

EDA_Qiang

Qiang Fang

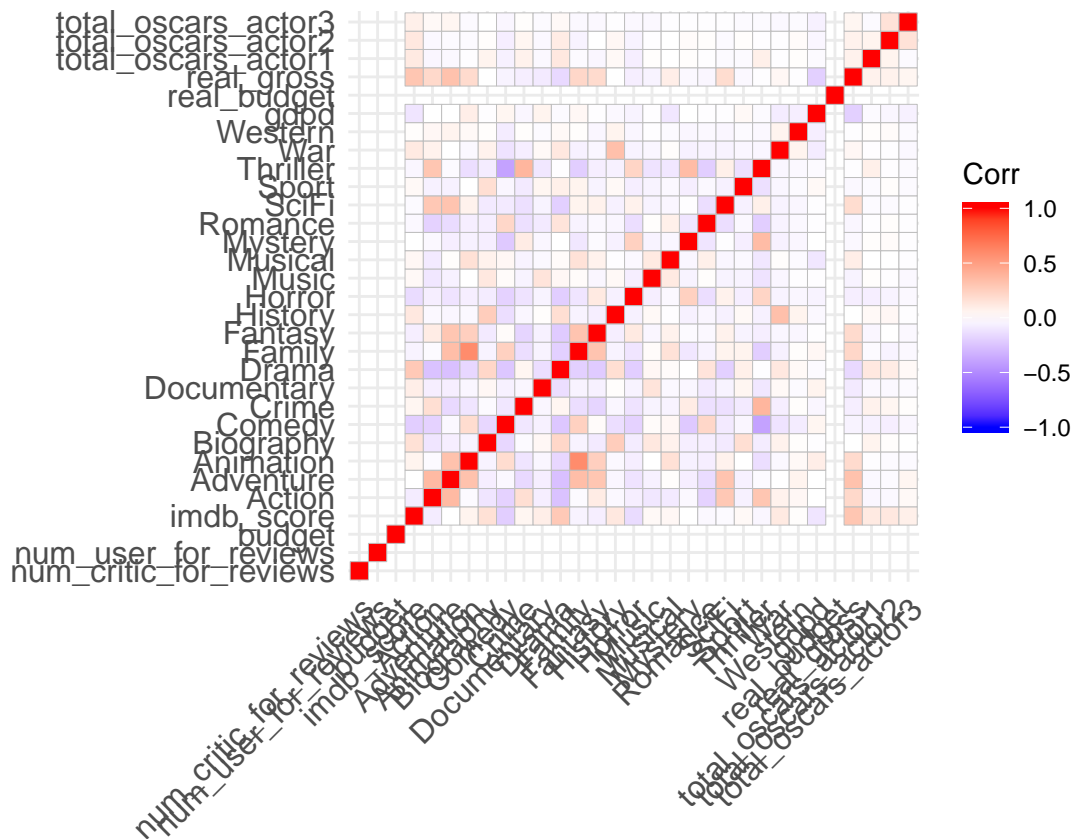
March 17, 2019

```
#load("D:/academic/DS 5110 Introduction to Data Management and Processing/project/katrina/proj_cleaned_
load(file = '~/DS5110/data/proj_cleaned_dta.RData')
```

```
library(ggcorrplot)
```

```
## Warning: package 'ggcorrplot' was built under R version 3.5.3
```

```
train_num<- select(train,num_critic_for_reviews,num_user_for_reviews,budget,imdb_score,
  Action,Adventure,Animation,Biography,Comedy,Crime,Documentary,Drama,Family,Fantasy,History,Horror,
corr <- cor(train_num)
ggcorrplot(corr)
```

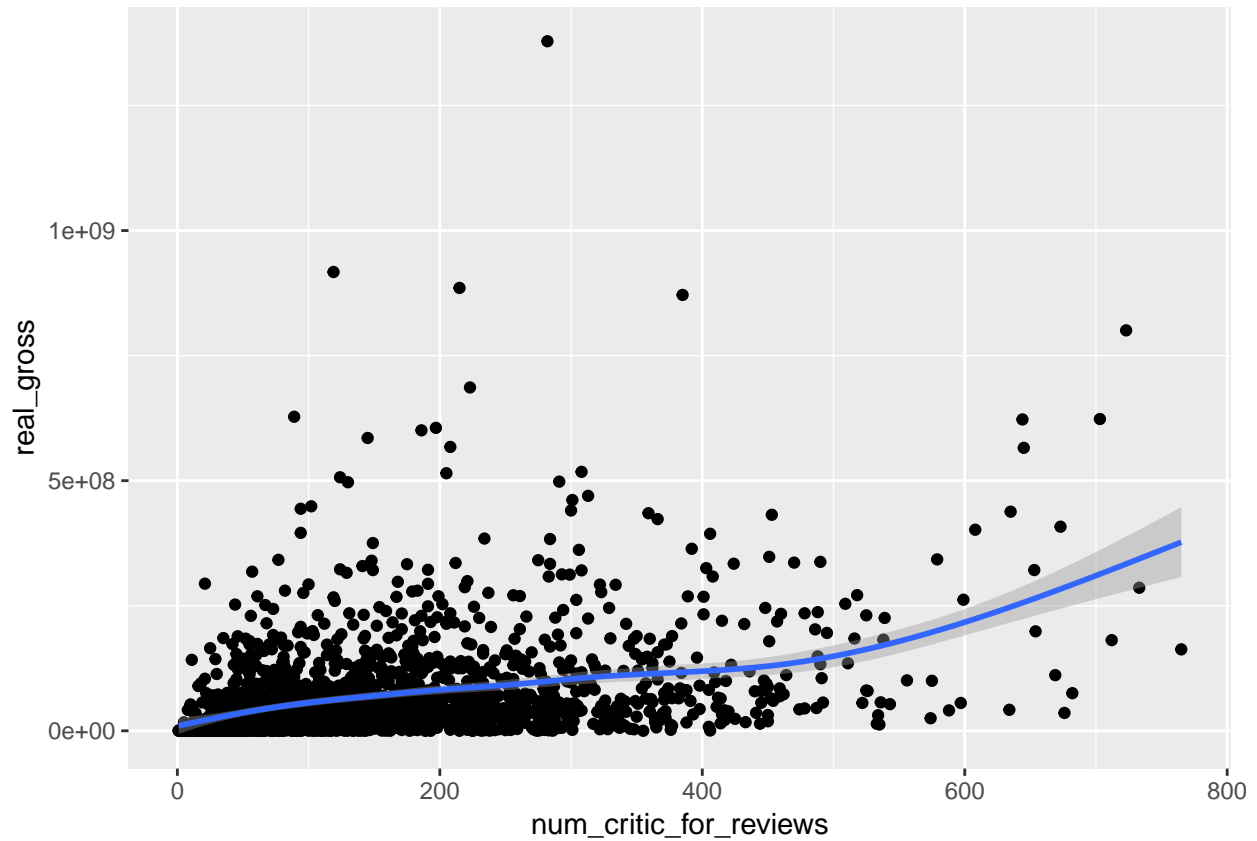


```
ggplot(train, aes(x=num_critic_for_reviews,y=real_gross)) +
  geom_point() +
  geom_smooth()
```

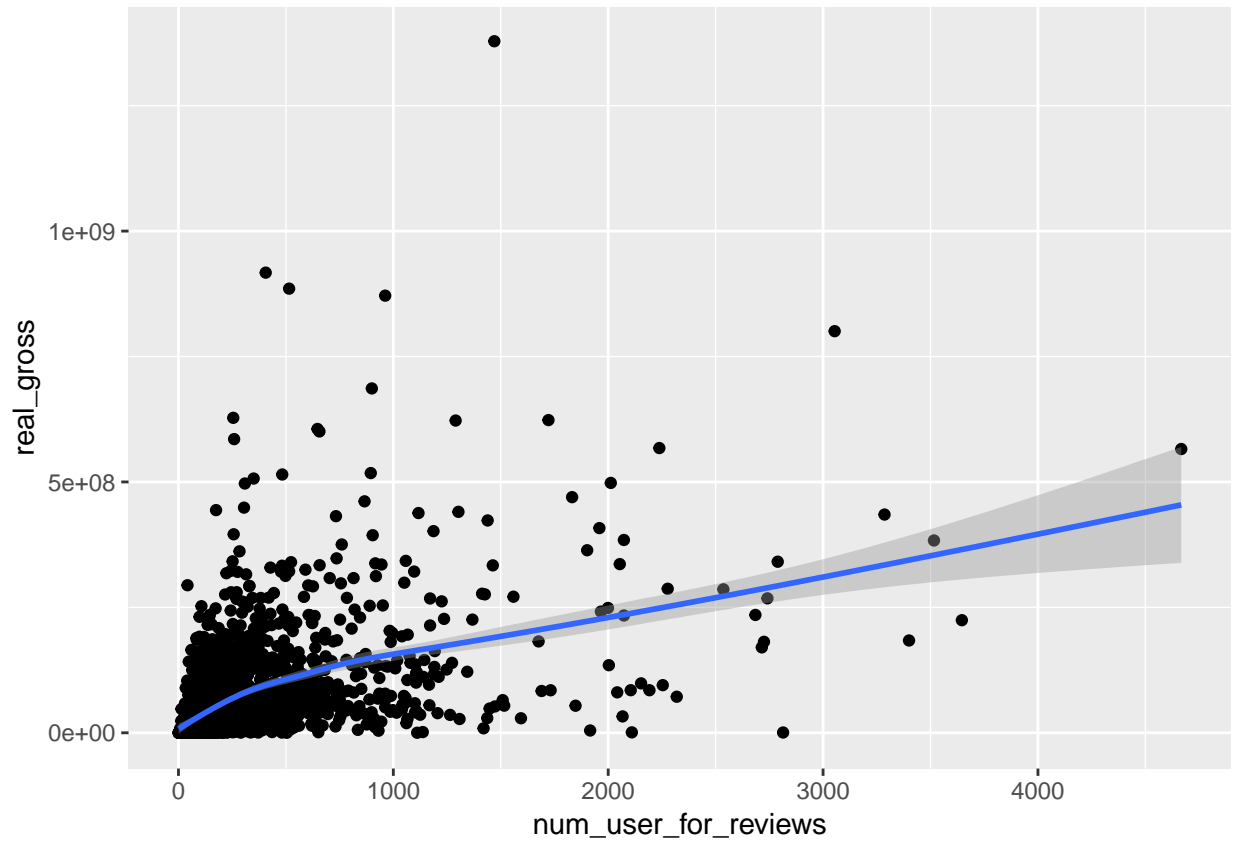
```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

```
## Warning: Removed 2 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 2 rows containing missing values (geom_point).
```

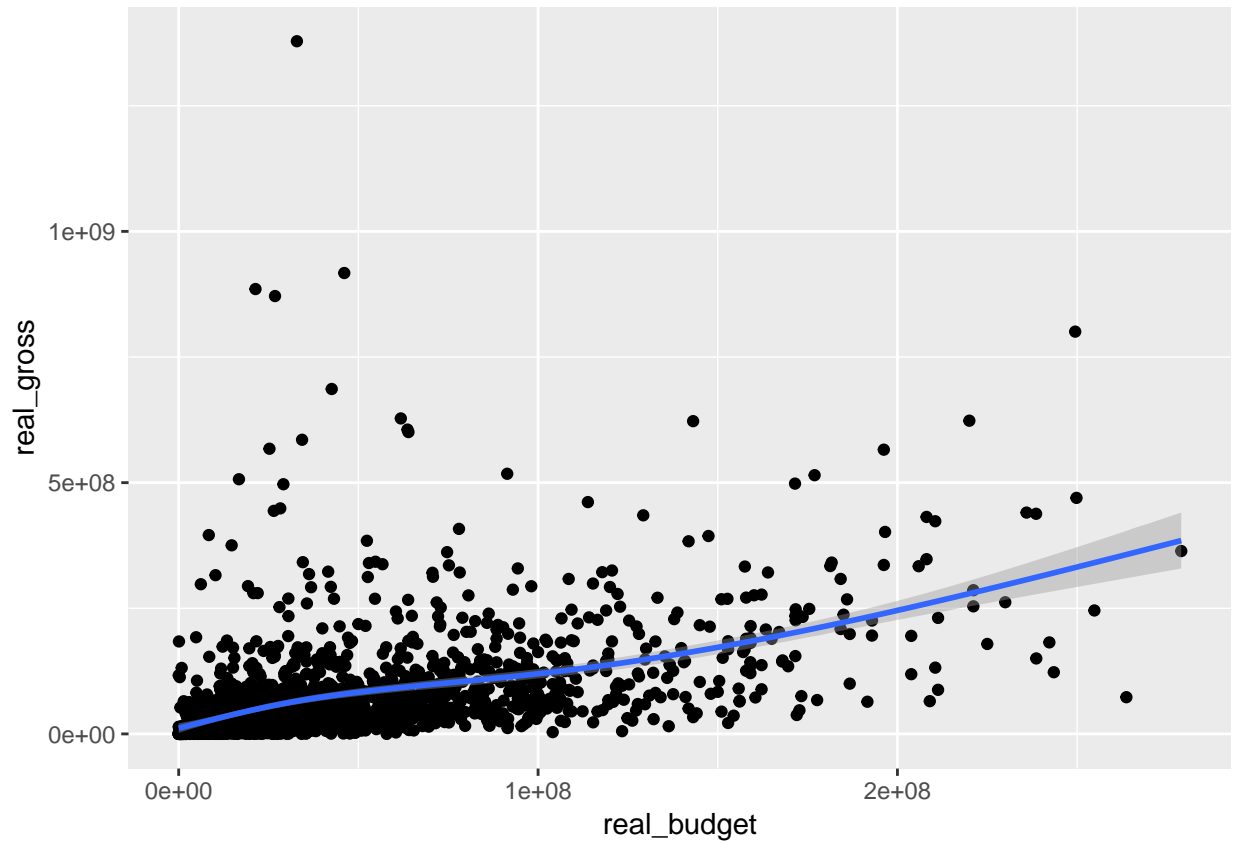


```
ggplot(train, aes(x=num_user_for_reviews,y=real_gross)) +  
  geom_point() +  
  geom_smooth()  
  
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'  
## Warning: Removed 1 rows containing non-finite values (stat_smooth).  
## Warning: Removed 1 rows containing missing values (geom_point).
```



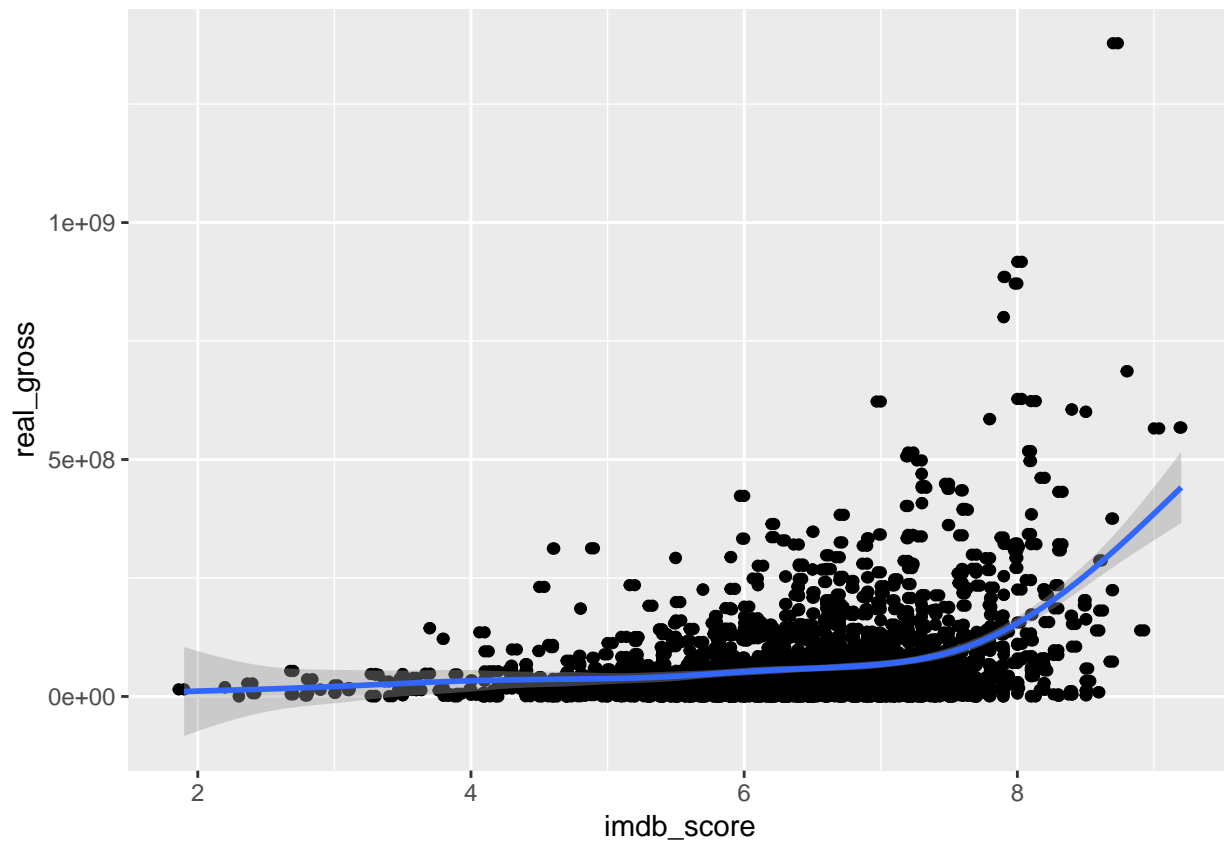
```
ggplot(train, aes(x=real_budget,y=real_gross)) +  
  geom_point() +  
  geom_smooth()
```

```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'  
## Warning: Removed 102 rows containing non-finite values (stat_smooth).  
## Warning: Removed 102 rows containing missing values (geom_point).
```



```
ggplot(train, aes(x=imdb_score,y=real_gross)) +  
  geom_point() +  
  geom_jitter() +  
  geom_smooth()
```

```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```



```
plotGenre <- function(df){
  for(i in 22:42){
    col_name <- colnames(train)[[i]]
    g <- df %>%
      group_by_(col_name) %>%
      summarize(avg_real_gross = mean(real_gross)) %>%
      ggplot(df,mapping = aes_string(x=col_name,y="avg_real_gross")) +
      geom_col()
    #print(i)
    print(g)
  }
}
plotGenre(train)
```

