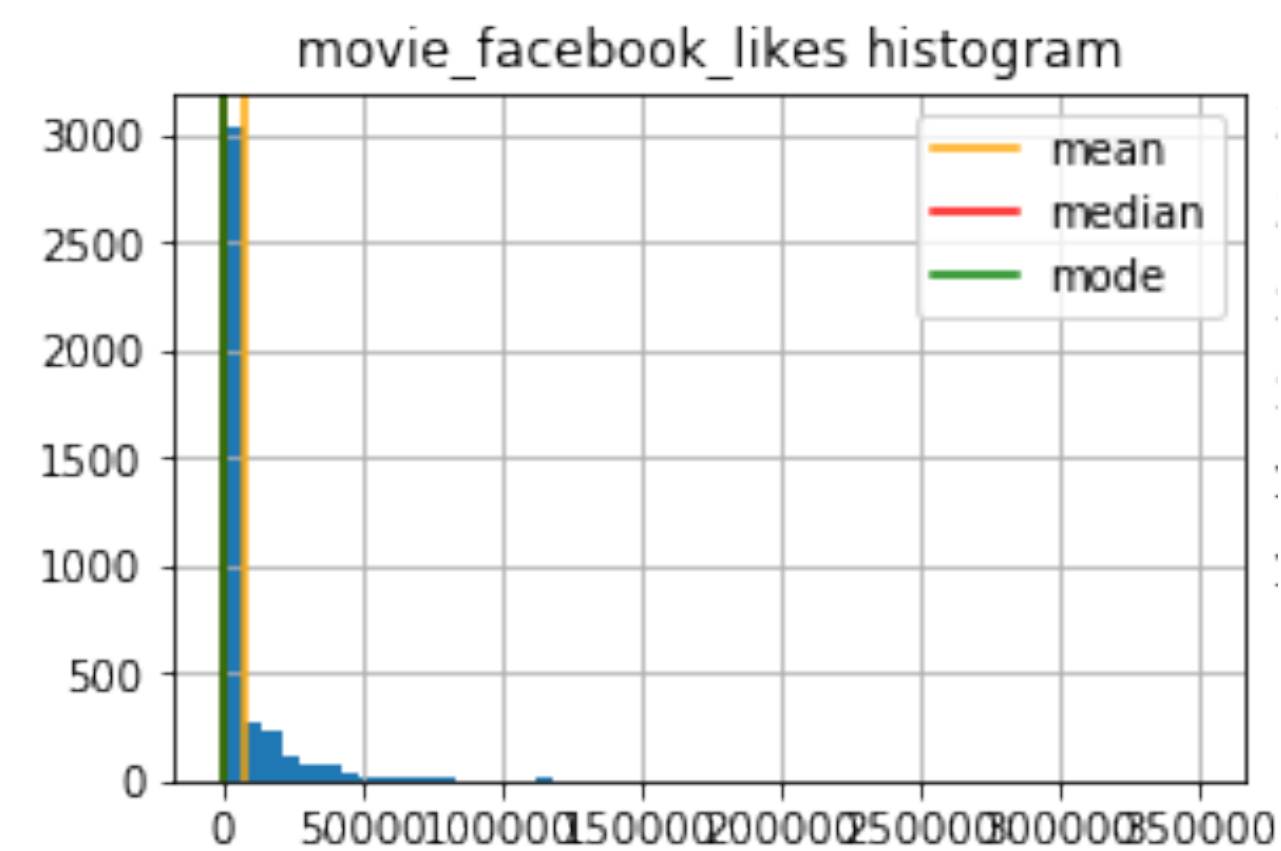
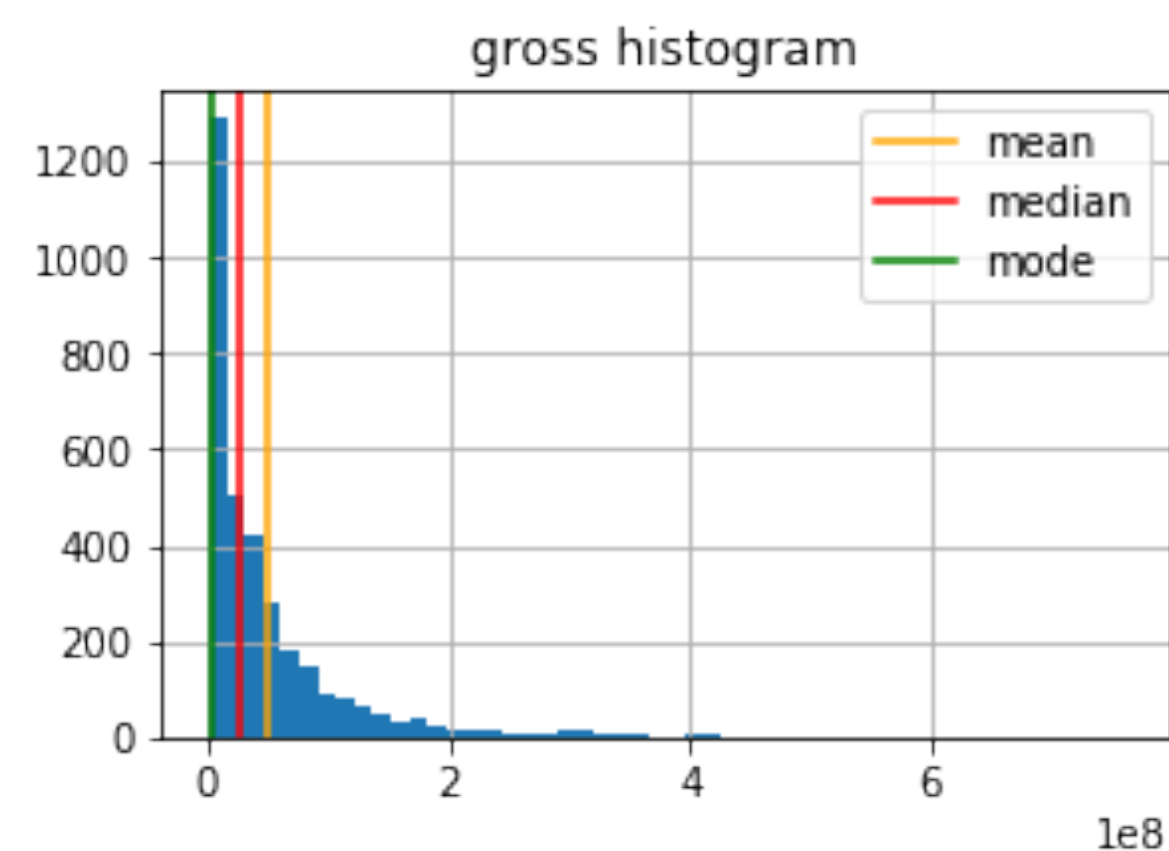
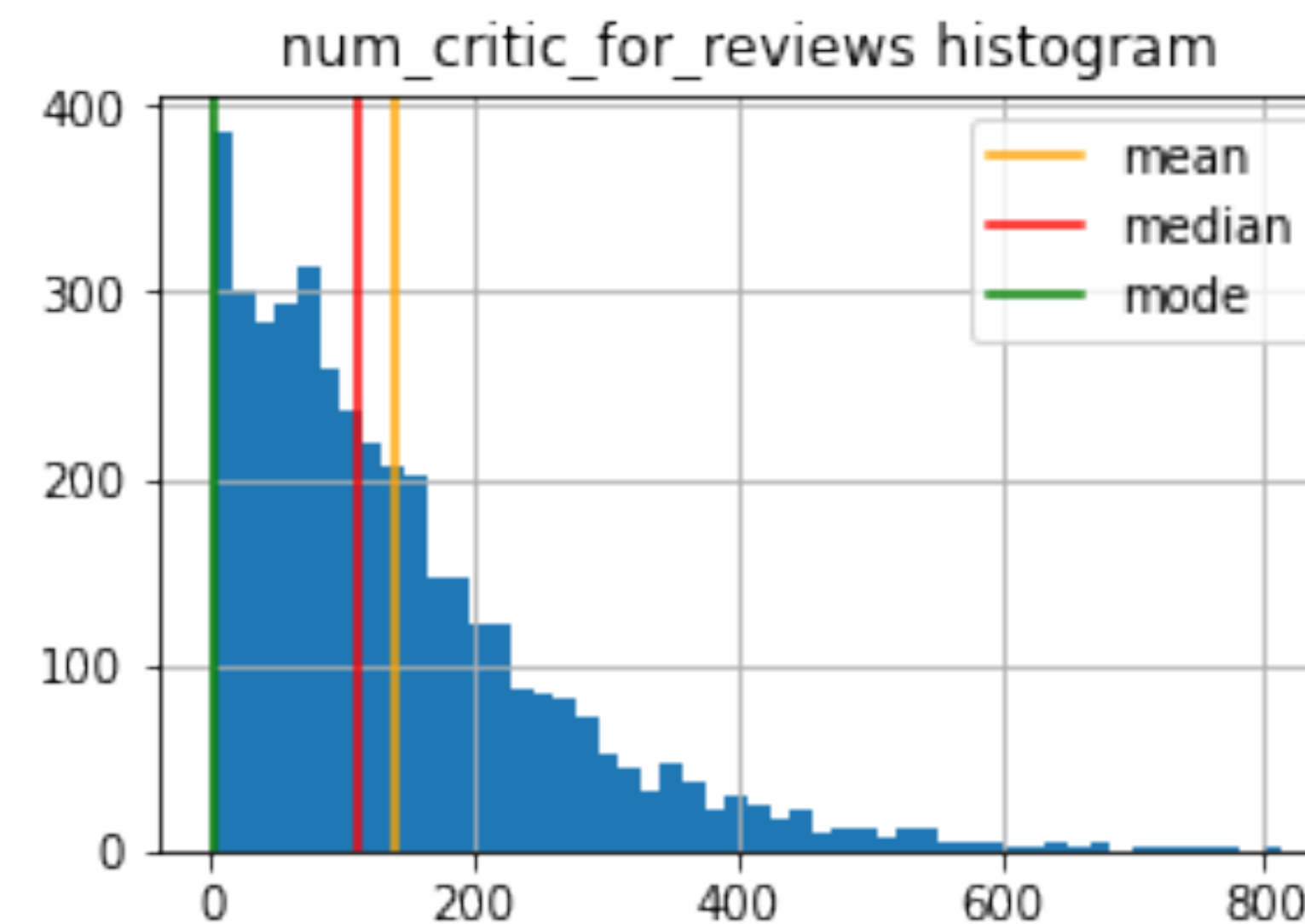
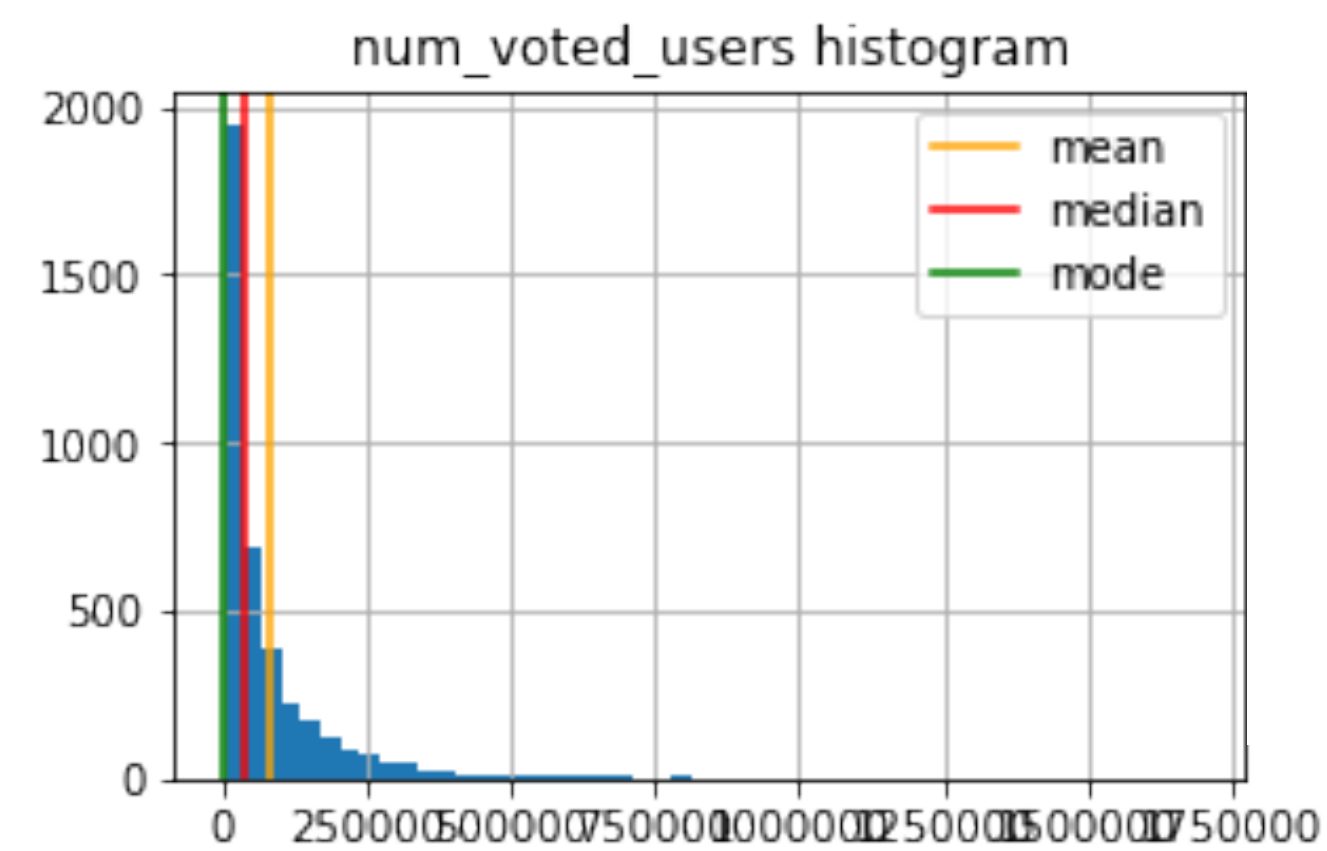
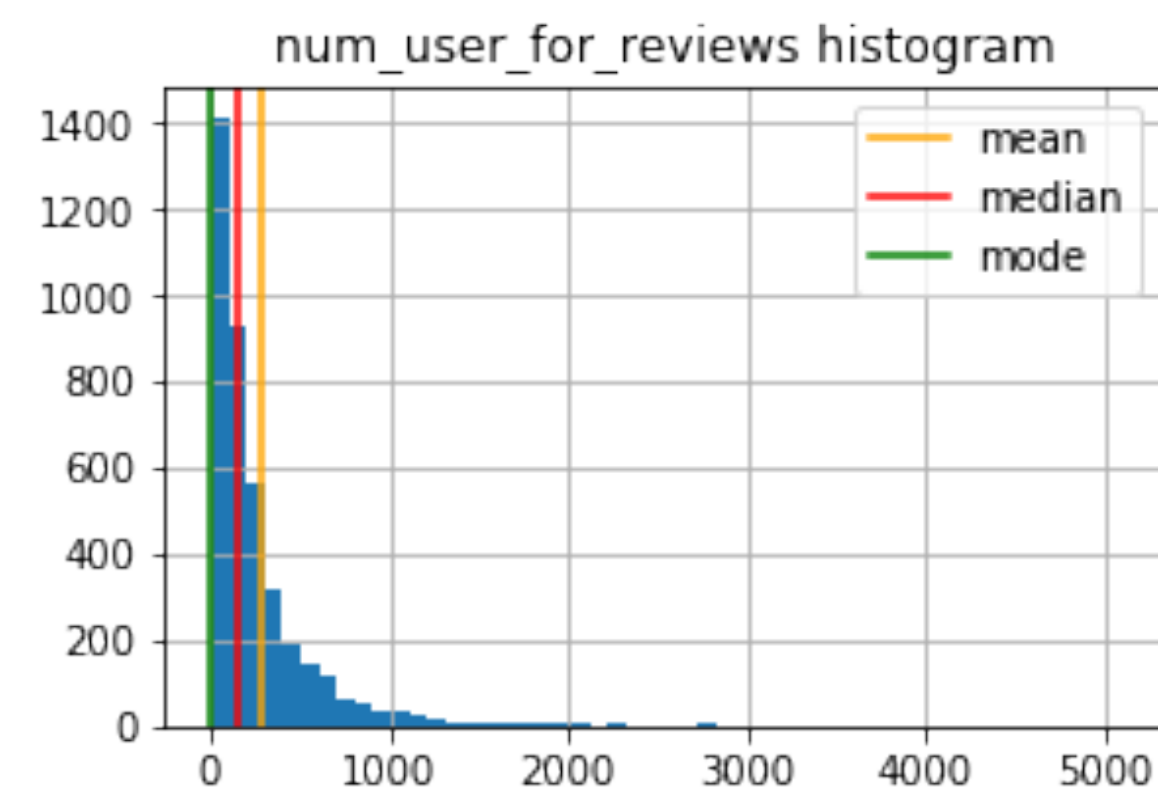
- Select all numerical features
- Cleaning data by **drop >3 NaN** in each column and **replace median**
- Scaling to log scale by both Min-Max normalisation and standardisation method
- Correlation cutout at **0.6**
- Divide **30%** of data for testing

# Scatter plot between Number of Critical Reviews

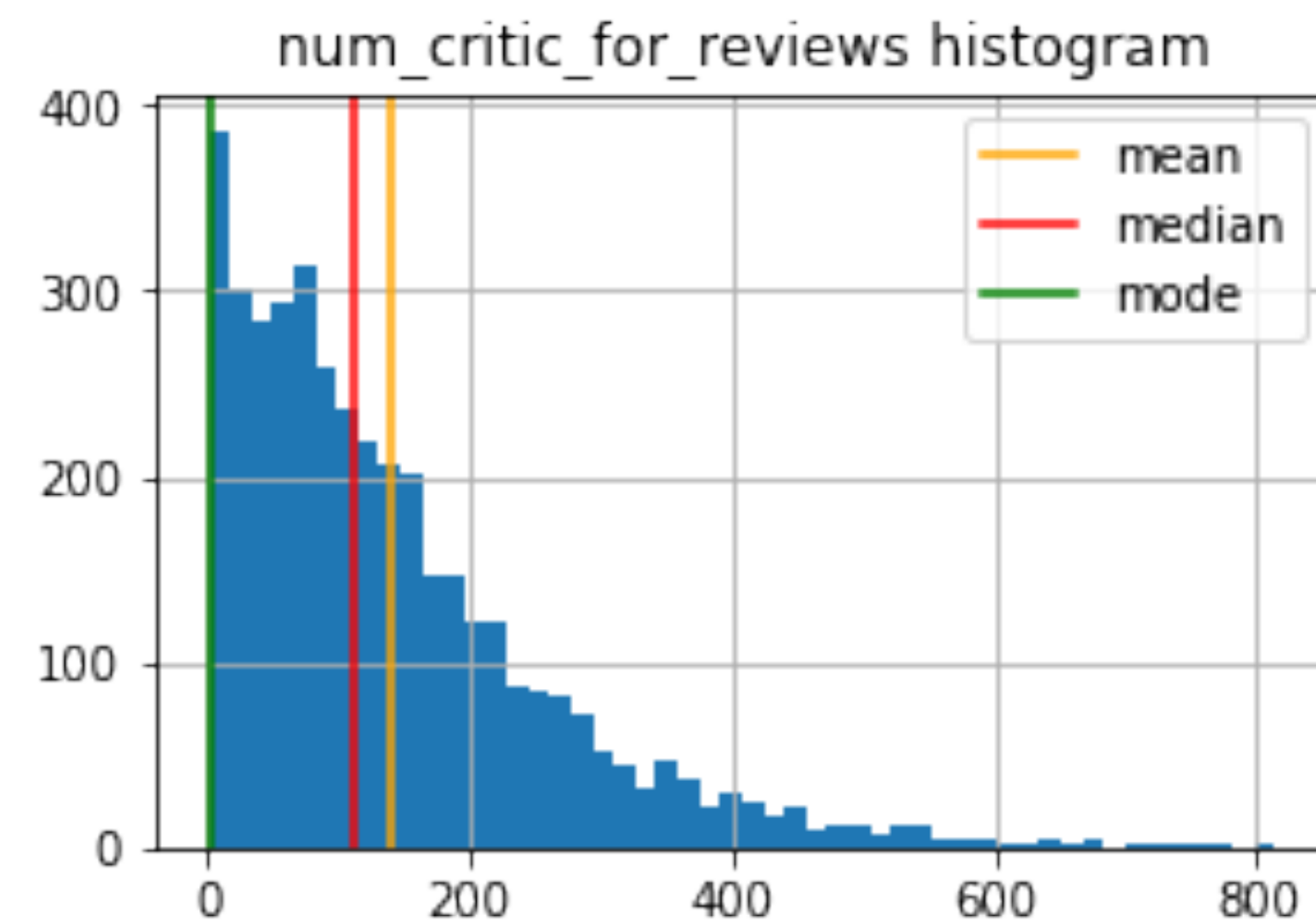
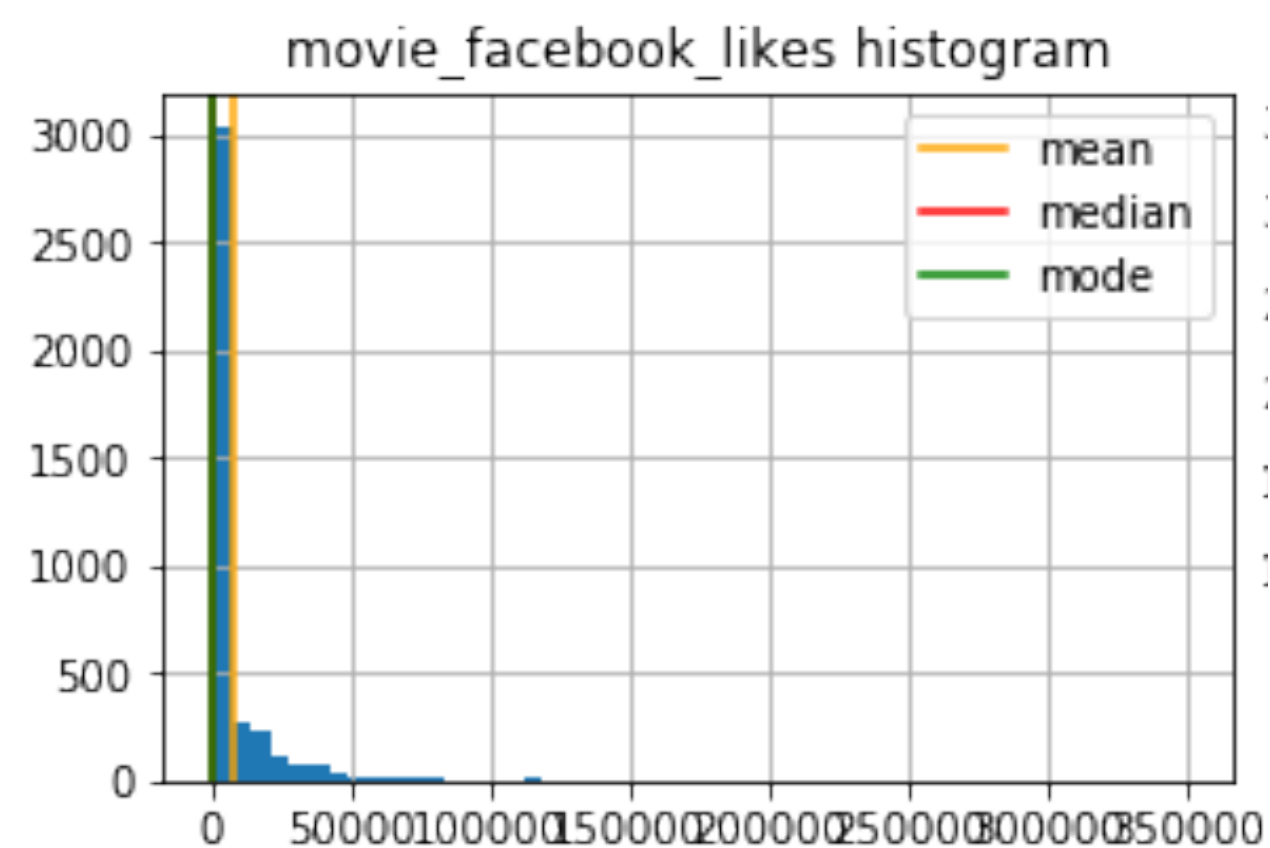
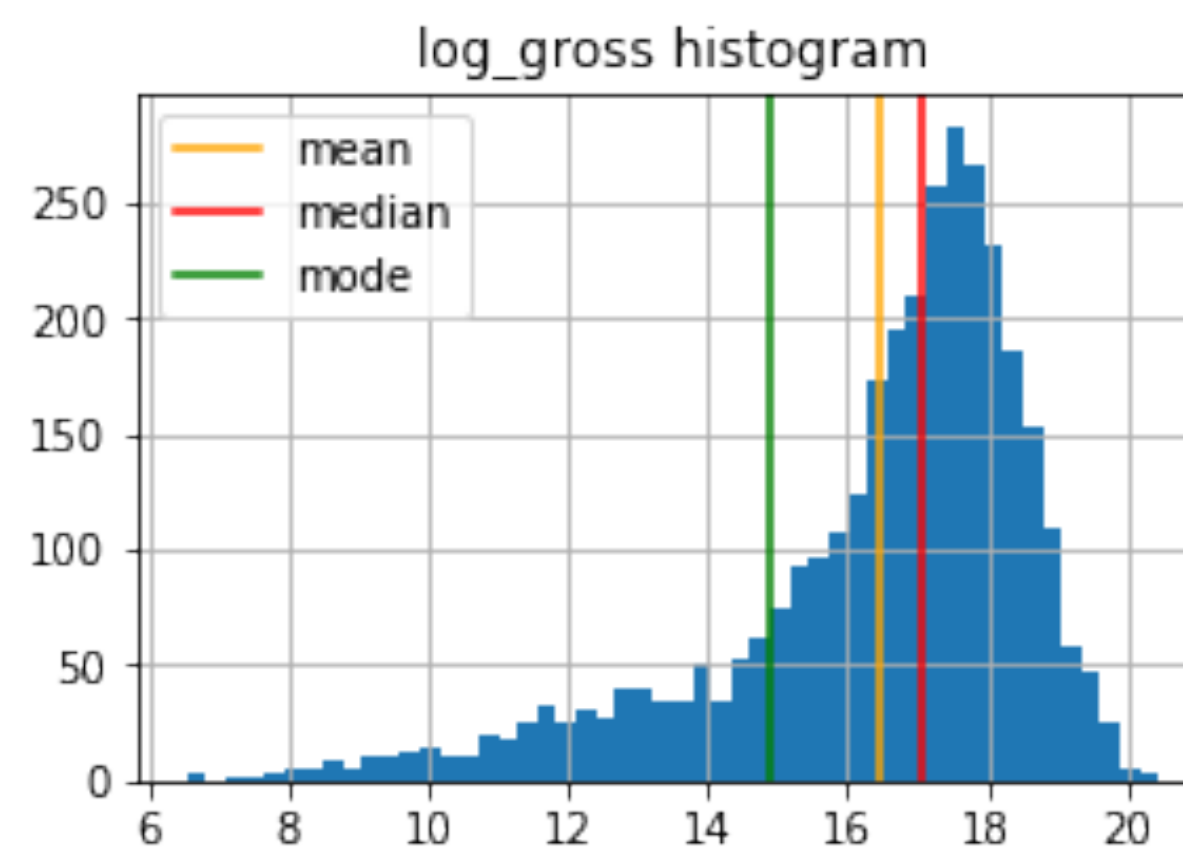
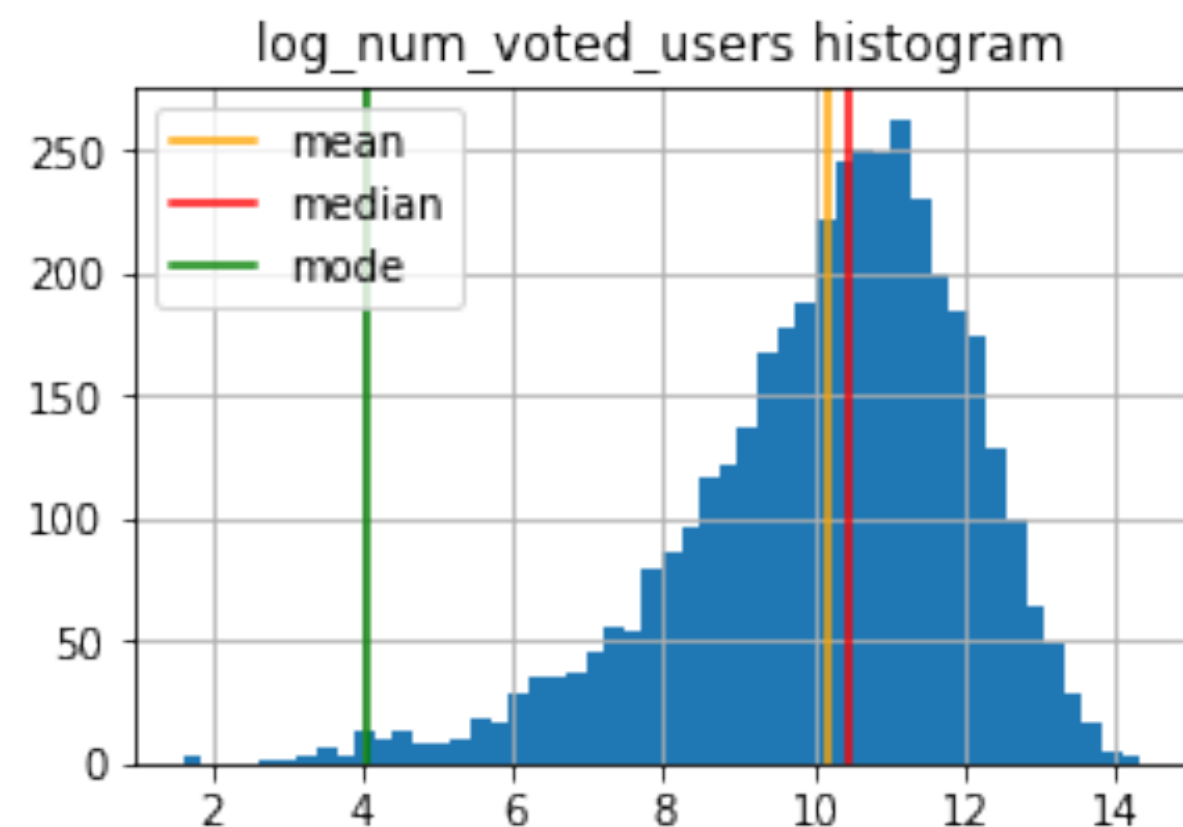
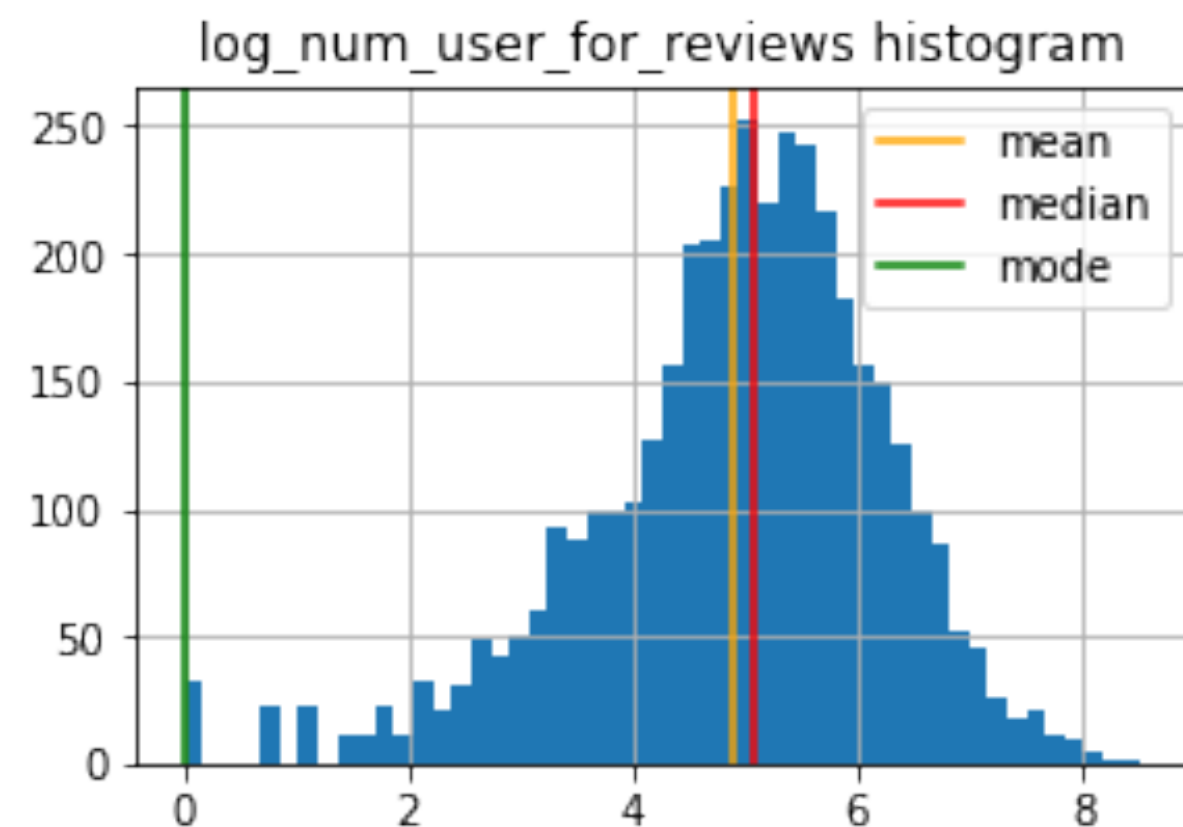
There are **4** variables  
show possible linear  
correlation

gross  
num\_voted\_users  
num\_users\_for\_reviews  
movie\_facebook\_likes





**Before taking log scaling**



**After taking log scaling**

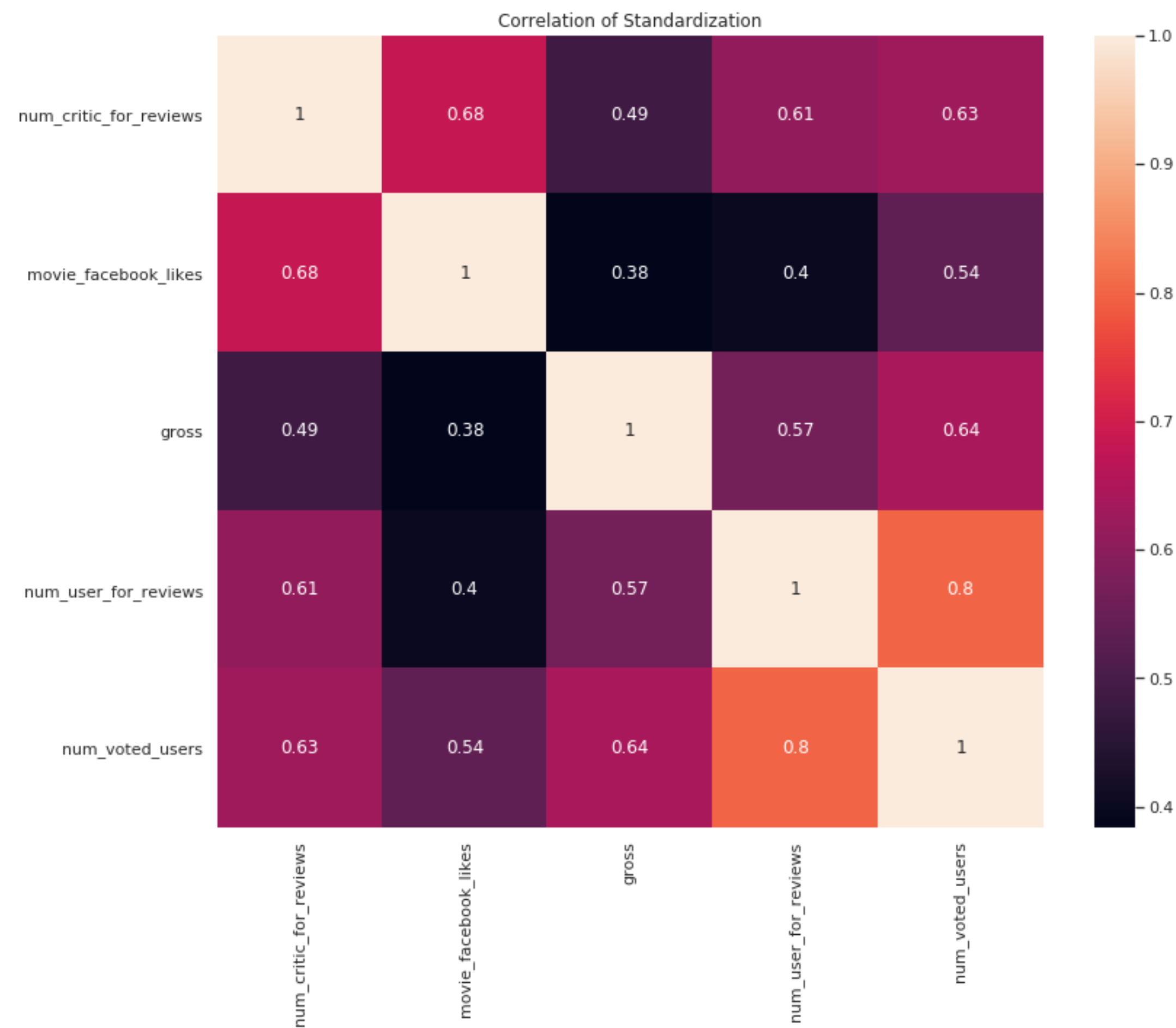
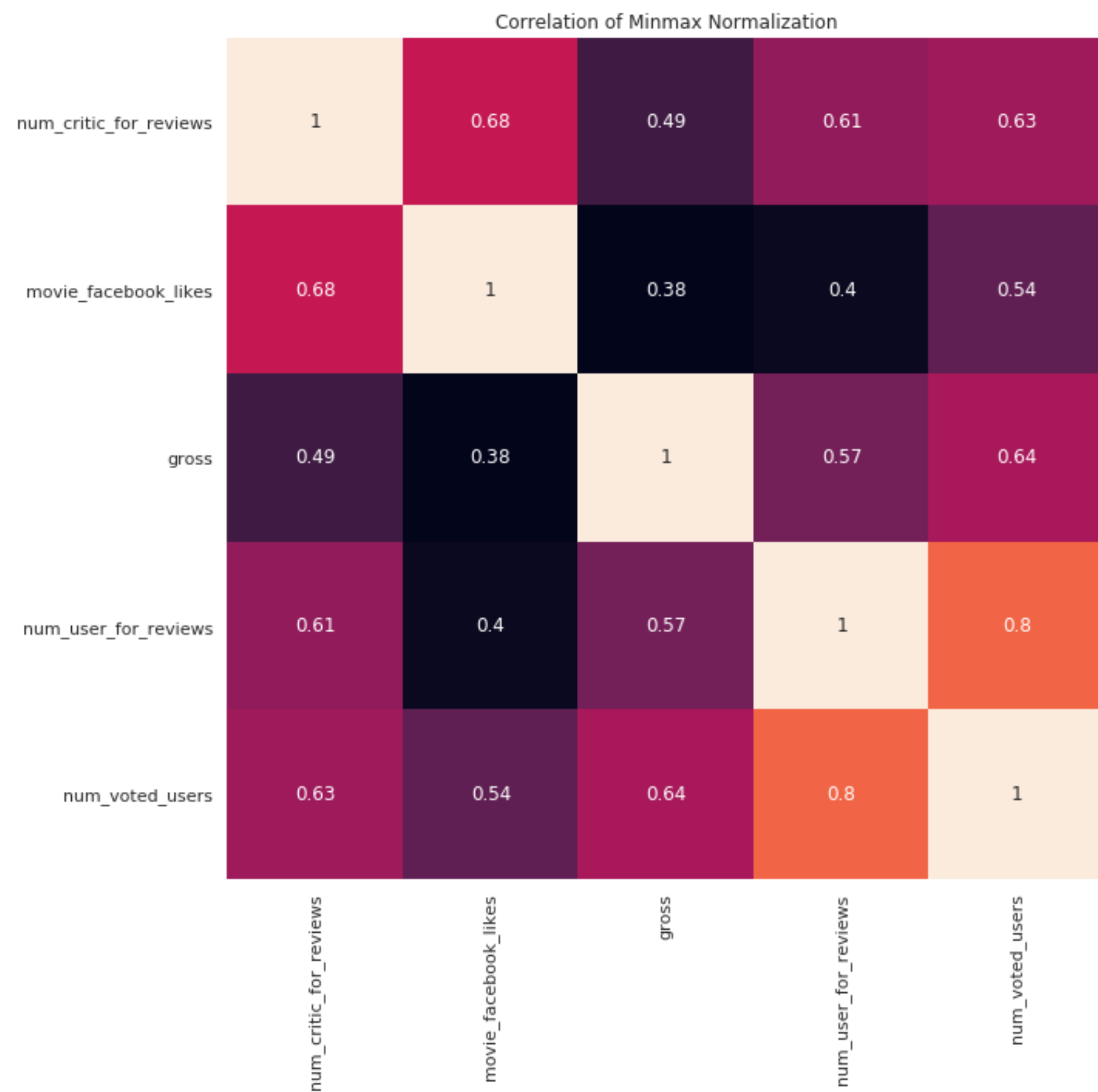


Judging from the feature's statistic, it is safe to **fill all NaN with median.**

I cannot taking log scale in **num\_critic\_for\_reviews** and **movies\_facebook\_likes** because there are some value might be **-infinity.**

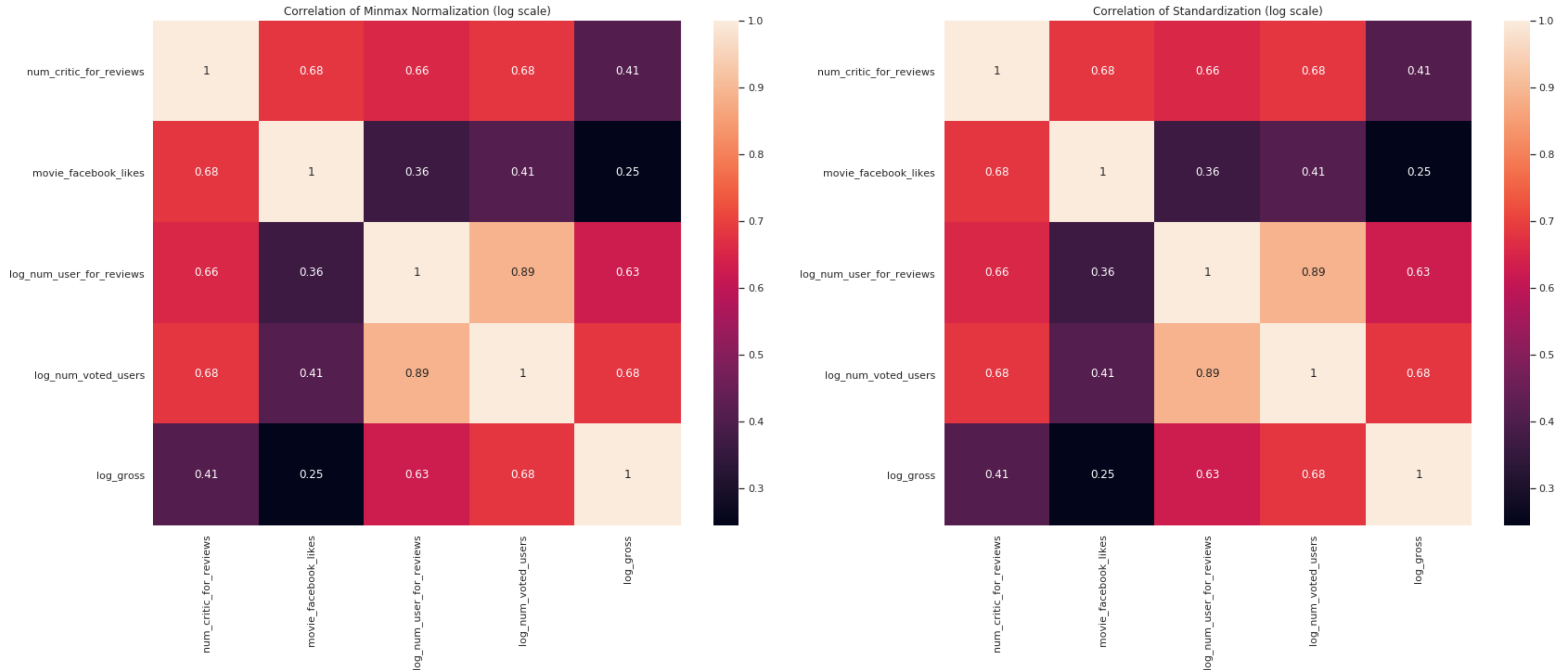
# CORRELATION ANALYSIS

Not scaled



# CORRELATION ANALYSIS

Scaled by log



From comparison, show that we should use log scale before training



Selected features correlated with **num\_critic\_for\_reviews** >0.6

		log_num_user_for_reviews	movie_facebook_likes	log_num_voted_users
log_num_user_for_reviews	0.656187	log_num_user_for_reviews	1.000000	0.363553
log_num_voted_users	0.682921			
movie_facebook_likes	0.683318			
		movie_facebook_likes	0.363553	1.000000
		log_num_voted_users	0.888058	0.411220
				1.000000

**gross** feature is eliminated since correlation <0.6

## Linear Regression Model Results (Normalized)

**With num\_voted\_users, movie\_facebook\_likes (cutout at 0.6):**

num\_critic\_for\_reviews = 0.5134 \* num\_voted\_users  
+ 1.4913 \* movie\_facebook\_likes  
+ 0.0957

(R<sup>2</sup> = 0.7828)

(RMSE : 122.14654)