

# Analysis of Lead Gender and Box Office

by Max Woolf (@minimaxir)

*This notebook is licensed under the MIT License. If you use the code or data visualization designs contained within this notebook, it would be greatly appreciated if proper attribution is given back to this notebook and/or myself. Thanks! :)*

```
1 options(warn = -1)
2
3 # IMPORTANT: This assumes that all packages in "Rstart.R" are installed,
4 # and the fonts "Source Sans Pro" and "Open Sans Condensed Bold" are installed
5 # via extrafont. If ggplot2 charts fail to render, you may need to change/remove the theme
6   call.
7 source("Rstart.R")
8 library(outliers)
9
10 sessionInfo()
```

```
1 Attaching package: 'dplyr'
2
3 The following objects are masked from 'package:stats':
4
5   filter, lag
6
7 The following objects are masked from 'package:base':
8
9   intersect, setdiff, setequal, union
10
11 Registering fonts with R
12
13 Attaching package: 'scales'
14
15 The following objects are masked from 'package:readr':
16
17   col_factor, col_numeric
18
19
20
21
22
23
24 R version 3.2.3 (2015-12-10)
25 Platform: x86_64-apple-darwin13.4.0 (64-bit)
26 Running under: OS X 10.11.4 (El Capitan)
27
28 locale:
29 [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
30
31 attached base packages:
32 [1] grid      stats    graphics  grDevices utils      datasets  methods
33 [8] base
34
35 other attached packages:
```

```

36 [1] outliers_0.14      stringr_1.0.0      digest_0.6.8      RColorBrewer_1.1-2
37 [5] scales_0.3.0       extrafont_0.17     ggplot2_2.0.0     dplyr_0.4.3
38 [9] readr_0.1.1
39
40 loaded via a namespace (and not attached):
41 [1] Rcpp_0.12.1      Rttf2pt1_1.3.3    magrittr_1.5      munsell_0.4.2
42 [5] uuid_0.1-2       colorspace_1.2-6  R6_2.1.1          plyr_1.8.3
43 [9] tools_3.2.3      parallel_3.2.3    gtable_0.1.2      DBI_0.3.1
44 [13] extrafontdb_1.0  assertthat_0.1    IRdisplay_0.3     repr_0.4
45 [17] base64enc_0.1-3  IRkernel_0.5      evaluate_0.8      rzmq_0.7.7
46 [21] stringi_0.5-5    jsonlite_0.9.19

```

## Process the Data

Take the movies data, load in R friendly format, and combine with Rotten Tomatoes data.

```

1 df <- read_delim("~/Downloads/omdb0316/omdbMovies.txt", "\t",
  col_types="icccccccidc_c_____")
2 df_tomatoes <- read_delim("~/Downloads/omdb0316/tomatoes.txt", "\t",
  col_types="i_diiiiidc_c_c")
3 df <- df %>% left_join(df_tomatoes, by="ID")
4 rm(df_tomatoes)

1 |=====| 100% 435
  MB

1 parseBoxOffice <- function(x) {
2   unit <- 0
3   if (is.na(x) | x=="") {return (NA)}
4   if (substr(x, nchar(x), nchar(x)) == "k") {unit <- 10^3}
5   else {unit <- 10^6}
6
7   number <- as.numeric(substr(x,2,nchar(x)-1))
8
9   return(number * unit)
10 }
11
12 df <- df %>% mutate(BoxOffice = as.numeric(sapply(BoxOffice, parseBoxOffice)))

1 df_dup <- df %>% select(Title, Year) %>% mutate(Title = gsub("The ", "", Title))
2 dup <- duplicated(df_dup) # find entry indices which are duplicates
3 rm(df_dup) # remove temp dataframe
4
5 df <- df %>% filter(!dup) # keep entries which are *not* dups

```

## Inflation

```

1 inflation <- read_csv("http://research.stlouisfed.org/fred2/data/CPIAUCSL.csv") %>%
2   group_by(Year = as.integer(substr(DATE, 1, 4))) %>%
3   summarize(Avg_Value = mean(VALUE)) %>% # average across all months
4   mutate(Adjust = tail(Avg_Value, 1) / Avg_Value) # normalize by
  most-recent year

```

```
1 df <- df %>% inner_join(inflation) %>% mutate(AdjBoxOffice = floor(BoxOffice * Adjust))
```

```
1 Joining by: "Year"
```

Select only data we need now.

```
1 df <- df %>% filter(Year >= 2000, AdjBoxOffice >= 10^7, Cast != '') %>%
2   select(imdbID, Title, Year, Cast, Meter, Metacritic, AdjBoxOffice) %>%
3   arrange(desc(AdjBoxOffice))
4
5 #write.csv(df, "test.csv", row.names=F)
6 print(nrow(df))
```

```
1 [1] 2048
```

## Determine Gender of Lead

```
1 # Helper function to get first actor given a string of actors
2 getLeadActor <- function(actors) {
3   return(unlist(strsplit(actors, ", "))[1])
4 }
5
6 # Unicode issues during testing, so use string w/ unicode as a test case
7 print(getLeadActor("Will Smith, Robert De Niro, Renée Zellweger, Jack Black"))
```

```
1 [1] "Will Smith"
```

```
1 df$LeadActor <- as.character(lapply(enc2utf8(df$Cast), getLeadActor))
2
3 print(head(df %>% select(Title, LeadActor)))
4 print(head(df %>% filter(imdbID=="tt0307453") %>% select(Title, LeadActor)))
```

```
1 Source: local data frame [6 x 2]
2
3           Title           LeadActor
4      (chr)         (chr)
5 1 Star Wars: Episode VII - The Force Awakens Harrison Ford
6 2                               Avatar Sam Worthington
7 3                               Jurassic World Chris Pratt
8 4                               The Avengers Robert Downey Jr.
9 5                               The Dark Knight Christian Bale
10 6                               Shrek 2 Mike Myers
11 Source: local data frame [1 x 2]
12
13      Title LeadActor
14    (chr)   (chr)
15 1 Shark Tale Will Smith
```

Attempt #1: Merge known gender data of actors. (via list from Matt Daniels)

```
1 actor_gender <- read_csv("actor_list.csv") %>% select(LeadActor = name, Gender = gender)
2
3 print(head(actor_gender))
```

```
1 Source: local data frame [6 x 2]
```

```
2
3       LeadActor Gender
4       (chr)   (chr)
5 1      Keir Dullea      m
6 2      Gary Lockwood     m
7 3 William Sylvester     m
8 4      Daniel Richter    m
9 5    Leonard Rossiter    m
10 6   Margaret Tyzack     f
```

```
1 df <- df %>% left_join(actor_gender)
2
3 print(head(df %>% select(Title, LeadActor, Gender)))
```

```
1 Joining by: "LeadActor"
```

```
2
3
4 Source: local data frame [6 x 3]
5
6           Title      LeadActor Gender
7           (chr)      (chr)   (chr)
8 1 Star Wars: Episode VII - The Force Awakens Harrison Ford      m
9 2           Avatar    Sam Worthington    NA
10 3      Jurassic World    Chris Pratt    NA
11 4      The Avengers Robert Downey Jr.    NA
12 5      The Dark Knight    Christian Bale    NA
13 6           Shrek 2      Mike Myers    NA
```

Attempt #2: Determine gender from most-likely guess from first name. (Using male and female lists from Carnegie-Mellon University)

```
1 male_names <- unlist(read_delim("male_names.txt", "\n", skip = 6, col_names=F))
2 female_names <- unlist(read_delim("female_names.txt", "\n", skip = 6, col_names=F))
3
4 print(head(male_names))
5 print(head(female_names))
```

```
1      X11      X12      X13      X14      X15      X16
2 "Aamir" "Aaron" "Abbey" "Abbie" "Abbot" "Abbott"
3      X11      X12      X13      X14      X15      X16
4 "Abagael" "Abigail" "Abbe" "Abbey" "Abbi" "Abbie"
```

```
1 getGenderFromFullName <- function (full_name) {
2   first_name <- unlist(strsplit(full_name, " "))[1]
3   gender <- ifelse(first_name %in% male_names, "m",
4                     ifelse(first_name %in% female_names, "f", "[EDIT ME]"))
5   return (gender)
6 }
7
8 print(getGenderFromFullName("Sam Worthington"))
9 print(getGenderFromFullName("Kristen Wiig"))
```

```
1 [1] "m"
2 [1] "f"
```

```

1 gender_guess <- as.character(lapply(as.character(df$LeadActor), getGenderFromFullName))
2
3 # if a known gender from IMDB is present, use that; else, use the gender guess
4 df$Gender <- ifelse(is.na(df$Gender), gender_guess, df$Gender)
5
6 print(head(df %>% select(Title, LeadActor, Gender)))
7 print(tail(df %>% select(Title, LeadActor, Gender)))

```

```

1 Source: local data frame [6 x 3]
2
3           Title           LeadActor Gender
4         (chr)         (chr)   (chr)
5 1 Star Wars: Episode VII - The Force Awakens Harrison Ford      m
6 2           Avatar      Sam Worthington      m
7 3       Jurassic World      Chris Pratt      m
8 4       The Avengers Robert Downey Jr.      m
9 5       The Dark Knight      Christian Bale      m
10 6          Shrek 2      Mike Myers      m
11 Source: local data frame [6 x 3]
12
13          Title          LeadActor   Gender
14        (chr)        (chr)   (chr)
15 1          The Man      Samuel L. Jackson      m
16 2 Aliens of the Deep Anatoly M. Sagalevitch      m
17 3       A Single Man      Colin Firth      m
18 4          Pollock      Ed Harris      m
19 5       Connie and Carla      Nia Vardalos [EDIT ME]
20 6 I Don't Know How She Does It Sarah Jessica Parker      f

```

Attempt #3: Manually edit edge cases in a GUI (not shown)

```

1 write.csv(df, "movie_gender_intermediate.csv", row.names=F)

```

## Begin the Analysis

Reload the updated dataset (a few rows were removed due to being dupes)

```

1 df <- read_csv("movie_gender_fixed.csv")
2
3 print(head(df %>% select(Title, LeadActor, Gender)))
4 print(nrow(df))

```

```

1 Source: local data frame [6 x 3]
2
3           Title           LeadActor Gender
4         (chr)         (chr)   (chr)
5 1 Star Wars: Episode VII - The Force Awakens Daisy Ridley      f
6 2           Avatar      Sam Worthington      m
7 3       Jurassic World      Chris Pratt      m
8 4       The Avengers Robert Downey Jr.      m
9 5       The Dark Knight      Christian Bale      m
10 6          Shrek 2      Mike Myers      m
11 [1] 2020

```

Can we remove any points (e.g. Star Wars) as outliers? (tests via R Explorations)

```

1 AdjRevenue <- unlist(head(df %>% select(AdjBoxOffice)))
2
3 dixon.test(AdjRevenue, opposite=F)
4 grubbs.test(AdjRevenue, opposite=F)
5 chisq.out.test(AdjRevenue, variance=var(AdjRevenue), opposite=F)

1     Dixon test for outliers
2
3 data: AdjRevenue
4 Q.AdjBoxOffice1 = 0.23695, p-value = 0.8916
5 alternative hypothesis: highest value 934381231 is an outlier
6
7
8
9
10
11
12
13
14     Grubbs test for one outlier
15
16 data: AdjRevenue
17 G.AdjBoxOffice1 = 1.52510, U = 0.44175, p-value = 0.2634
18 alternative hypothesis: highest value 934381231 is an outlier
19
20
21
22
23
24
25
26
27     chi-squared test for outlier
28
29 data: AdjRevenue
30 X-squared.AdjBoxOffice1 = 2.326, p-value = 0.1272
31 alternative hypothesis: highest value 934381231 is an outlier

```

No outlier detection test supports it.

## Plot Box Office Revenues

```

1 df_summary <- df %>%
2     group_by(Gender) %>%
3     summarize(count = n(),
4               perc = n()/nrow(df),
5               mean = mean(AdjBoxOffice),
6               median = median(AdjBoxOffice))
7
8 color_m <- "#2980b9"
9 color_f <- "#27ae60"
10
11 print(df_summary)

```

```

1 Source: local data frame [2 x 5]
2
3   Gender count      perc      mean  median
4   (chr) (int)      (dbl)      (dbl)   (int)
5 1      f   467 0.2311881 65586882 44144648
6 2      m  1553 0.7688119 79786060 49841069

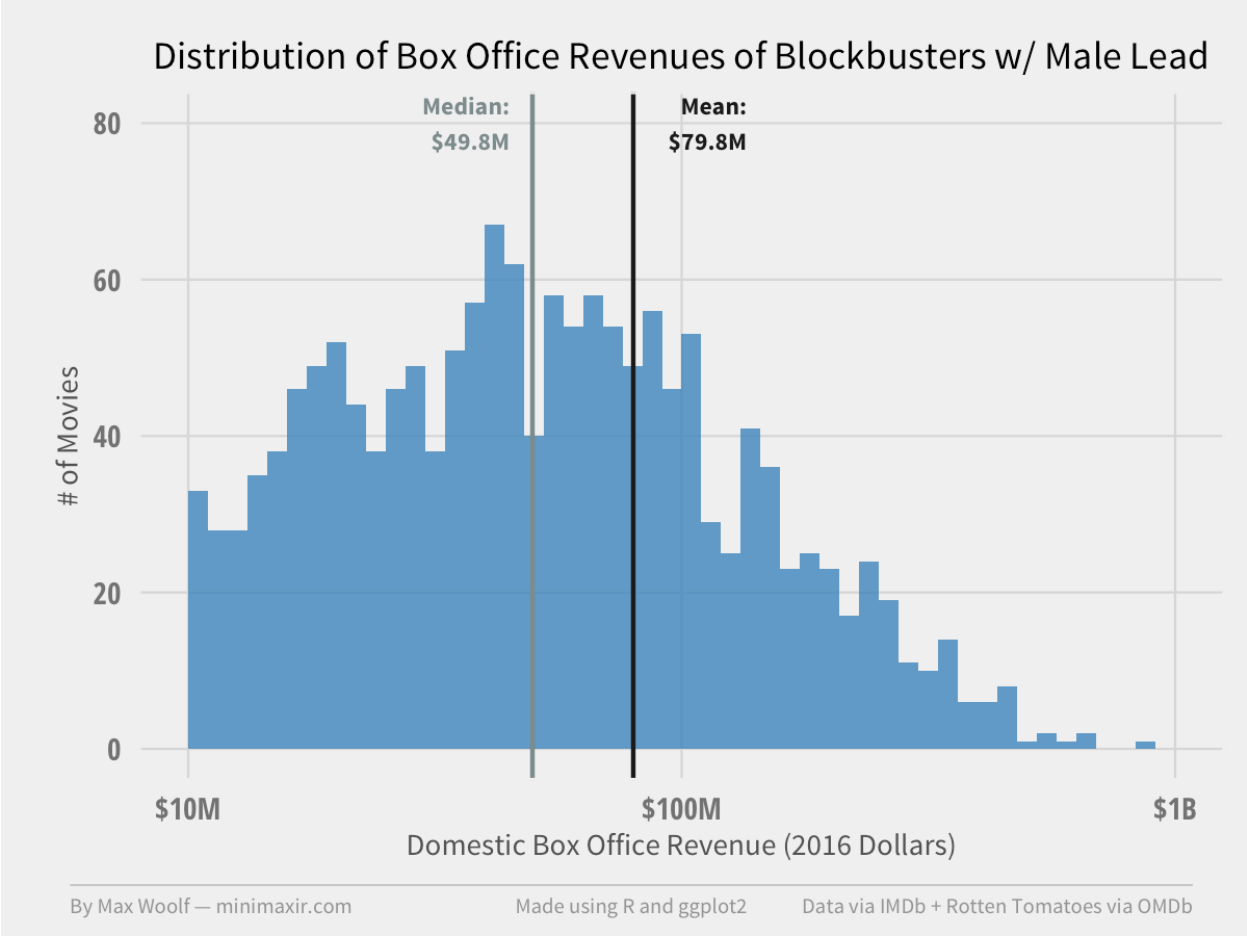
```

Plot Male and Female distributions separately, since medians are too close.

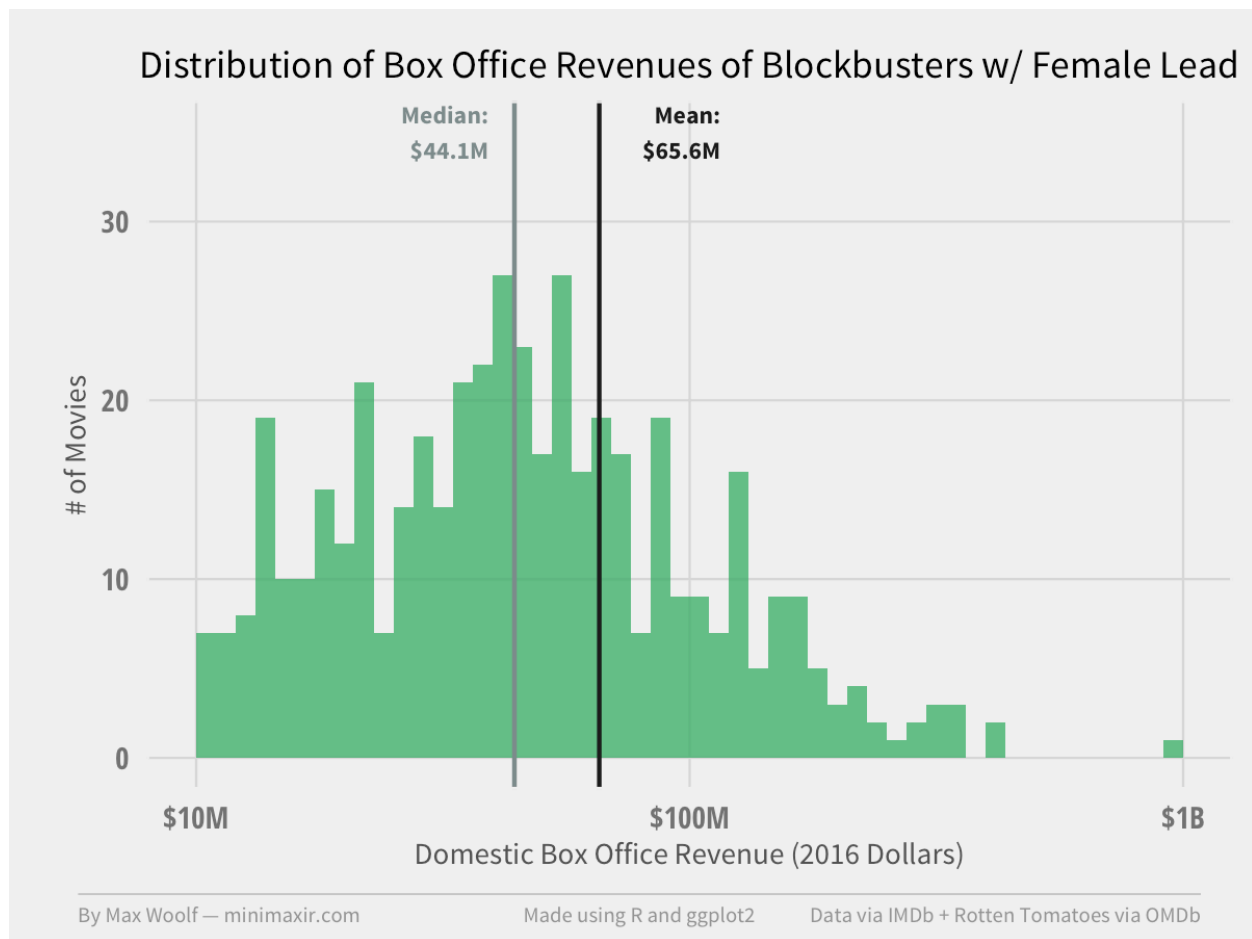
```

1 df_summary_m <- df_summary %>% filter(Gender=="m")
2
3 plot <- ggplot(df %>% filter(Gender=="m"), aes(x=AdjBoxOffice)) +
4   geom_histogram(fill=color_m, bins=50, alpha=0.75) +
5   fte_theme() +
6   scale_x_log10(limits=c(10^7, 10^9), breaks=10^c(7:9), labels=c("$10M", "$100M",
7     "$1B")) +
8   geom_vline(xintercept=df_summary_m$mean, color="#1a1a1a") +
9   geom_vline(xintercept=df_summary_m$median, color="#7f8c8d") +
10  annotate(geom="text", label = "Mean:\n$79.8M", x=df_summary_f$mean+7*10^7, y=80,
11    color="#1a1a1a", family="Source Sans Pro Bold", hjust=1, size=2) +
12  annotate(geom="text", label = "Median:\n$49.8M", x=df_summary_m$median-0.5*10^7,
13    y=80, color="#7f8c8d", family="Source Sans Pro Bold", hjust=1, size=2) +
14  labs(title="Distribution of Box Office Revenues of Blockbusters w/ Male Lead",
15    x="Domestic Box Office Revenue (2016 Dollars)", y="# of Movies")
16
17 max_save(plot, "movie-gender-1", "IMDb + Rotten Tomatoes via OMDb")
18
19 df_summary_f <- df_summary %>% filter(Gender=="f")
20
21 plot <- ggplot(df %>% filter(Gender=="f"), aes(x=AdjBoxOffice)) +
22   geom_histogram(fill=color_f, bins=50, alpha=0.75) +
23   fte_theme() +
24   scale_x_log10(limits=c(10^7, 10^9), breaks=10^c(7:9), labels=c("$10M", "$100M",
25     "$1B")) +
26   geom_vline(xintercept=df_summary_f$mean, color="#1a1a1a") +
27   geom_vline(xintercept=df_summary_f$median, color="#7f8c8d") +
28   annotate(geom="text", label = "Mean:\n$65.6M", x=df_summary_f$mean+5*10^7, y=35,
29     color="#1a1a1a", family="Source Sans Pro Bold", hjust=1, size=2) +
30   annotate(geom="text", label = "Median:\n$44.1M", x=df_summary_f$median-0.5*10^7,
31     y=35, color="#7f8c8d", family="Source Sans Pro Bold", hjust=1, size=2) +
32   labs(title="Distribution of Box Office Revenues of Blockbusters w/ Female Lead",
33     x="Domestic Box Office Revenue (2016 Dollars)", y="# of Movies")
34
35 max_save(plot, "movie-gender-2", "IMDb + Rotten Tomatoes via OMDb")

```







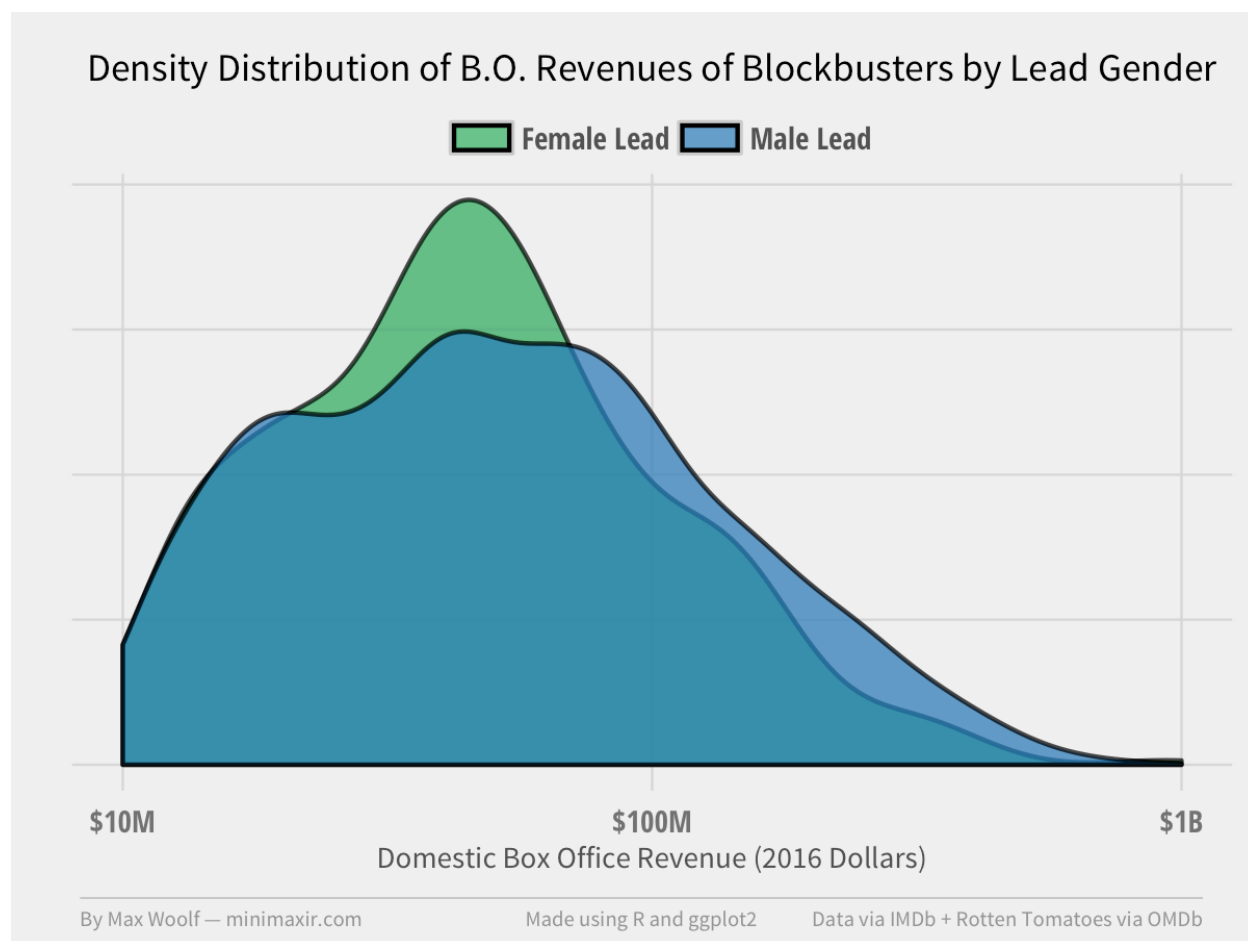


Figure 1:

```

6
7 ## check if log-scaling changes the result
8
9 ks_test <- ks.test(
10     log10(unlist(df %>% filter(Gender=="m") %>% select(AdjBoxOffice))),
11     log10(unlist(df %>% filter(Gender=="f") %>% select(AdjBoxOffice))))
12
13 print(ks_test)

1     Two-sample Kolmogorov-Smirnov test
2
3 data:  unlist(df %>% filter(Gender == "m") %>% select(AdjBoxOffice)) and unlist(df %>%
4     filter(Gender == "f") %>% select(AdjBoxOffice))
5 D = 0.10585, p-value = 0.0006411
6 alternative hypothesis: two-sided
7
8     Two-sample Kolmogorov-Smirnov test
9
10 data:  log10(unlist(df %>% filter(Gender == "m") %>% select(AdjBoxOffice))) and
11     log10(unlist(df %>% filter(Gender == "f") %>% select(AdjBoxOffice)))
12 D = 0.10585, p-value = 0.0006411
13 alternative hypothesis: two-sided

The distribution is different! Are the differences in means statistically significant?

1 wilcox_test <- wilcox.test(
2     unlist(df %>% filter(Gender=="m") %>% select(AdjBoxOffice)),
3     unlist(df %>% filter(Gender=="f") %>% select(AdjBoxOffice)),
4     alternative = "g")
5
6 print(wilcox_test)
7
8 ## check if log-scaling changes the result
9
10 wilcox_test <- wilcox.test(
11     log10(unlist(df %>% filter(Gender=="m") %>% select(AdjBoxOffice))),
12     log10(unlist(df %>% filter(Gender=="f") %>% select(AdjBoxOffice))),
13     alternative = "g")
14
15 print(wilcox_test)

1     Wilcoxon rank sum test with continuity correction
2
3 data:  unlist(df %>% filter(Gender == "m") %>% select(AdjBoxOffice)) and unlist(df %>%
4     filter(Gender == "f") %>% select(AdjBoxOffice))
5 W = 390070, p-value = 0.006514
6 alternative hypothesis: true location shift is greater than 0
7
8     Wilcoxon rank sum test with continuity correction
9
10 data:  log10(unlist(df %>% filter(Gender == "m") %>% select(AdjBoxOffice))) and
11     log10(unlist(df %>% filter(Gender == "f") %>% select(AdjBoxOffice)))
12 W = 390070, p-value = 0.006514

```

12 alternative hypothesis: true location shift is greater than 0

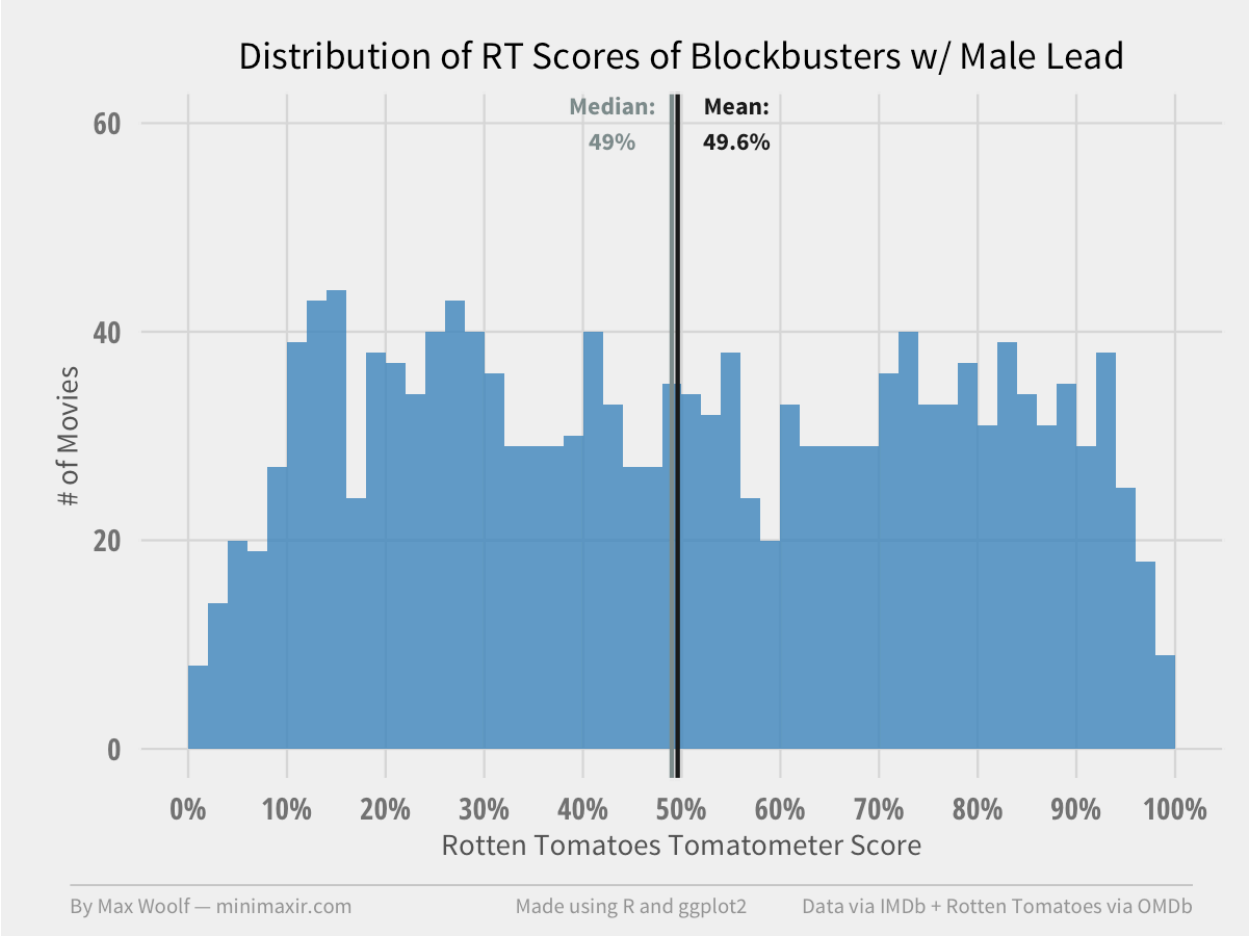
## Plot Rotten Tomatoes Meter

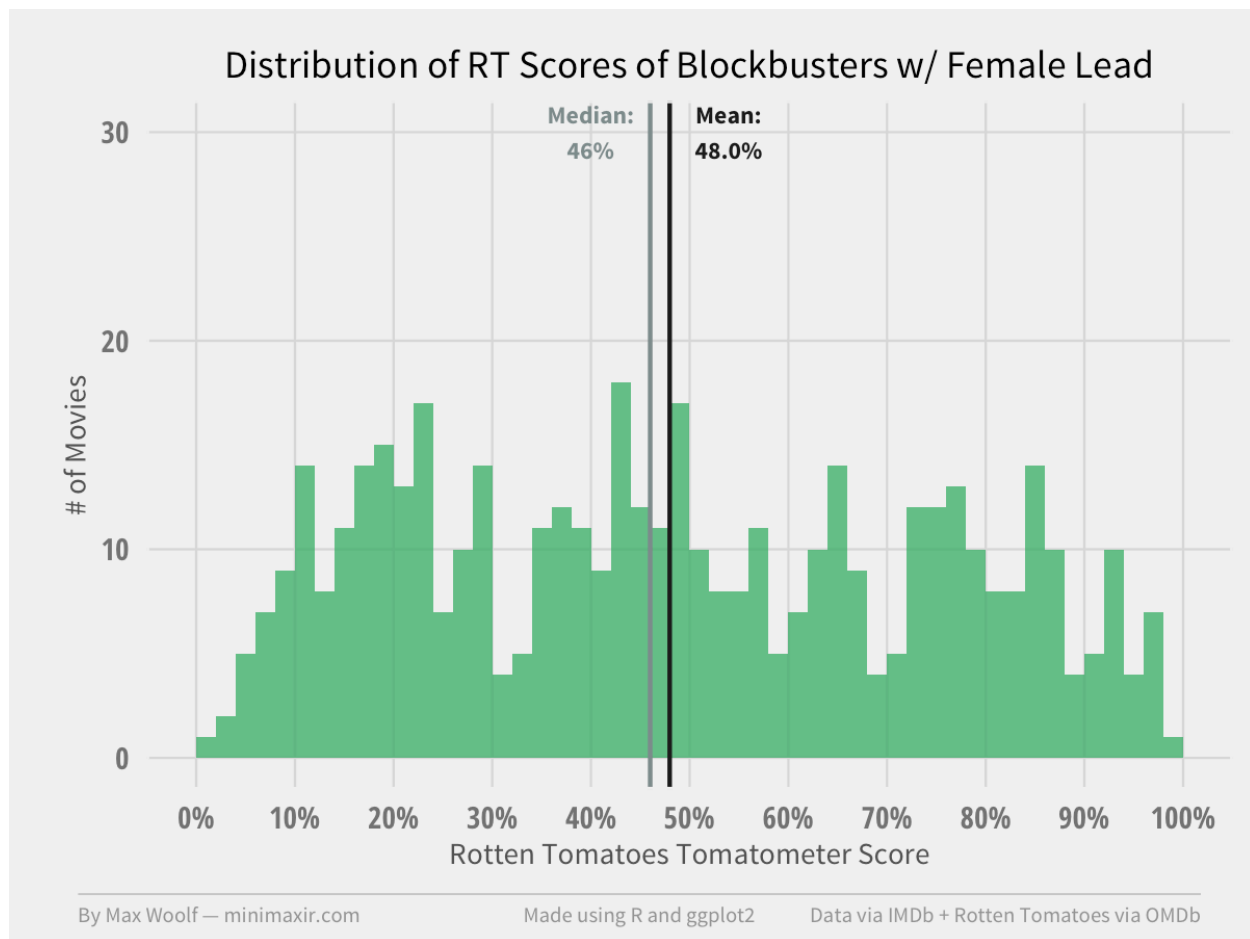
Can reuse most of the code, unfortunately have to violate DRY for ad-hoc fixes.

```
1 df_summary <- df %>%
2   group_by(Gender) %>%
3   summarize(mean = mean(Meter, na.rm=T), median = median(Meter, na.rm=T))
4
5
6 print(df_summary)

1 Source: local data frame [2 x 3]
2
3   Gender      mean median
4   (chr)    (dbl)  (int)
5 1      f 47.97859     46
6 2      m 49.59381     49

1 df_summary_m <- df_summary %>% filter(Gender=="m")
2
3 plot <- ggplot(df %>% filter(Gender=="m"), aes(x=Meter)) +
4   geom_histogram(fill=color_m, bins=50, alpha=0.75) +
5   fte_theme() +
6   scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100),
7     labels=paste0(seq(0,100, by=10),"%")) +
8   geom_vline(xintercept=df_summary_m$mean, color="#1a1a1a") +
9   geom_vline(xintercept=df_summary_m$median, color="#7f8c8d") +
10  annotate(geom="text", label = "Mean:\n49.6%", x=df_summary_m$mean+6, y=60,
11    color="#1a1a1a", family="Source Sans Pro Bold", size=2) +
12  annotate(geom="text", label = "Median:\n49%", x=df_summary_m$median-6, y=60,
13    color="#7f8c8d", family="Source Sans Pro Bold", size=2) +
14  labs(title="Distribution of RT Scores of Blockbusters w/ Male Lead", x="Rotten
15    Tomatoes Tomatometer Score", y="# of Movies")
16
17 max_save(plot, "movie-gender-4", "IMDb + Rotten Tomatoes via OMDb")
18
19 df_summary_f <- df_summary %>% filter(Gender=="f")
20
21 plot <- ggplot(df %>% filter(Gender=="f"), aes(x=Meter)) +
22   geom_histogram(fill=color_f, bins=50, alpha=0.75) +
23   fte_theme() +
24   scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100),
25     labels=paste0(seq(0,100, by=10),"%")) +
26   geom_vline(xintercept=df_summary_f$mean, color="#1a1a1a") +
27   geom_vline(xintercept=df_summary_f$median, color="#7f8c8d") +
28   annotate(geom="text", label = "Mean:\n48.0%", x=df_summary_f$mean+6, y=30,
29     color="#1a1a1a", family="Source Sans Pro Bold", size=2) +
30   annotate(geom="text", label = "Median:\n46%", x=df_summary_f$median-6, y=30,
31     color="#7f8c8d", family="Source Sans Pro Bold", size=2) +
32   labs(title="Distribution of RT Scores of Blockbusters w/ Female Lead", x="Rotten
33     Tomatoes Tomatometer Score", y="# of Movies")
34
35 max_save(plot, "movie-gender-5", "IMDb + Rotten Tomatoes via OMDb")
```





```

1 plot <- ggplot(df, aes(x=Meter, fill=Gender)) +
2   geom_density(alpha=0.75) +
3   fte_theme() +
4   scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100),
5     labels=paste0(seq(0,100, by=10),"%")) +
6   theme(legend.title = element_blank(), legend.position="top",
7     legend.direction="horizontal", legend.key.width=unit(0.5, "cm"),
8     legend.key.height=unit(0.25, "cm"), legend.margin=unit(0,"cm"),
9     axis.title.y=element_blank(), axis.text.y=element_blank()) +
10  scale_fill_manual(labels=c("Female Lead", "Male Lead"),
11    values=c(color_f,color_m)) +
12  labs(title="Density Distribution of RT Scores of Blockbusters by Lead Gender",
13    x="Rotten Tomatoes Tomatometer Score")
14
15 max_save(plot, "movie-gender-6", "IMDb + Rotten Tomatoes via OMDb")
16
17 ks_test <- ks.test(
18   unlist(df %>% filter(Gender=="m") %>% select(Meter)),
19   unlist(df %>% filter(Gender=="f") %>% select(Meter)))
20
21 print(ks_test)
22
23 wilcox_test <- wilcox.test(
24   unlist(df %>% filter(Gender=="m") %>% select(Meter)),

```

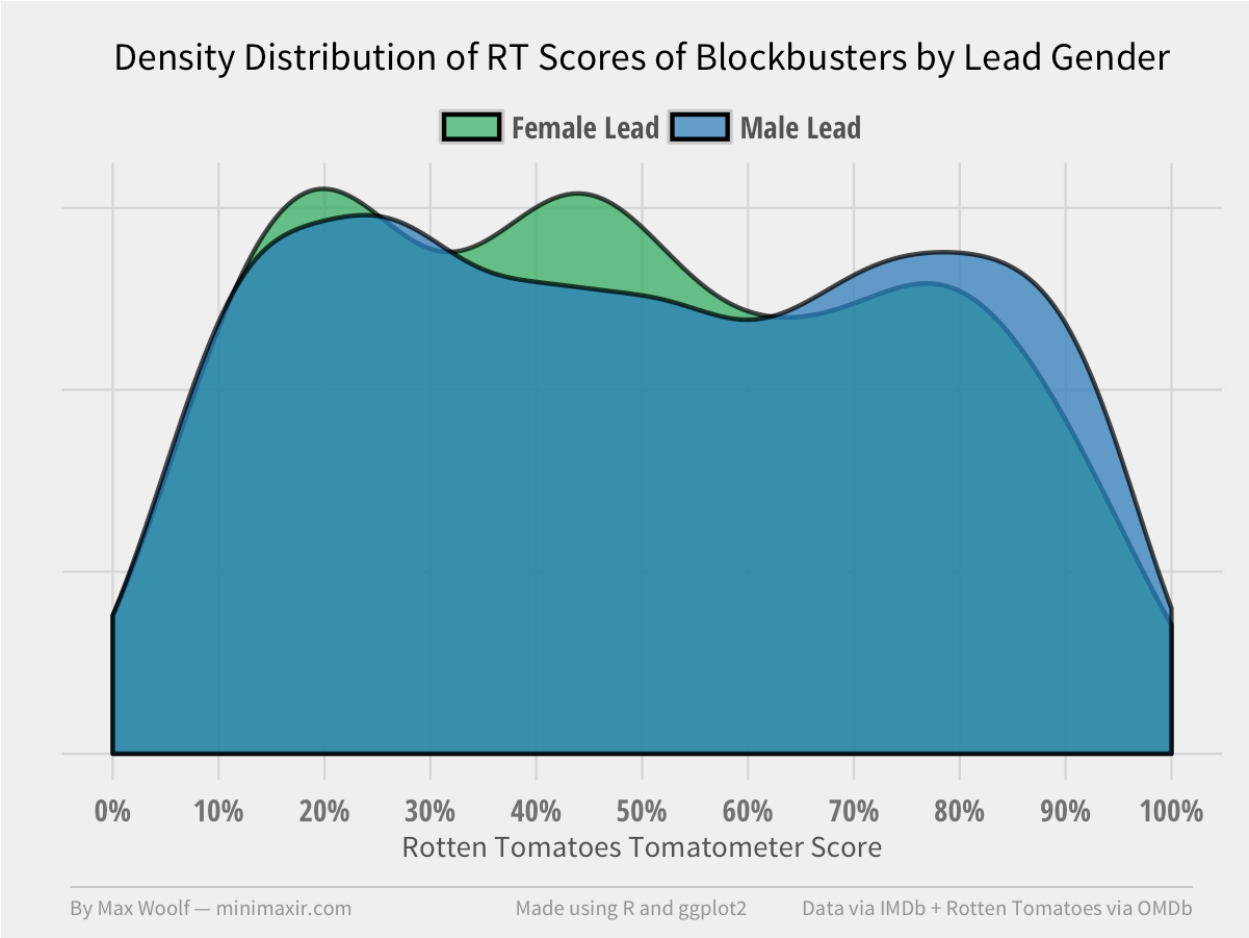


Figure 2:

```

9         unlist(df %>% filter(Gender=="f") %>% select(Meter)),
10         alternative="g")
11
12 print(wilcox_test)

1     Two-sample Kolmogorov-Smirnov test
2
3 data:  unlist(df %>% filter(Gender == "m") %>% select(Meter)) and unlist(df %>%
         filter(Gender == "f") %>% select(Meter))
4 D = 0.048455, p-value = 0.3684
5 alternative hypothesis: two-sided
6
7
8     Wilcoxon rank sum test with continuity correction
9
10 data:  unlist(df %>% filter(Gender == "m") %>% select(Meter)) and unlist(df %>%
         filter(Gender == "f") %>% select(Meter))
11 W = 374460, p-value = 0.1326
12 alternative hypothesis: true location shift is greater than 0

```

## Plot Metacritic

```

1 df_summary <- df %>%
2     group_by(Gender) %>%
3     summarize(mean = mean(Metacritic, na.rm=T), median = median(Metacritic,
4         na.rm=T))
5
6 print(df_summary)

1 Source: local data frame [2 x 3]
2
3   Gender    mean median
4   (chr)    (dbl)  (dbl)
5 1     f 50.78523    50
6 2     m 51.76032    51

1 df_summary_m <- df_summary %>% filter(Gender=="m")
2
3 plot <- ggplot(df %>% filter(Gender=="m"), aes(x=Metacritic)) +
4     geom_histogram(fill=color_m, bins=50, alpha=0.75) +
5     fte_theme() +
6     scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100)) +
7     geom_vline(xintercept=df_summary_m$mean, color="#1a1a1a") +
8     geom_vline(xintercept=df_summary_m$median, color="#7f8c8d") +
9     annotate(geom="text", label = "Mean:\n51.8", x=df_summary_m$mean+6, y=80,
10         color="#1a1a1a", family="Source Sans Pro Bold", size=2) +
11     annotate(geom="text", label = "Median:\n51", x=df_summary_m$median-6, y=80,
12         color="#7f8c8d", family="Source Sans Pro Bold", size=2) +
13     labs(title="Distribution of Metacritic Scores of Blockbusters w/ Male Lead",
14         x="Metacritic Score", y="# of Movies")
12
13 max_save(plot, "movie-gender-7", "IMDb + Rotten Tomatoes via OMDb")

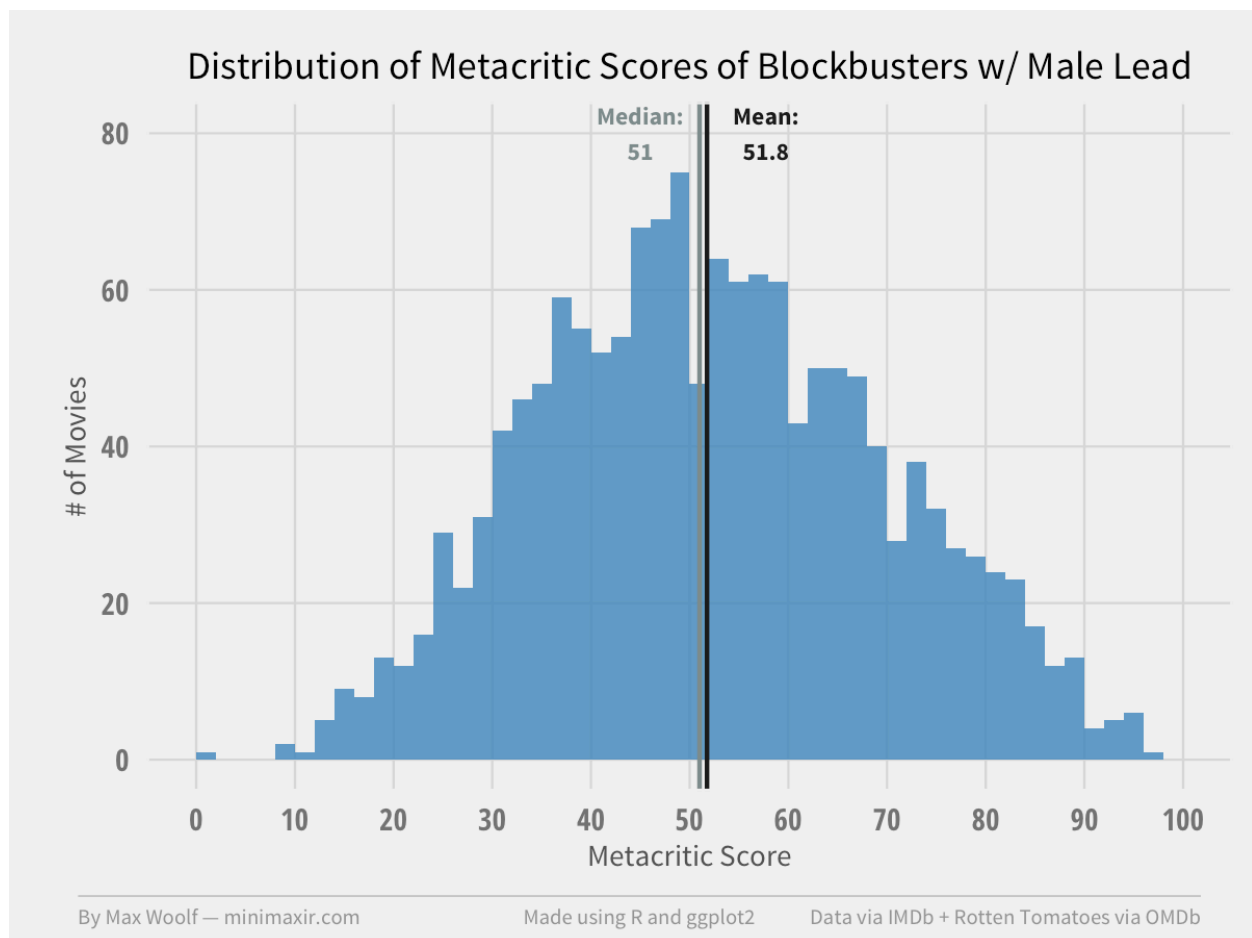
```

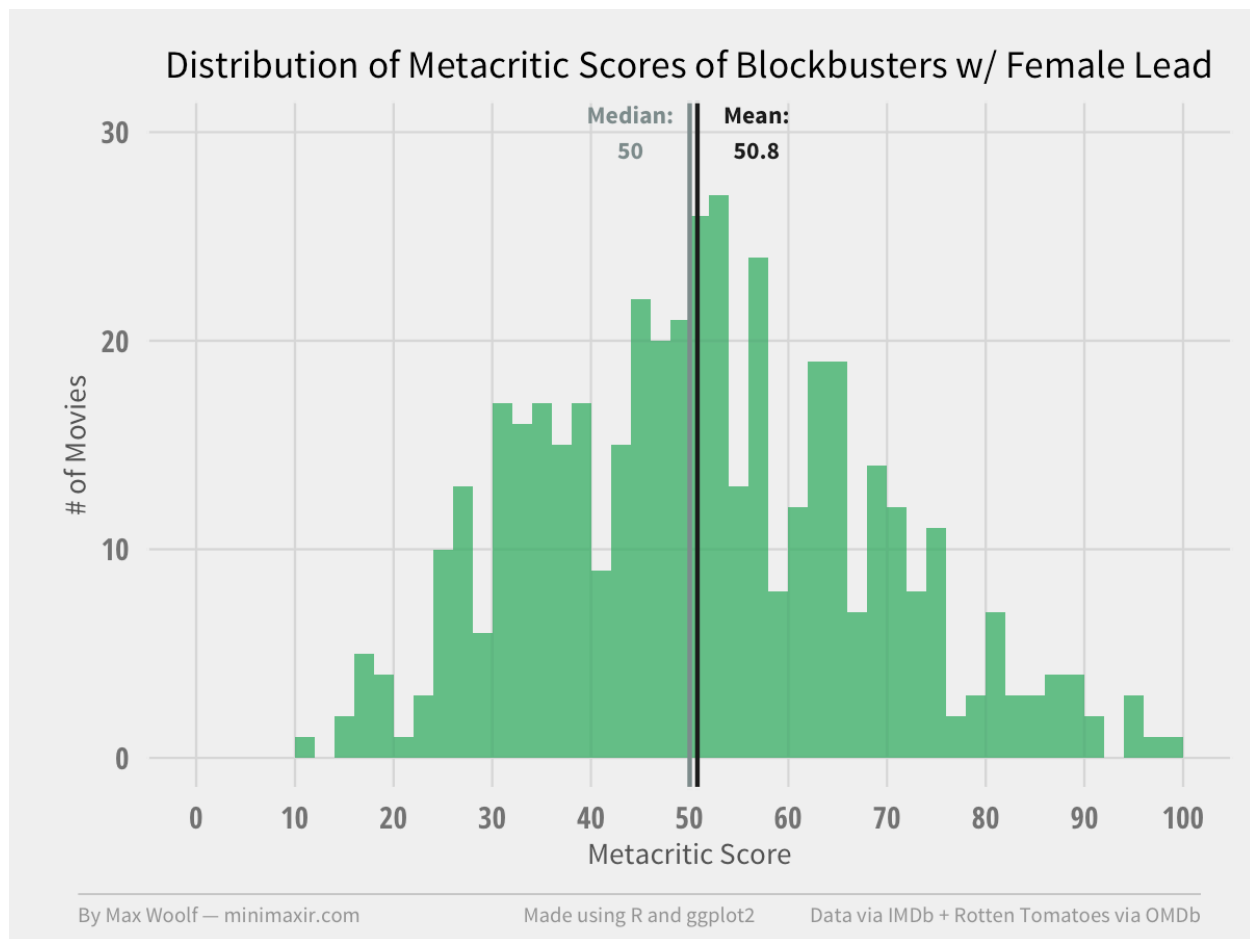


```

14
15 df_summary_f <- df_summary %>% filter(Gender=="f")
16
17 plot <- ggplot(df %>% filter(Gender=="f"), aes(x=Metacritic)) +
18   geom_histogram(fill=color_f, bins=50, alpha=0.75) +
19   fte_theme() +
20   scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100)) +
21   geom_vline(xintercept=df_summary_f$mean, color="#1a1a1a") +
22   geom_vline(xintercept=df_summary_f$median, color="#7f8c8d") +
23   annotate(geom="text", label = "Mean:\n50.8", x=df_summary_f$mean+6, y=30,
24     color="#1a1a1a", family="Source Sans Pro Bold", size=2) +
25   annotate(geom="text", label = "Median:\n50", x=df_summary_f$median-6, y=30,
26     color="#7f8c8d", family="Source Sans Pro Bold", size=2) +
27   labs(title="Distribution of Metacritic Scores of Blockbusters w/ Female Lead",
28     x="Metacritic Score", y="# of Movies")

```





```

1 plot <- ggplot(df, aes(x=Metacritic, fill=Gender)) +
2   geom_density(alpha=0.75) +
3   fte_theme() +
4   scale_x_continuous(breaks=seq(0,100, by=10), limits=c(0, 100)) +
5   theme(legend.title = element_blank(), legend.position="top",
6         legend.direction="horizontal", legend.key.width=unit(0.5, "cm"),
7         legend.key.height=unit(0.25, "cm"), legend.margin=unit(0, "cm"),
8         axis.title.y=element_blank(), axis.text.y=element_blank()) +
9   scale_fill_manual(labels=c("Female Lead", "Male Lead"),
10                    values=c(color_f,color_m)) +
11   labs(title="Density Distribution of Metacritic Scores of Blockbusters by Lead
12         Gender", x="Metacritic Score")
13
14 max_save(plot, "movie-gender-9", "IMDb + Rotten Tomatoes via OMDb")
15
16 ks_test <- ks.test(
17   unlist(df %>% filter(Gender=="m") %>% select(Metacritic)),
18   unlist(df %>% filter(Gender=="f") %>% select(Metacritic)))
19
20 print(ks_test)
21
22 wilcox_test <- wilcox.test(
23   unlist(df %>% filter(Gender=="m") %>% select(Metacritic)),
24   unlist(df %>% filter(Gender=="f") %>% select(Metacritic)),

```

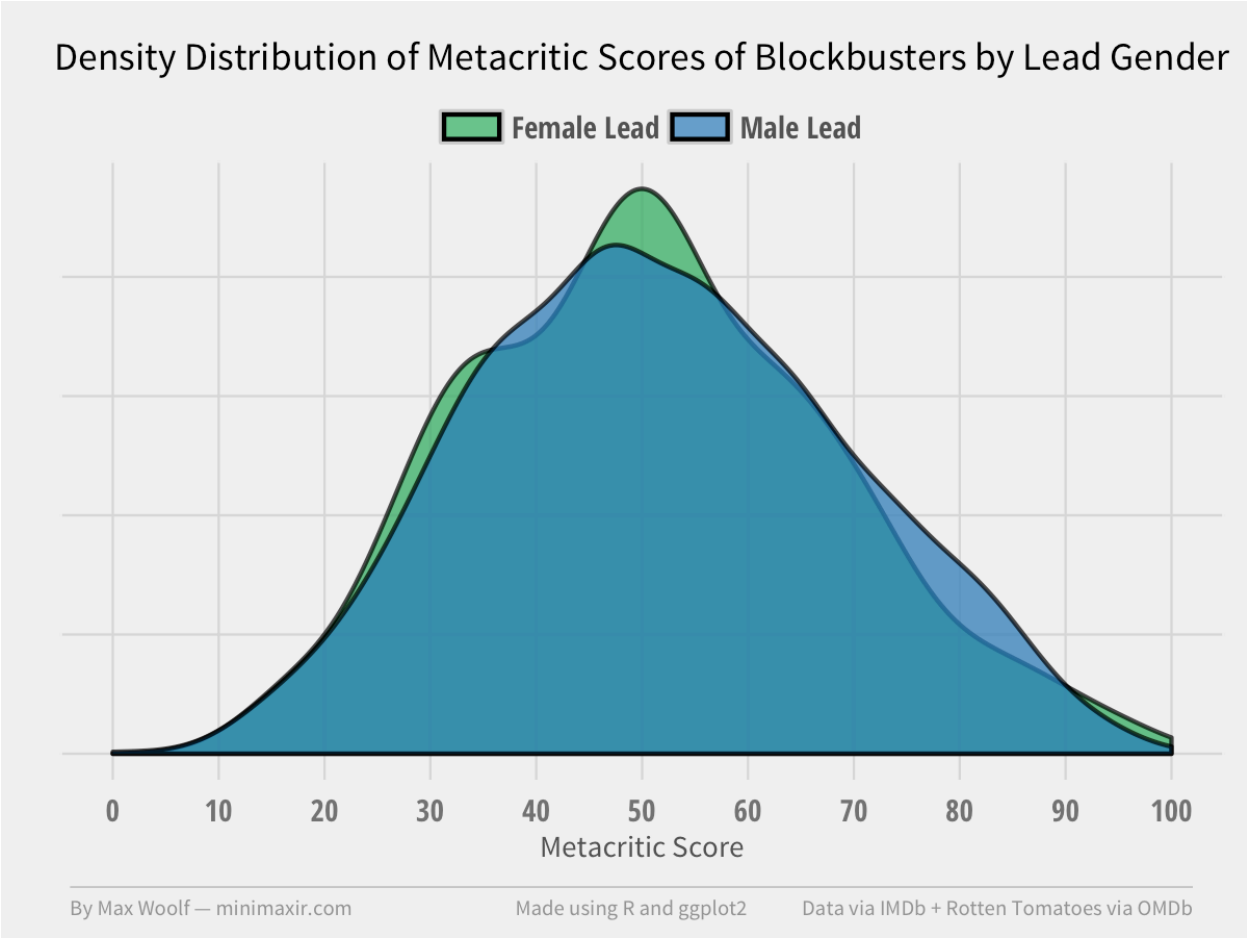


Figure 3:

```

10         alternative="g")
11
12 print(wilcox_test)

1   Two-sample Kolmogorov-Smirnov test
2
3 data:  unlist(df %>% filter(Gender == "m") %>% select(Metacritic)) and unlist(df %>%
         filter(Gender == "f") %>% select(Metacritic))
4 D = 0.046268, p-value = 0.4521
5 alternative hypothesis: two-sided
6
7
8   Wilcoxon rank sum test with continuity correction
9
10 data:  unlist(df %>% filter(Gender == "m") %>% select(Metacritic)) and unlist(df %>%
         filter(Gender == "f") %>% select(Metacritic))
11 W = 347130, p-value = 0.1368
12 alternative hypothesis: true location shift is greater than 0

```

## Bootstrap Resample Means

```

1 resampleMeans <- function(df) {
2   df_new <- df %>% sample_frac(replace=T)
3
4   summary <- df_new %>%
5     group_by(Gender) %>%
6     summarize(AdjBoxOffice_m = mean(AdjBoxOffice),
7               Meter_m = mean(Meter, na.rm=T),
8               Metacritic_m = mean(Metacritic, na.rm=T))
9
10  return (summary)
11 }
12
13 set.seed(4)
14 print(resampleMeans(df))

```

```

1 Source: local data frame [2 x 4]
2
3   Gender AdjBoxOffice_m Meter_m Metacritic_m
4   (chr)      (dbl)      (dbl)      (dbl)
5 1     f        67986152 45.93206      49.40440
6 2     m        82138781 49.89780      52.02617

```

Pre-allocate space per this Stack Overflow answer.

```

1 resampleMovieData <- function(n) {
2
3   df_resample_summary <- data.frame(Gender = character(n*2), AdjBoxOffice_m = numeric(n*2),
4                                     Meter_m = numeric(n*2), Metacritic_m = numeric(n*2), stringsAsFactors =
5                                     FALSE)
6
7   for (i in seq(1,n*2 - 1, by = 2)) {
8     df_resample_summary[c(i,i+1),] <- resampleMeans(df)
9   }

```

```

10 return(tbl_df(df_resample_summary))
11
12 }
13
14 set.seed(4)
15 print(resampleMovieData(4))

```

```

1 Source: local data frame [8 x 4]
2
3   Gender AdjBoxOffice_m Meter_m Metacritic_m
4   (chr)      (dbl)      (dbl)      (dbl)
5 1      f      67986152 45.93206      49.40440
6 2      m      82138781 49.89780      52.02617
7 3      f      65405109 47.98039      51.02036
8 4      m      80331357 49.79551      51.76354
9 5      f      65073479 47.72557      50.11063
10 6      m      78003310 48.89669      51.36024
11 7      f      65424463 48.83369      51.77605
12 8      m      80139428 49.59100      51.92642

```

```

1 system.time( df_boot <- resampleMovieData(10000))
2
3 print(head(df_boot))
4 print(nrow(df_boot)) # expect 10000 * 2

```

```

1   user  system elapsed
2 37.460   4.230  42.968
3
4
5
6 Source: local data frame [6 x 4]
7
8   Gender AdjBoxOffice_m Meter_m Metacritic_m
9   (chr)      (dbl)      (dbl)      (dbl)
10 1      f      67708627 48.58747      51.07289
11 2      m      79596707 50.01286      52.30313
12 3      f      70793578 47.57675      50.28372
13 4      m      77681742 50.53171      52.40521
14 5      f      72614163 47.77322      50.12472
15 6      m      78020424 49.21093      51.50498
16 [1] 20000

```

```

1 df_boot_agg <- df_boot %>%
2   group_by(Gender) %>%
3   summarize(
4     AdjBoxOffice_res_m = mean(AdjBoxOffice_m),
5     AdjBoxOffice_low_ci = quantile(AdjBoxOffice_m, 0.025),
6     AdjBoxOffice_high_ci = quantile(AdjBoxOffice_m, 0.975)
7   )
8
9 print(df_boot_agg)

```

```

1 Source: local data frame [2 x 4]
2

```

	Gender	AdjBoxOffice_res_m	AdjBoxOffice_low_ci	AdjBoxOffice_high_ci
	(chr)	(dbl)	(dbl)	(dbl)
5 1	f	65531567	59238519	72485603
6 2	m	79758350	75590754	84168300

## Plot Final Bootstrap

```

1 df_summary_means <- df %>% group_by(Gender) %>% summarize(mean = mean(AdjBoxOffice))
2
3 plot <- ggplot(df_boot, aes(x=AdjBoxOffice_m, fill=Gender)) +
4   scale_x_continuous(limits=c(5*10^7, 10^8), breaks=seq(5*10^7, 9*10^7, by=10^7),
5     labels=paste0("$", seq(50, 90, by=10), "M")) +
6   scale_y_continuous(breaks=pretty_breaks(4)) +
7   geom_histogram(bins=100, alpha=0.75, position="identity") +
8   geom_point(mapping=aes(x=mean, y=0), data=df_summary_means, show.legend=F,
9     color="black") +
10  geom_errorbarh(mapping=aes(x=AdjBoxOffice_res_m, xmin=AdjBoxOffice_low_ci,
11    xmax=AdjBoxOffice_high_ci, y=0), data=df_boot_agg, show.legend=F, color="black",
12    height=0) +
13  fte_theme() +
14  theme(legend.title = element_blank(), legend.position="top",
15    legend.direction="horizontal", legend.key.width=unit(0.5, "cm"),
16    legend.key.height=unit(0.25, "cm"), legend.margin=unit(0, "cm")) +
17  scale_fill_manual(labels=c("Female Lead", "Male Lead"), values=c(color_f, color_m)) +
18  labs(title=sprintf("Resampled Avg. B.O. Revenues by Movie Lead Gender (n = %2d)",
19    nrow(df_boot)/2), x="Average Domestic Box Office Revenue for Blockbusters (2016
20    Dollars)", y="# of Resampled Averages")
21
22 max_save(plot, "movie-gender-10", "IMDb + Rotten Tomatoes via OMDb")

```

## Determine P-Value of Final Bootstrap

Calculate the difference between the bootstrapped means; the P-value is the proportion of values where  $m - f < 0$ .

```

1 n <- 10000
2
3 means_vector <- unlist(df_boot$AdjBoxOffice)
4 means_diff <- c()
5
6 for (i in seq(1, n*2 - 1, by = 2)) {
7   means_diff <- c(means_diff, means_vector[i+1] - means_vector[i])
8 }
9
10 print(means_diff[1:4])
11
12 print(sum(means_diff <= 0)/n) # p-value of difference

```

```

1 [1] 11888080 6888164 5406261 8136011
2 [1] 2e-04

```

## Bootstrap Movie!

Render each frame of the resample; composite into GIF later.

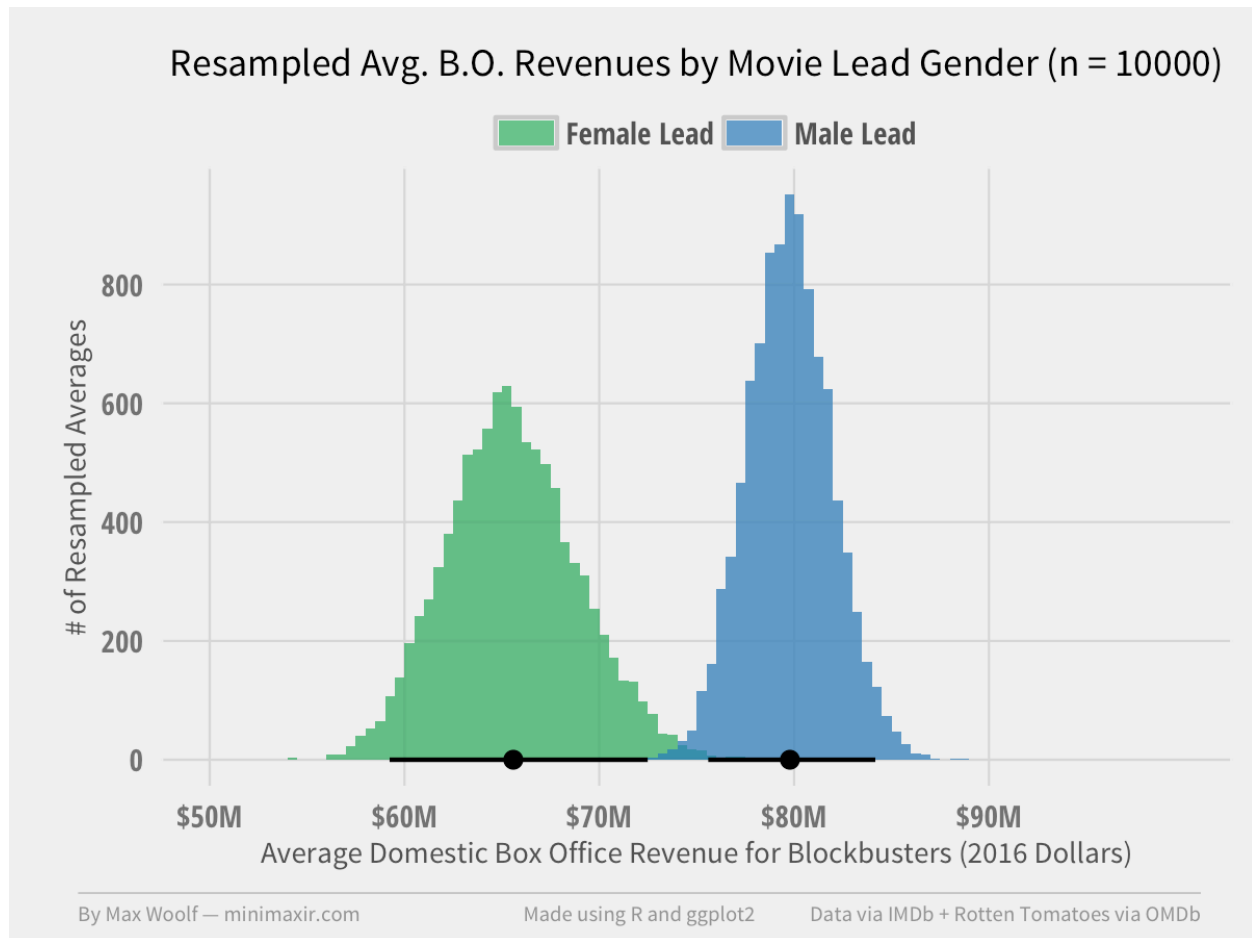


Figure 4:

```

1 system("mkdir -p movie_frames")
2
3 movie_frames <- function(size) {
4   df_boot_sub <- df_boot %>% head(size*2)
5
6   df_boot_agg_sub <- df_boot_sub %>%
7     group_by(Gender) %>%
8     summarize(
9       AdjBoxOffice_res_m = mean(AdjBoxOffice_m),
10      AdjBoxOffice_low_ci = quantile(AdjBoxOffice_m, 0.025),
11      AdjBoxOffice_high_ci = quantile(AdjBoxOffice_m, 0.975)
12    )
13
14   plot <- ggplot(df_boot_sub, aes(x=AdjBoxOffice_m, fill=Gender)) +
15     scale_x_continuous(limits=c(5*10^7, 10^8), breaks=seq(5*10^7, 9*10^7, by=10^7),
16       labels=paste0("$", seq(50,90, by=10), "M")) +
17     scale_y_continuous(breaks=pretty_breaks(4)) +
18     geom_histogram(bins=100, alpha=0.75, position="identity") +
19     geom_point(mapping=aes(x=mean, y=0), data=df_summary_means, show.legend=F,
20       color="black") +
21     geom_errorbarh(mapping=aes(x=AdjBoxOffice_res_m, xmin=AdjBoxOffice_low_ci,
22       xmax=AdjBoxOffice_high_ci, y=0), data=df_boot_agg_sub, show.legend=F, color="black",
23       height=0) +
24     fte_theme() +
25     theme(legend.title = element_blank(), legend.position="top",
26       legend.direction="horizontal", legend.key.width=unit(0.5, "cm"),
27       legend.key.height=unit(0.25, "cm"), legend.margin=unit(0, "cm")) +
28     scale_fill_manual(labels=c("Female Lead", "Male Lead"), values=c(color_f,color_m)) +
29     labs(title=sprintf("Resampled Avg. B.O. Revenues by Movie Lead Gender (n = %2d)", size),
30       x="Average Domestic Box Office Revenue for Blockbusters (2016 Dollars)", y="# of
31       Resampled Averages")
32
33   max_save(plot, sprintf("movie_frames/movie_%06d", size), "IMDb + Rotten Tomatoes via OMDb")
34 }
35
36 system.time( x <- lapply(seq(100,10000,100), movie_frames) )
37
38 user system elapsed
39 54.031 2.565 60.278

```

## The MIT License (MIT)

Copyright (c) 2016 Max Woolf

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FIT-



NESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.