

Statpadders

Toma Shigaki-Than, CJ Frederickson, Camden Reeves, Sam Kakarla

2025-03-17

```
library(tidyverse)
library(tidymodels)
library(dplyr)
library(corrplot)

imdb_top_1000 <- read_csv("data/imdb_top_1000.csv") |>
  drop_na()
```

To begin our EDA, we first had to deal with the NA values in our data. Some observations had NA values in their Gross Revenues. After examining these observations, there were no discernible patterns or connections between the NA values; they were random. As such, we were able to drop these values without compromising our data set or losing important observations.

```
imdb_top_1000 <- imdb_top_1000 |>
  mutate(no_votes_scaled = No_of_Votes / 10^6,
         gross_scaled = Gross / 10^6,
         IMDB_scaled = IMDB_Rating * 10,
         Released_Year = if_else(Series_Title == "Apollo 13", "1995",
                                Released_Year),
         difference = Meta_score - IMDB_scaled,
         runtime = as.numeric(str_remove(Runtime, " min")),
         Released_Year = as.numeric(Released_Year),
         decade = case_when(Released_Year < 1940 ~ "1930s",
                             Released_Year >= 1940 & Released_Year < 1950 ~ "1940s",
                             Released_Year >= 1950 & Released_Year < 1960 ~ "1950s",
                             Released_Year >= 1960 & Released_Year < 1970 ~ "1960s",
                             Released_Year >= 1970 & Released_Year < 1980 ~ "1970s",
                             Released_Year >= 1980 & Released_Year < 1990 ~ "1980s",
```

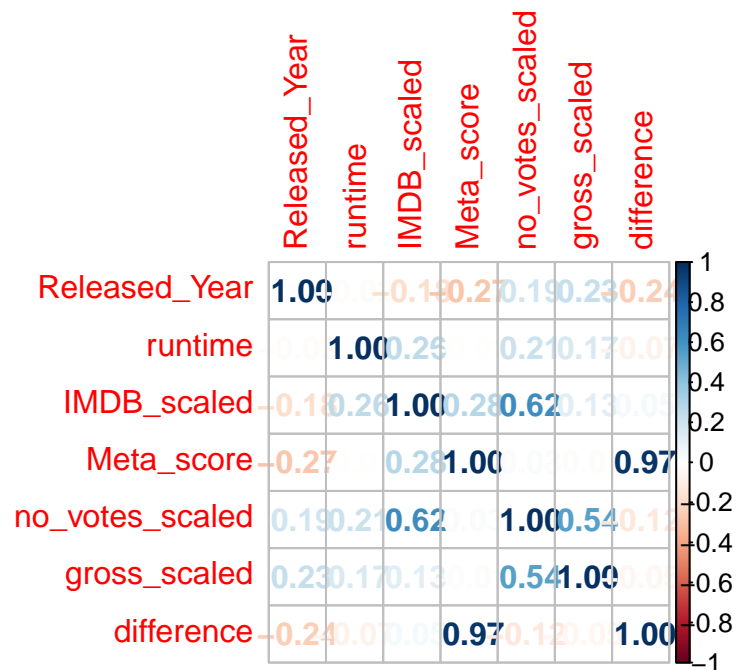
```
Released_Year >= 1990 & Released_Year < 2000 ~ "1990s",
Released_Year >= 2000 & Released_Year < 2010 ~ "2000s",
Released_Year >= 2010 ~ "2010s",))
```

Further, we turned our year into a categorical variable by creating a new variable: **decade**. Since there is a very wide range of values in **Released_Year** for the movies selected, that variable itself is not productive. Very few observations even had the same released year, and the differences between one unit in that variable were arbitrary for some movies (for example a movie released in 1966 vs 1967 does not give much insight) For data cleaning, and for better interpretability, we changed this variable into a categorical variable **decade**, where all of the years released are divided into decades (i.e. 1950s, 1960s, etc.) This gives better interpretability for that variable.

Additionally, the variable **Runtime** gave the runtime of each movie in the following template: “(number) mins”. As such, it was a categorical variable. We changed this to remove the “mins” label, and make the runtime into a numerical variable, which only shows the number of minutes in that movie’s runtime.

```
#geeksforgeeks.org/correlation-matrix-in-r-programming

matrix <- imdb_top_1000 %>%
  select(Released_Year, runtime, IMDB_scaled, Meta_score, no_votes_scaled,
         gross_scaled, difference)
c <- cor(matrix)
corrplot(c, method = "number")
```



Here, we created a correlation matrix to see which of our numerical predictors may be highly correlated, thus indicating potential multicollinearity.

! Important

Before you submit, make sure your code chunks are turned off with `echo: false` and there are no warnings or messages with `warning: false` and `message: false` in the YAML.