

# Adrian's

## Initializing Packages

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0
```

```
## v ggplot2 3.3.0    v purrr  0.3.4
## v tibble  3.0.1    v dplyr  1.0.1
## v tidyr   1.0.2    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.5.0
```

```
## -- Conflicts ----- tidyverse_conflict_1.3.0
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(ggplot2)
library(dplyr)
library(readr)
library(broom)
library(highcharter)
```

```
## Registered S3 method overwritten by 'quantmod':
##   method      from
##   as.zoo.data.frame zoo
```

```
library(readxl)
library(lubridate)
```

```
##
## Attaching package: 'lubridate'
##
## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union
```

```
library(magrittr)
```

```
##
## Attaching package: 'magrittr'
##
## The following object is masked from 'package:purrr':
##
##   set_names
##
## The following object is masked from 'package:tidyr':
##
##   extract
```

```
library(stargazer)
```

```
##
## Please cite as:
```

## Hlavac, Marek (2018). stargazer: Well-Formatted Regression and Summary Statistics Tables.

## R package version 5.2.2. <https://CRAN.R-project.org/package=stargazer>

Dataset

```
options(scipen = 999)
box_office <- read_excel("data/Recruitment x Box Office x Q1 2019.xlsx", sheet = 1)
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4266 / R4266C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4267 / R4267C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4268 / R4268C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4269 / R4269C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4270 / R4270C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4271 / R4271C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4272 / R4272C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4273 / R4273C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4274 / R4274C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4275 / R4275C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4276 / R4276C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4277 / R4277C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4278 / R4278C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4279 / R4279C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4280 / R4280C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4281 / R4281C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4287 / R4287C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4288 / R4288C1: got '-'
```

```
## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4289 / R4289C1: got '-'
```



[illegible]





[illegible]

[illegible]



```

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4433 / R4433C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4434 / R4434C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4435 / R4435C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4436 / R4436C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4437 / R4437C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4438 / R4438C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4439 / R4439C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4440 / R4440C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4441 / R4441C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4442 / R4442C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4443 / R4443C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4444 / R4444C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4445 / R4445C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4446 / R4446C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4447 / R4447C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4448 / R4448C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4449 / R4449C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4450 / R4450C1: got '-'

## Warning in read_fun(path = enc2native(normalizePath(path)), sheet_i = sheet, :
## Expecting numeric in A4451 / R4451C1: got '-'

title_performance <- read_excel("data/Recruitment x Box Office x Q1 2019.xlsx", sheet = 2)
twitter <- read_csv("data/Recruitment x Box Office x Q1 2019 - Twitter.csv", sep="\t", header=TRUE, fil

```

Data Cleaning

```

box_office$Title <- as.factor(box_office$Title) #classifying titles into factors
box_office$Studio <- as.factor(box_office$Studio) #classifying studios into factors
box_office$Year <- as.factor(box_office$Year) #classifying years into factors
box_office$`Weekend Gross` <- str_remove_all(box_office$`Weekend Gross`, "[AU$,]") #removing obstructive
box_office$`Weekend Gross` <- as.numeric(box_office$`Weekend Gross`) #transforming column into numeric
box_office$`This Week's Rank` <- str_remove_all(box_office$`This Week's Rank`, "N-") #removing obstructive
box_office$`Last Week's Rank` <- str_remove_all(box_office$`Last Week's Rank`, "[N-]") #removing obstructive
box_office$`Last Week's Rank` <- sub(".0+$", "", as.character(box_office$`Last Week's Rank`)) #removing
box_office$`Weeks Into Release` <- sub(".0+$", "", as.character(box_office$`Weeks Into Release`)) #removing
levels(box_office$Title) <- gsub("Dr. Seuss' The Grinch (2018)", "The Grinch", levels(box_office$Title))

title_performance$Title <- as.factor(title_performance$Title) #transforming title into factor
title_performance$Date <- gsub("(.*)-.*", "\\1", title_performance$Date) #removing obstructive strings
title_performance$Date <- str_remove_all(title_performance$Date, "[.-]") #removing obstructive strings
title_performance$`Weekend Gross` <- str_remove_all(title_performance$`Weekend Gross`, "[.0]") #removing
title_performance$`Weekend Gross` <- as.numeric(title_performance$`Weekend Gross`)

## Warning: NAs introduced by coercion

title_performance$`Week # of release` <- as.numeric(title_performance$`Week # of release`) #transforming
title_performance$`Gross-to-Date` <- as.numeric(title_performance$`Gross-to-Date`)
title_performance$new_date <- paste(title_performance$Year, title_performance$Date, sep="-") %>% ymd()
title_performance$Week <- format(title_performance$new_date, "%V") #Week of the year
levels(title_performance$Title) <- gsub("Avengers: Infinty War", "Avengers Infinity War", levels(title_performance$Title))

twitter$post_date_time <- gsub("T.*", "", twitter$post_date_time) #remove details on time
twitter$post_date_time <- ymd(twitter$post_date_time)
twitter$title <- as.factor(twitter$title)
twitter$Year <- format(twitter$post_date_time, "%G") #creating a year column
twitter$Week <- format(twitter$post_date_time, "%V") #creating a week of the year column
levels(twitter$title) <- gsub("beauty and the beast", "Beauty and the Beast", levels(twitter$title)) #transforming
twitter$ID <- as.numeric(twitter$ID) #making column into numeric

```

Creating a New Dataframe to understand twitter statistics + merging box office sales, weeks, etc.

```

twitter_stats <- twitter %>%
  select(Title = title, ID, Year, Week) #Creating a new dataframe to work with Twitter Data
twitter_stats <- twitter_stats %>%
  group_by(Title) %>%
  mutate(cumulative_tweet_count = sum(ID))
twitter_stats <- twitter_stats %>%
  group_by(Title, Week) %>%
  mutate(weekly_tweet_count = sum(ID)) #weekly tweet count for film
twitter_stats <- unique(twitter_stats) #summarizing to show weekly features by selecting unique traits

twitter_stats_combined <- merge(twitter_stats, box_office, by.x=c("Title", "Week"), by.y=c("Title", "Week"))
twitter_stats_combined <- twitter_stats_combined %>%
  select(!Year.y) %>%
  select(!ID) %>%
  rename(Year = Year.x) #removing unnecessary columns and renaming columns

```

Looking for Trends: Visualising relationship between Weekend Gross and Number of Tweets

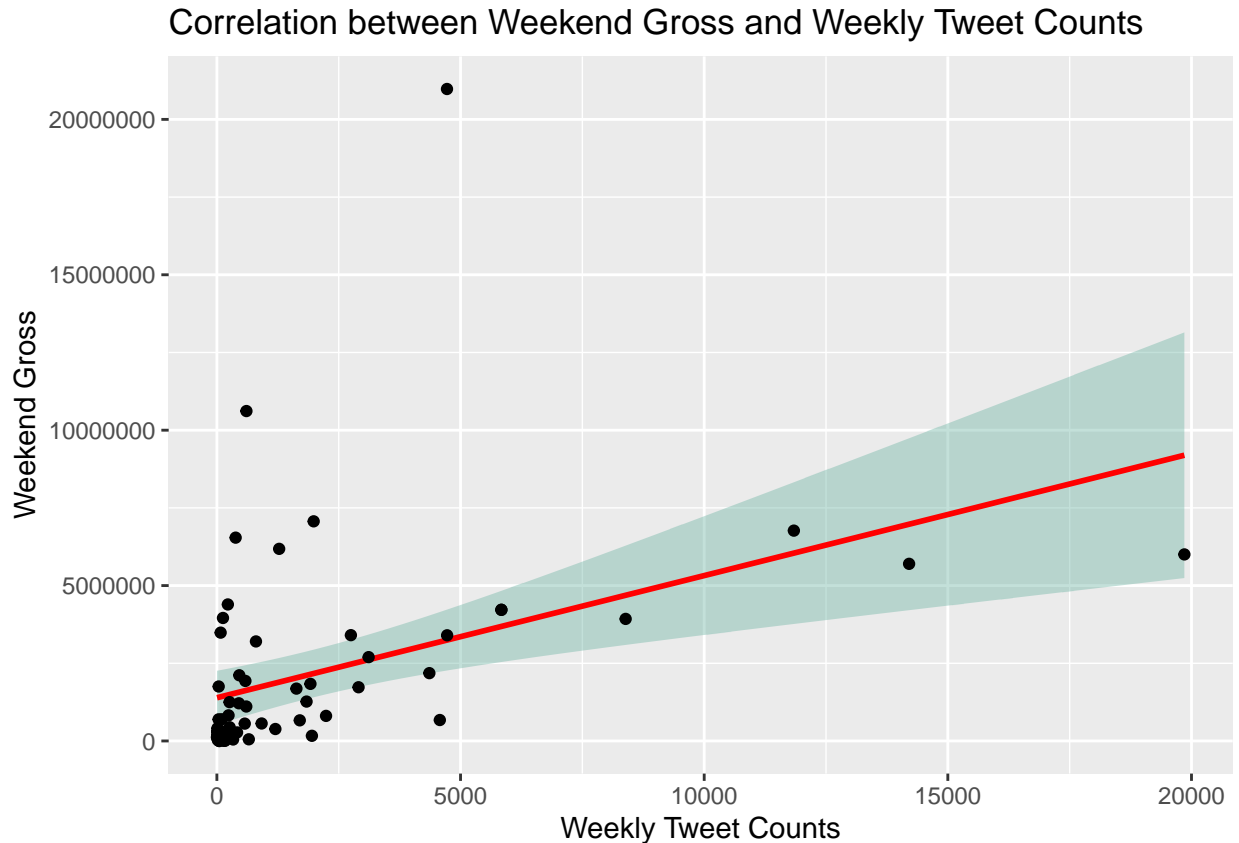
```

twitter_stats_combined %>%
  ggplot() +
  aes(y=`Weekend Gross`, x=weekly_tweet_count) +

```

```
geom_smooth(method=lm , color="red", fill="#69b3a2", se=TRUE) +
geom_point() +
labs(x="Weekly Tweet Counts", title = "Correlation between Weekend Gross and Weekly Tweet Counts")
```

```
## `geom_smooth()` using formula 'y ~ x'
```



The visualisation suggests that there seem to be a positive linear correlation between film performance (represented by “Weekend Gross”) and social buzz (represented by “Weekly Tweet Counts”). However, the error in estimate seems large as we can see that there is visible gap between data points and regression line. Moreover, there is a cluster that the line does not seem to interact with.

```
stargazer(twitter_stats_combined)
```

```
##
## % Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
## % Date and time: Wed, Aug 12, 2020 - 18:20:12
## \begin{table}[!htbp] \centering
##   \caption{}
##   \label{}
##   \begin{tabular}{@{\extracolsep{5pt}}lcccccc}
##     \hline[-1.8ex]
##     \hline \hline[-1.8ex]
##     Statistic & \multicolumn{1}{c}{N} & \multicolumn{1}{c}{Mean} & \multicolumn{1}{c}{St. Dev.} & \multicolumn{2}{c}{} \\
##     \hline \hline[-1.8ex]
##     cumulative\_tweet\_count & 63 & 24,882.590 & 24,462.610 & 110 & 2,723 & 52,530 & 62,964 \\
##     weekly\_tweet\_count & 63 & 1,891.667 & 3,575.983 & 3 & 134.5 & 1,935 & 19,854 \\
##     Weekend Gross & 63 & 2,134,783.000 & 3,326,059.000 & 0 & 179,632.5 & 3,301,566 & 20,977,461 \\
##     \hline \hline[-1.8ex]
```

```
## \end{tabular}  
## \end{table}
```