

Statistical Analysis

Kathryn Sam

2025-12-11

Data setup

Load file

```
dfSpacy<-read_excel('MovieData_Spacy.xlsx')
```

Change variable types and recode variables.

```
# Recode character gender for interpretation
dfSpacy$CharIsWoman <- ifelse(dfSpacy$Gender == "M", 0, 1)

# Convert variable to factor
dfSpacy$CharIsWoman<-as.factor(dfSpacy$CharIsWoman)

# Recode writer gender for interpretation
# Man = 0, Both = 1, Woman = 2
dfSpacy$WriterGender <- ifelse(dfSpacy$WriterGender == "M", 0,
                              ifelse(dfSpacy$WriterGender == "B", 1, 2))
dfSpacy$WriterGender= as.factor(dfSpacy$WriterGender) # Convert to factor

dfSpacy$MovieYear = as.numeric(dfSpacy$MovieYear) # Convert variable to numeric
```

Check and remove movies which have no text in either character gender category

```
# Replace intensifiers with NA for rows with no word count
dfSpacy$intensifiers <- ifelse(dfSpacy$NW == 0, NA, dfSpacy$intensifiers)

# Inspect NAs
sum(is.na(dfSpacy$intensifiers))
```

```
## [1] 132
```

```
# Check if NAs are related to Char Sex code
table(is.na(dfSpacy$intensifiers), dfSpacy$CharIsWoman)
```

```
##
##           0    1
## FALSE 570 524
##  TRUE   43   89
```

```

# Count movies with no Text for both Male and Female
movies_no_text <- dfSpacy %>%
  group_by(MovieID) %>%
  summarise(
    no_text_male = all(is.na(Text[CharIsWoman == 0])),
    no_text_female = all(is.na(Text[CharIsWoman == 1]))
  ) %>%
  filter(no_text_male & no_text_female)

# Number of movies with no Text data
print(paste("Number of movies with no Text data:", nrow(movies_no_text)))

```

```
## [1] "Number of movies with no Text data: 19"
```

```

#Drop if movies have no text for either M or F
dfSpacy2 <- dfSpacy %>%
  filter(!MovieID %in% movies_no_text$MovieID)

# Final data
length(unique(dfSpacy2$MovieID))

```

```
## [1] 594
```

Descriptives

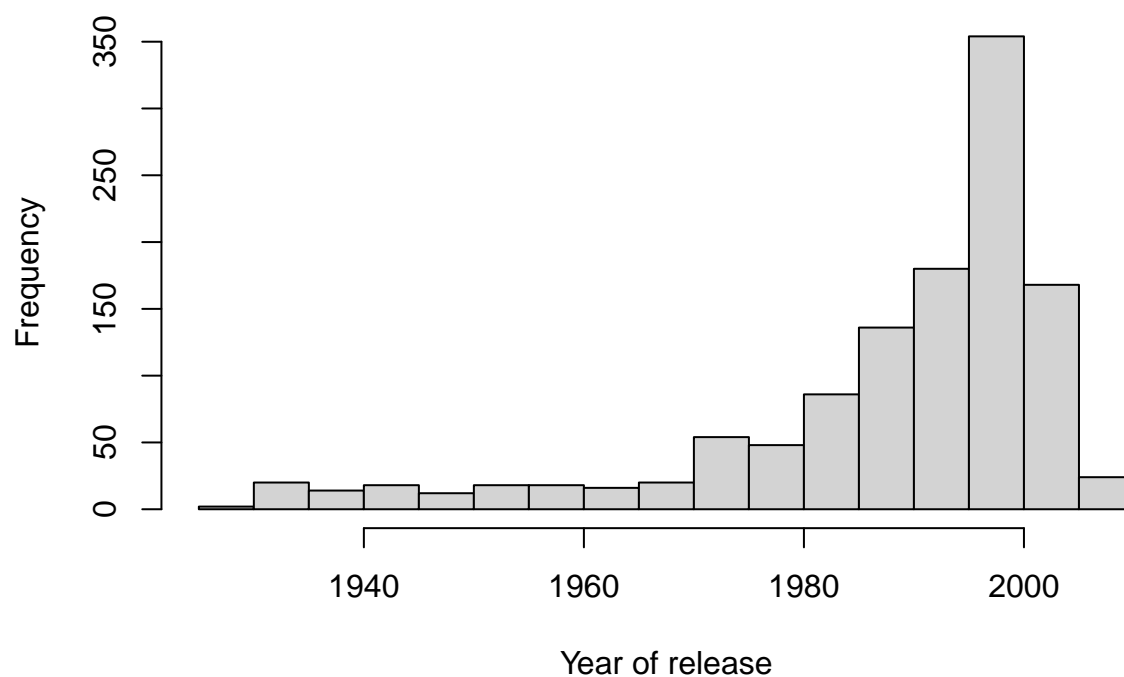
Movie Year

```
describe(dfSpacy2$MovieYear) # Descriptives
```

```
##      vars      n    mean      sd median trimmed mad  min  max range  skew kurtosis
## X1      1 1188 1988.53 16.36   1994 1991.63  8.9 1927 2009    82 -1.75     2.73
##      se
## X1 0.47
```

```
hist(dfSpacy2$MovieYear, xlab = "Year of release") # Histogram
```

Histogram of dfSpacy2\$MovieYear

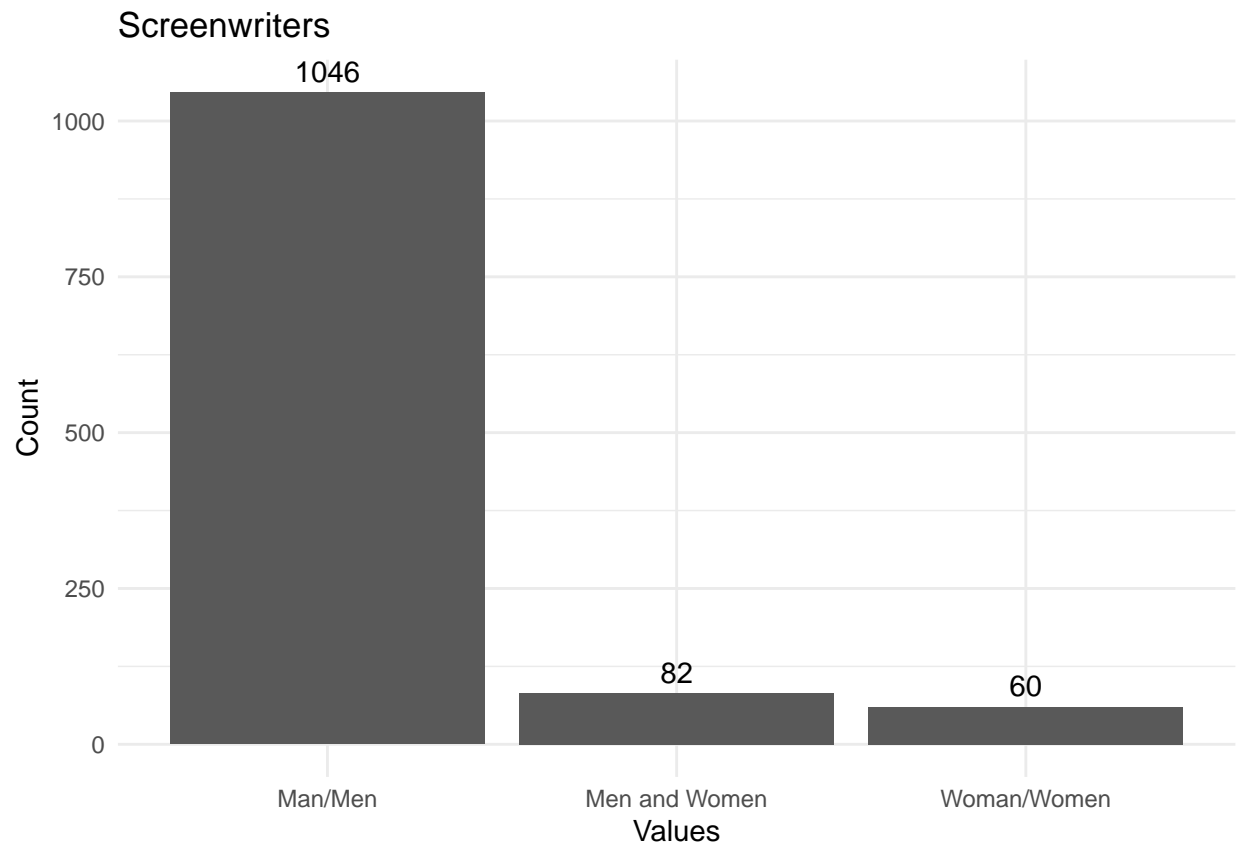


Writer Gender

```
describe(dfSpacy2$WriterGender) # Descriptives
```

```
##      vars      n mean  sd median trimmed mad min max range skew kurtosis  se
## X1*      1 1188 1.17 0.49      1    1.03  0  1  3      2 2.89      7.21 0.01
```

```
ggplot(dfSpacy2, aes(x = WriterGender)) + # Bar chart
  geom_bar() +
  geom_text(stat = "count",
            aes(label = after_stat(count)),
            vjust = -0.5) +
  labs(x = "Values", y = "Count", title = "Screenwriters") +
  scale_x_discrete(labels = c("Man/Men", "Men and Women", "Woman/Women")) +
  theme_minimal()
```



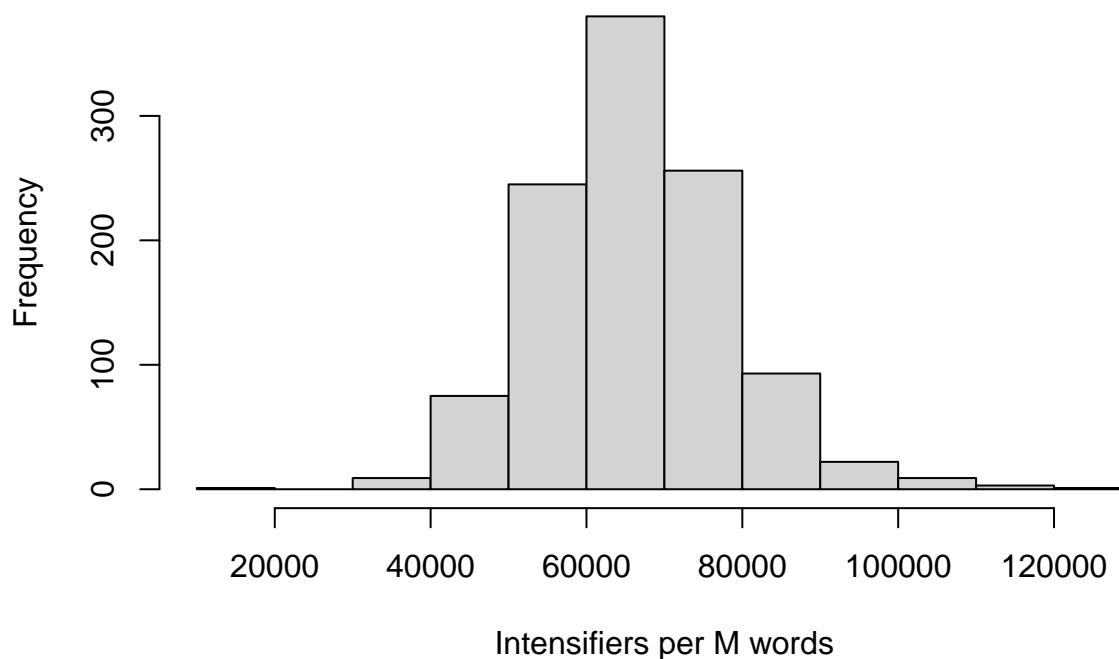
Intensifiers

```
describe(dfSpacy2$intensifiers) # Descriptives
```

```
##      vars      n    mean      sd  median trimmed    mad    min    max
## X1      1 1094 66358.61 12195.42 65511.27 66015.54 11583.49 13574.66 127423.8
##      range skew kurtosis      se
## X1 113849.2 0.41      1.29 368.71
```

```
hist(dfSpacy2$intensifiers, xlab = "Intensifiers per M words") # Histogram
```

Histogram of dfSpacy2\$intensifiers

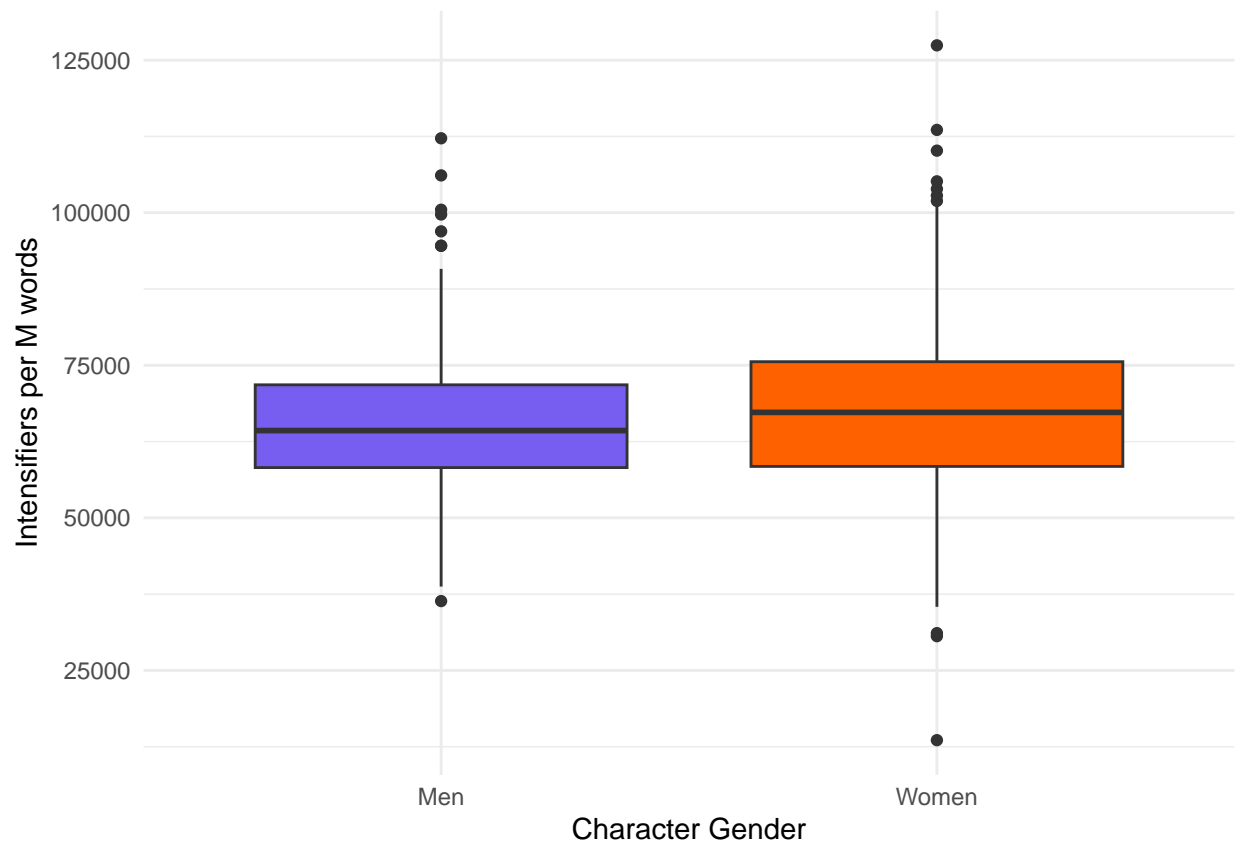


Plots

Intensifiers by Character Gender

```
ggplot(dfSpacy2, aes(x=CharIsWoman, y=intensifiers,  
                     fill=CharIsWoman)) +  
  geom_boxplot() + scale_x_discrete(labels = c("Men", "Women")) +  
  scale_fill_manual(values = c("#785EF0", "#FE6100")) +  
  labs(x = "Character Gender", y = "Intensifiers per M words") +  
  theme(axis.text = element_text(color = "black")) +  
  theme_minimal() + theme(legend.position = "none")
```

```
## Warning: Removed 94 rows containing non-finite outside the scale range  
## ('stat_boxplot()').
```



```
ggsave("IntByChar.svg") # Save plot
```

```
## Saving 6.5 x 4.5 in image
```

```
## Warning: Removed 94 rows containing non-finite outside the scale range
## ('stat_boxplot()').
```

```
dfSpacy2 %>%
  group_by(CharIsWoman) %>%
  summarise(median(intensifiers, na.rm = TRUE))
```

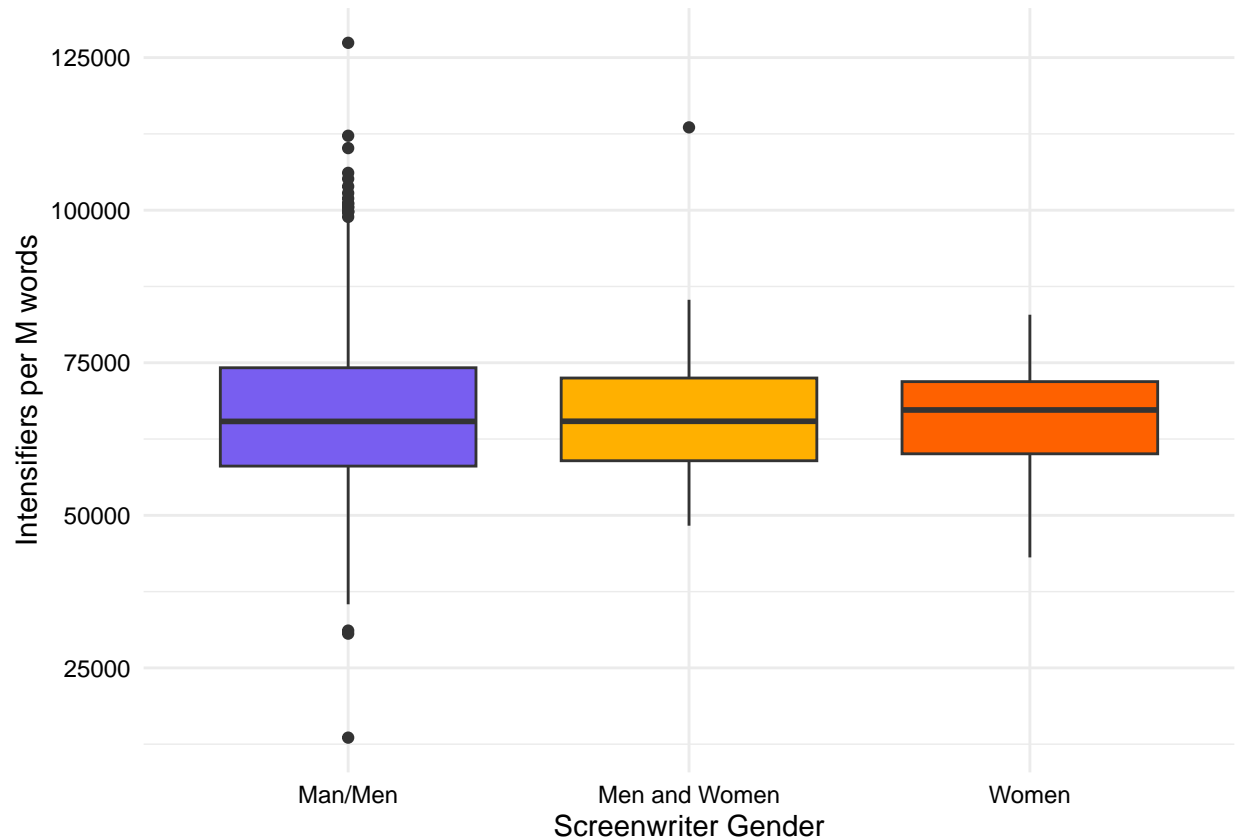
```
## # A tibble: 2 x 2
##   CharIsWoman 'median(intensifiers, na.rm = TRUE)'
##   <fct>                <dbl>
## 1 0                    64302.
## 2 1                    67280.
```

Intensifiers by Writer Gender

```
ggplot(dfSpacy2, aes(x=WriterGender, y=intensifiers,
  fill=as.factor(WriterGender))) +
  scale_fill_manual(values = c("#785EF0", "#FFB000", "#FE6100"))+
  geom_boxplot() +
  scale_x_discrete(labels = c("Man/Men", "Men and Women", "Women"))+
```

```
labs(x = "Screenwriter Gender", y = "Intensifiers per M words") +
theme_minimal() + theme(axis.text = element_text(color = "black")) +
theme(legend.position = "none")
```

```
## Warning: Removed 94 rows containing non-finite outside the scale range
## ('stat_boxplot()').
```



```
ggsave("IntByWrGender.svg") # Save plot
```

```
## Saving 6.5 x 4.5 in image
```

```
## Warning: Removed 94 rows containing non-finite outside the scale range
## ('stat_boxplot()').
```

```
dfSpacy2 %>%
  group_by(WriterGender) %>%
  summarise(median(intensifiers, na.rm = TRUE))
```

```
## # A tibble: 3 x 2
##   WriterGender 'median(intensifiers, na.rm = TRUE)'
##   <fct>                <dbl>
## 1 0                65391.
## 2 1                65410.
## 3 2                67263.
```

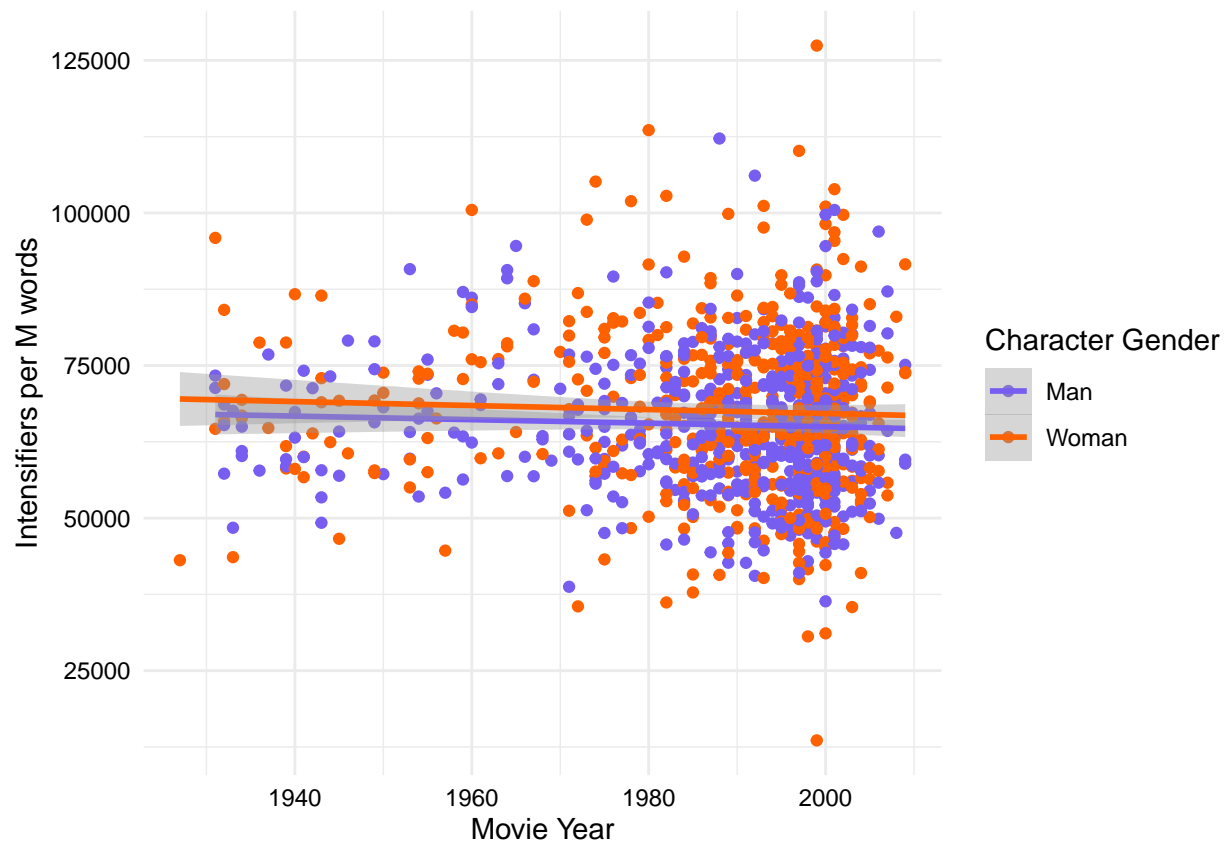
Intensifiers over time by character

```
ggplot(dfSpacy2, aes(x = MovieYear, y = intensifiers,
                     color = CharIsWoman)) +
  geom_point() +
  labs(color="Character Gender", x = "Movie Year",
       y = "Intensifiers per M words") +
  geom_smooth(method=lm) +
  scale_color_manual(values = c("#785EF0", "#FE6100"),
                    labels = c("Man", "Woman")) + theme_minimal() +
  theme(axis.text = element_text(color = "black"))
```

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

```
## Warning: Removed 94 rows containing non-finite outside the scale range
## ('stat_smooth()').
```

```
## Warning: Removed 94 rows containing missing values or values outside the scale range
## ('geom_point()').
```



```
ggsave("timeplot.svg")
```

```
## Saving 6.5 x 4.5 in image
## 'geom_smooth()' using formula = 'y ~ x'
```



```
## Warning: Removed 94 rows containing non-finite outside the scale range
## ('stat_smooth()').
## Removed 94 rows containing missing values or values outside the scale range
## ('geom_point()').
```

Statistical Models

Null model

```
nullmodel <- lmer (intensifiers~(1|MovieID), data = dfSpacy, REML = FALSE)
summary(nullmodel)

## Linear mixed model fit by maximum likelihood . t-tests use Satterthwaite's
## method [lmerModLmerTest]
## Formula: intensifiers ~ (1 | MovieID)
## Data: dfSpacy
##
##          AIC          BIC      logLik -2*log(L)  df.resid
##  23641.0    23656.0   -11817.5    23635.0      1091
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.7475 -0.5642 -0.0544  0.5350  4.4086
##
## Random effects:
## Groups   Name            Variance Std.Dev.
## MovieID (Intercept)  47648077   6903
## Residual                100840389 10042
## Number of obs: 1094, groups: MovieID, 594
##
## Fixed effects:
##              Estimate Std. Error      df t value Pr(>|t|)
## (Intercept)  66304.8      417.7    587.6   158.7  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# tab_model(nullmodel,, show.se=TRUE, show.ci=FALSE, show.aicc=TRUE)
```

Model with predictors

A random slopes model is not possible given the data.

```
# model1.1<- lmer (intensifiers~(CharIsWoman|MovieID) + CharIsWoman*MovieYear +WriterGender,
# data = dfSpacy2, REML = FALSE, lmerControl(autoscale = TRUE))

#Error: number of observations (=1094) <= number of random effects (=1188) for term (CharIsWoman | Movi
```

Random intercept model

```
model1.2 <- lmer (intensifiers~(1|MovieID) + CharIsWoman*MovieYear +WriterGender,
data = dfSpacy2, REML = FALSE, lmerControl(autoscale = TRUE))
summary(model1.2)
```

```
## Linear mixed model fit by maximum likelihood . t-tests use Satterthwaite's
## method [lmerModLmerTest]
## Formula: intensifiers ~ (1 | MovieID) + CharIsWoman * MovieYear + WriterGender
## Data: dfSpacy2
## Control: lmerControl(autoscale = TRUE)
##
##      AIC      BIC    logLik -2*log(L)  df.resid
## 23636.7 23676.7 -11810.4 23620.7    1086
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.8557 -0.5863 -0.0375  0.5297  4.3287
##
## Random effects:
## Groups Name Variance Std.Dev.
## MovieID (Intercept) 48489464 6963
## Residual 98506834 9925
## Number of obs: 1094, groups: MovieID, 594
##
## Fixed effects:
##              Estimate Std. Error      df t value Pr(>|t|)
## (Intercept) 6.632e+04 6.194e+04 2.865e+11 1.071 0.284
## CharIsWoman1 8.774e+03 7.325e+04 8.774e+03 0.120 0.905
## MovieYear -4.346e+02 3.115e+01 1.369e-02 -13.951 0.928
## WriterGender1 -1.477e+02 1.632e+03 1.302e+05 -0.090 0.928
## WriterGender2 -6.922e+01 1.936e+03 2.873e+05 -0.036 0.971
## CharIsWoman1:MovieYear -7.681e+03 3.683e+01 5.609e-10 -208.515 1.000
```

```
# tab_model(model1.2, df.method = "satterthwaite", show.se=TRUE, show.ci=FALSE,
# show.icc=FALSE, show.aicc=TRUE)
```

Model assumptions

```
plot_model(model1.2, type='diag')
```

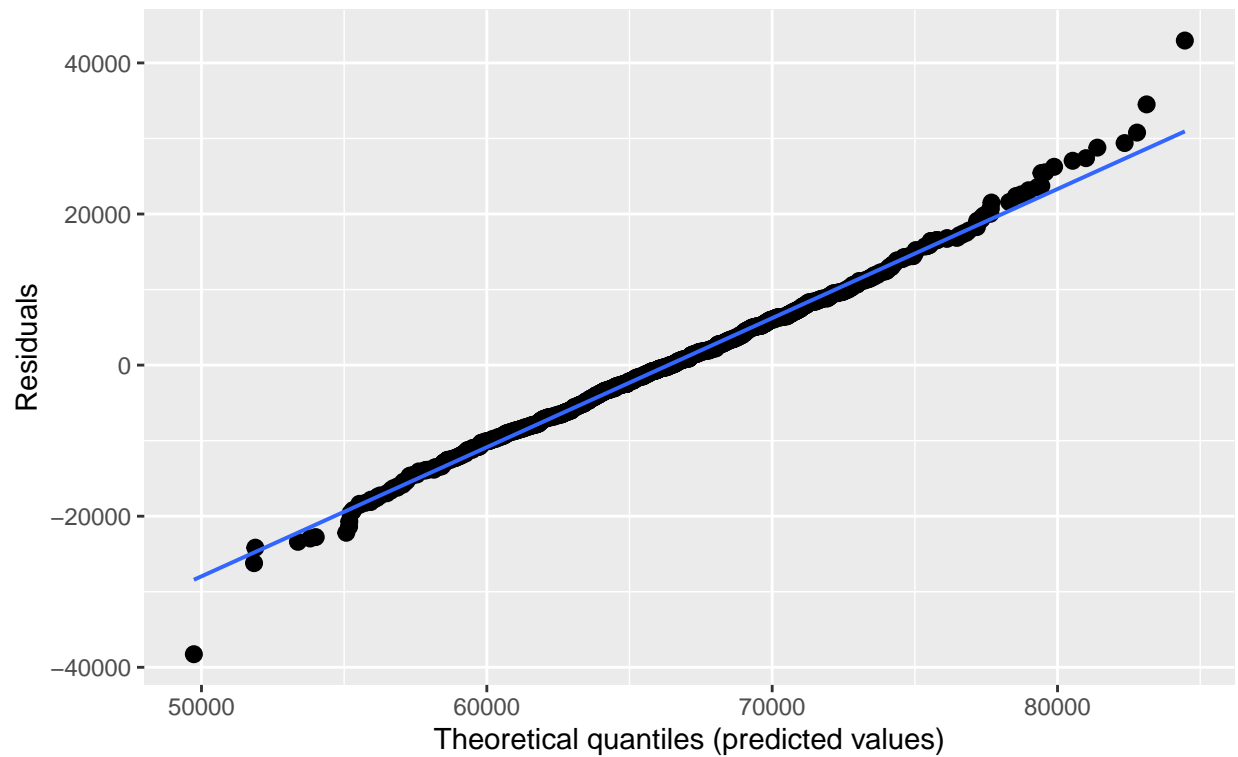
```
## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## i The deprecated feature was likely used in the sjPlot package.
## Please report the issue at <https://github.com/strengejacke/sjPlot/issues>.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

```
## [[1]]
```

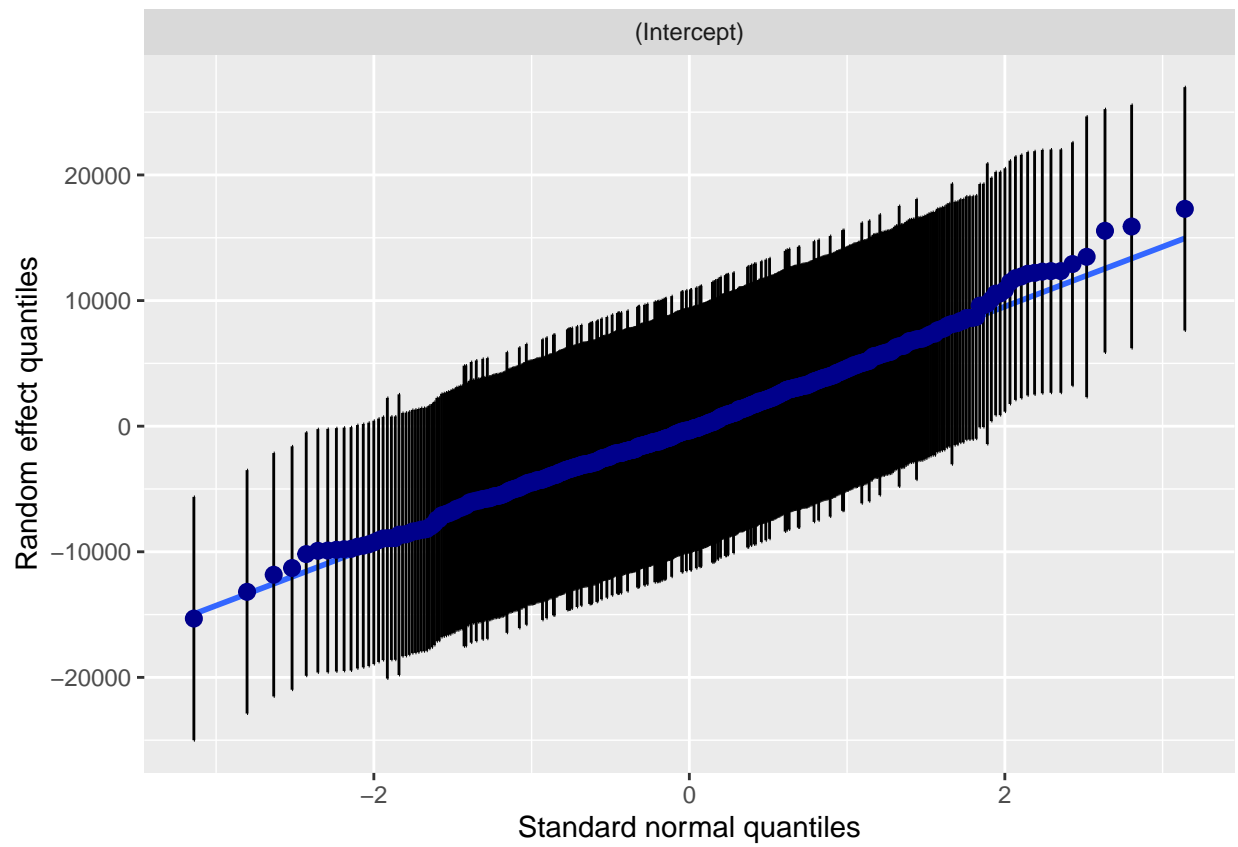
```
## 'geom_smooth()' using formula = 'y ~ x'
```

Non-normality of residuals and outliers

Dots should be plotted along the line



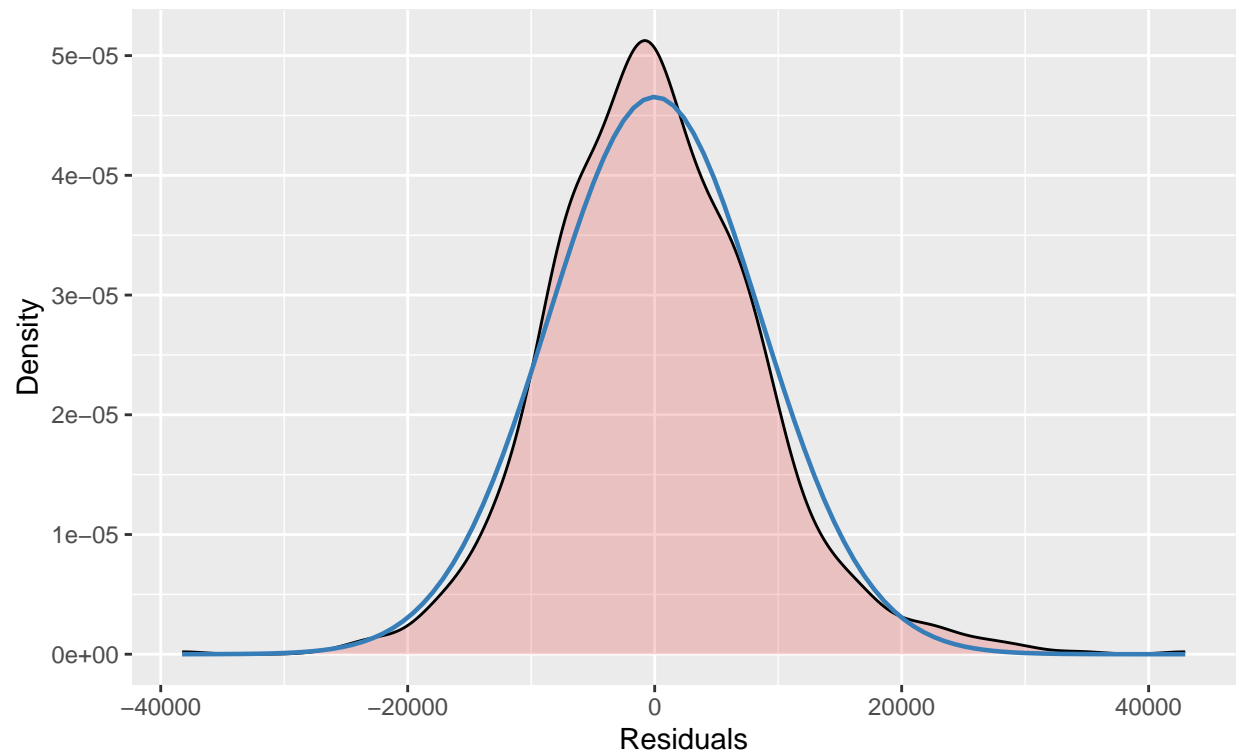
```
##  
## [[2]]  
## [[2]]$MovieID  
  
## 'geom_smooth()' using formula = 'y ~ x'
```



```
##  
##  
## [[3]]
```

Non-normality of residuals

Distribution should look like normal curve



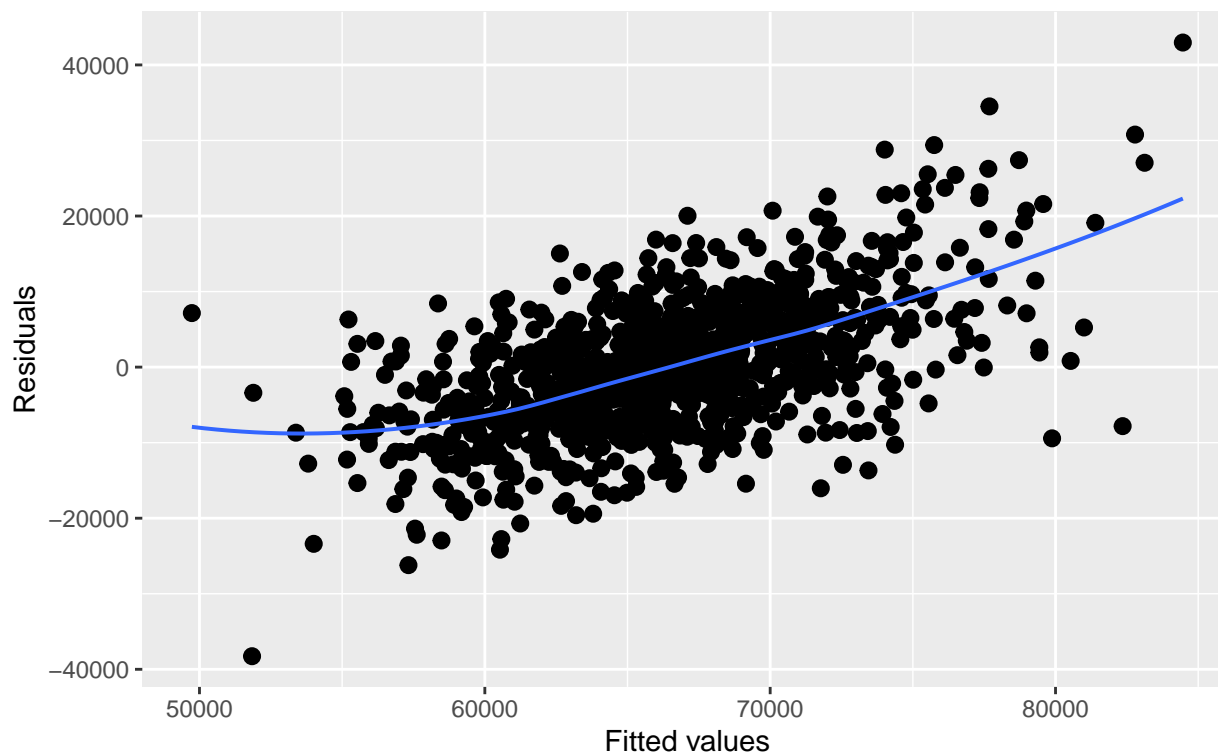
```
##
```

```
## [[4]]
```

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

Homoscedasticity (constant variance of residuals)

Amount and distance of points scattered above/below line is equal or randomly spread



As residuals are not homoscedastic, try running a model using robust lmer. Overall conclusions are the same, so retain model 1.2.

```
model2<- rlmer(intensifiers~(1|MovieID)+CharIsWoman*MovieYear+WriterGender,  
              data = dfSpacy2,  
              REML = FALSE, lmerControl(autoscale = TRUE))  
summary(model2)
```

```
## Robust linear mixed model fit by DASTau  
## Formula: intensifiers ~ (1 | MovieID) + CharIsWoman * MovieYear + WriterGender  
## Data: dfSpacy2  
##  
## Scaled residuals:  
##      Min       1Q   Median       3Q      Max   
## -4.8302 -0.6120 -0.0380  0.6107  5.5142   
##  
## Random effects:  
## Groups Name Variance Std.Dev.  
## MovieID (Intercept) 23198000 4816  
## Residual 104288686 10212  
## Number of obs: 1094, groups: MovieID, 594  
##  
## Fixed effects:  
##  
##              Estimate Std. Error t value  
## (Intercept) 66065.3 377.7 174.91  
## CharIsWoman1 9005.2 38436.3 0.23
```

```

## MovieYear          -479.5      491.7   -0.98
## WriterGender1      -161.0      379.8   -0.42
## WriterGender2       105.7      378.0    0.28
## CharIsWoman1:MovieYear -7921.8   38438.4  -0.21
##
## Correlation of Fixed Effects:
##           (Intr) ChrIW1 MoviYr WrtrG1 WrtrG2
## CharIsWomn1  0.001
## MovieYear    0.001  0.636
## WriterGndr1  0.002 -0.009  0.010
## WriterGndr2 -0.002 -0.021 -0.061  0.061
## ChrIsWm1:MY  0.000 -1.000 -0.636  0.009  0.021
##
## Robustness weights for the residuals:
## 912 weights are ~= 1. The remaining 182 ones are summarized as
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.244  0.632  0.830  0.769  0.936  0.999
##
## Robustness weights for the random effects:
## 588 weights are ~= 1. The remaining 6 ones are
##   30  147  316  401  410  536
## 0.993 0.996 0.991 0.996 0.999 0.994
##
## Rho functions used for fitting:
##   Residuals:
##     eff: smoothed Huber (k = 1.345, s = 10)
##     sig: smoothed Huber, Proposal 2 (k = 1.345, s = 10)
##   Random Effects, variance component 1 (MovieID):
##     eff: smoothed Huber (k = 1.345, s = 10)
##     vcp: smoothed Huber, Proposal 2 (k = 1.345, s = 10)

# tab_model(model2, show.se=TRUE, df.method = "satterthwaite", show.ci=FALSE,
# show.icc=FALSE) view model

```