

Midterm Project Report

Runxin Yu

12/6/2019

Introduction

The dataset is called **Movie Industry: Three Decades of Movies** that contains 6820 movies in total (220 movies per year from 1986-2016) collected from 57 regions all over the world. The goal is to see how score of a movie is affected by different variables such as budget of the movie, revenue of the movie, duration of the movie, etc.

Columns of the Dataset

Each movie has the following attributes:

- **budget**: the budget of a movie. Some movies don't have this, so it appears as 0
- **company**: the production company
- **country**: country of origin
- **director**: the director
- **genre**: main genre of the movie
- **gross**: revenue of the movie
- **name**: name of the movie
- **rating**: rating of the movie (R, PG, etc.)
- **released**: release date (YYYY-MM-DD)
- **runtime**: duration of the movie
- **score**: IMDb user rating
- **votes**: number of user votes
- **star**: main actor/actress
- **writer**: writer of the movie
- **year**: year of release

Exploratory Data Analysis(EDA)

Data Cleaning

Reset the cells with a budget of 0 into NA and remove those rows from the data frame. New data frame without missing values is created in order for further exploration. Notice that some of the movies are produced by the same film company but are named differently, rename those companies so that their names are consistent. Another thing observed from the data is that there are 2179 film production companies, which is too large to set to be a random effect. Therefore, 7 of the world most famous film production companies, Twentieth Century Fox, Columbia Pictures, Universal Pictures, Warner Bros., Paramount Pictures, Walt Disney, and Metro-Goldwyn-Mayer, as well as companies with counts greater than 50 are kept their original names, otherwise others is being labeled for the **company** column. This process reduces that number of different type of companies into 16, which can be used as a random effect.

##	budget	company
## Min.	: 6000	Universal Pictures : 265
## 1st Qu.	: 10000000	Warner Bros. : 257
## Median	: 23000000	Paramount Pictures : 222
## Mean	: 36145602	Twentieth Century Fox Film Corporation: 174

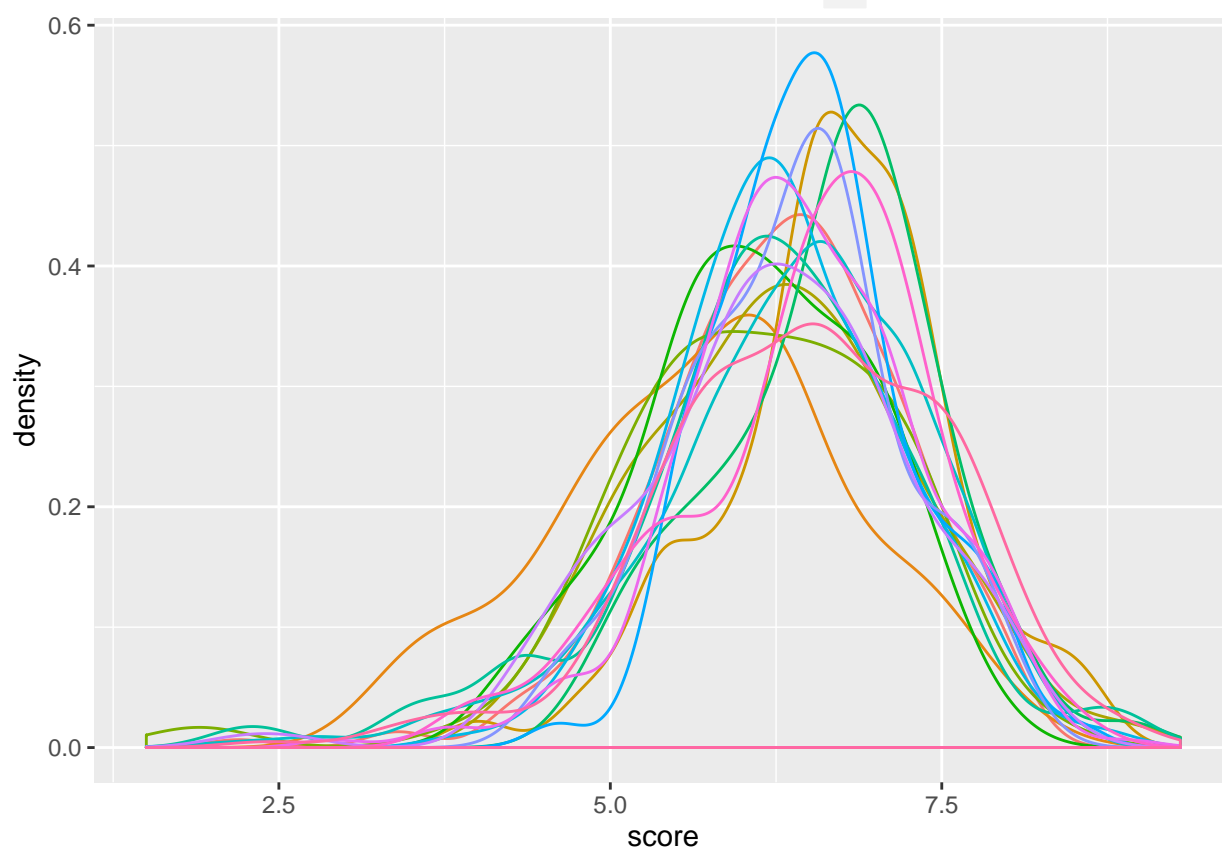
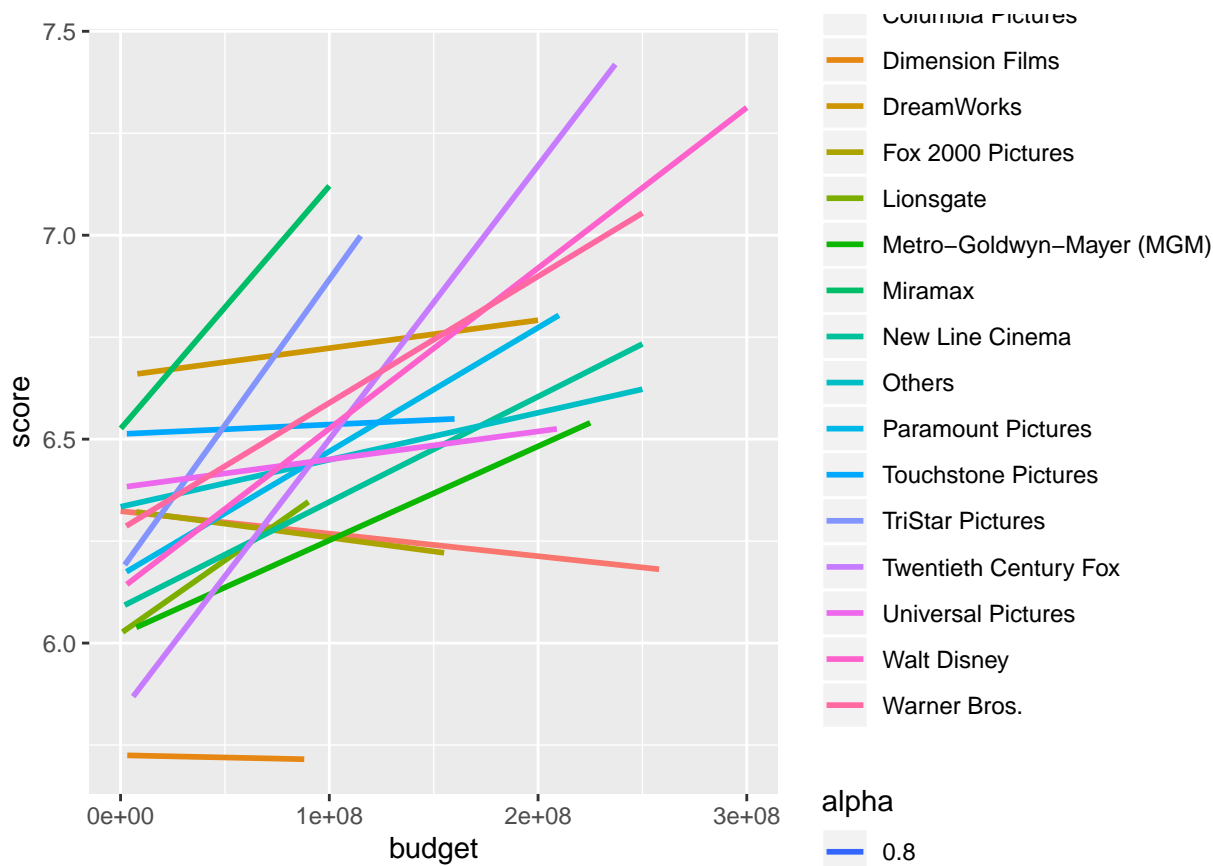
```

## 3rd Qu.: 46000000 New Line Cinema : 145
## Max. :300000000 Columbia Pictures Corporation : 137
## (Other) :3438
## country director genre gross
## USA :3726 Woody Allen : 30 Comedy :1310 Min. : 309
## UK : 366 Clint Eastwood : 24 Action :1099 1st Qu.: 6290905
## France : 108 Steven Soderbergh: 21 Drama : 793 Median : 23455506
## Germany : 93 Steven Spielberg : 21 Crime : 356 Mean : 46074694
## Canada : 79 Ron Howard : 20 Adventure: 291 3rd Qu.: 57782434
## Australia: 36 Ridley Scott : 19 Biography: 239 Max. :936662225
## (Other) : 230 (Other) :4503 (Other) : 550
## name rating released runtime
## Bulletproof : 2 R :2247 2007-10-19: 8 Min. : 69.0
## Cocktail : 2 PG-13 :1561 2012-04-20: 8 1st Qu.: 96.0
## Deadfall : 2 PG : 659 2000-10-20: 7 Median :104.0
## Death at a Funeral: 2 G : 100 2002-10-11: 7 Mean :107.6
## Fair Game : 2 NOT RATED: 38 2003-11-14: 7 3rd Qu.:117.0
## Fantastic Four : 2 UNRATED : 20 2008-12-25: 7 Max. :280.0
## (Other) :4626 (Other) : 13 (Other) :4594
## score star votes writer
## Min. :1.500 Nicolas Cage : 38 Min. : 183 Woody Allen : 29
## 1st Qu.:5.800 Bruce Willis : 33 1st Qu.: 16110 Stephen King: 20
## Median :6.400 Robert De Niro : 32 Median : 43940 John Hughes : 14
## Mean :6.356 Denzel Washington: 31 Mean : 95702 Luc Besson : 13
## 3rd Qu.:7.100 Tom Hanks : 31 3rd Qu.: 109393 Wes Craven : 12
## Max. :9.300 Johnny Depp : 30 Max. :1861666 Joel Coen : 11
## (Other) :4443 (Other) :4539
## year
## Min. :1986
## 1st Qu.:1996
## Median :2003
## Mean :2002
## 3rd Qu.:2010
## Max. :2016
##

```

Visually Explore the Data

The plot below shows that there is a relationship between budget and score varying by different film production companies. The density plot shows that the score follows a approximately normal distribution.



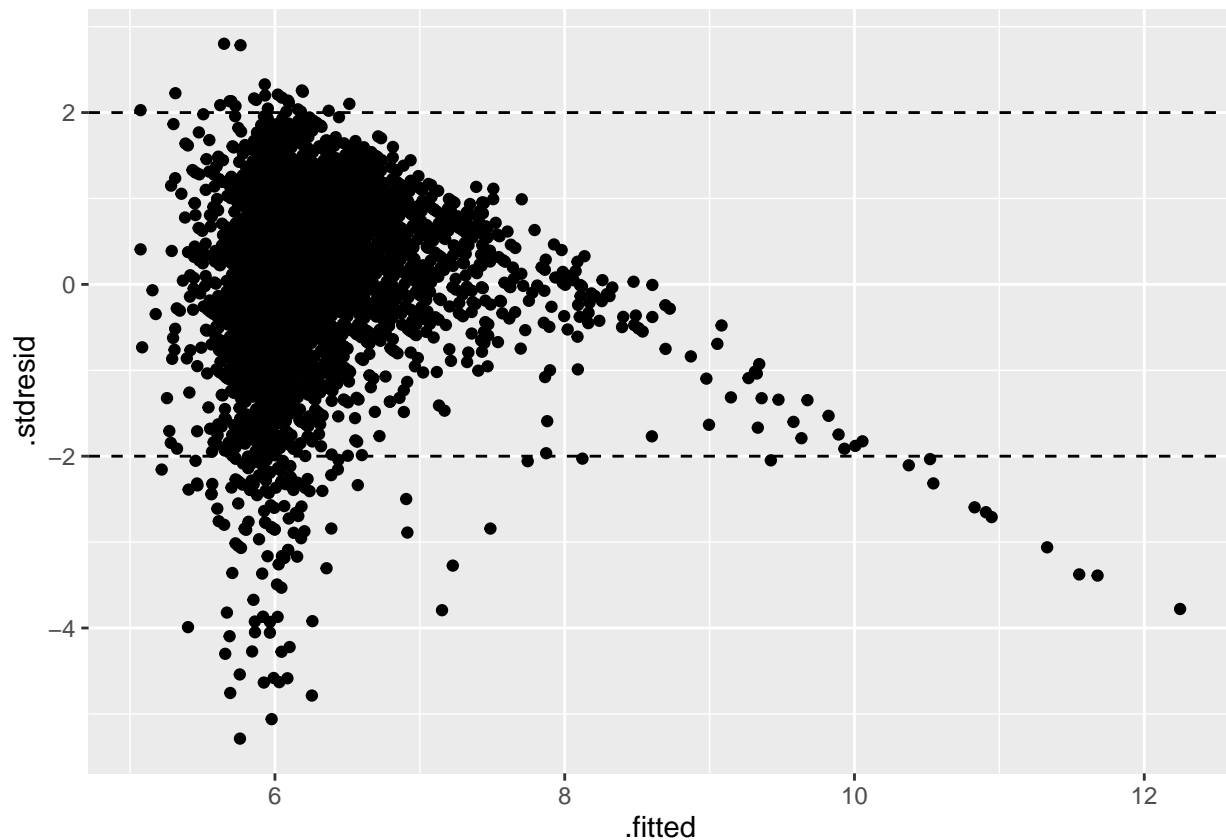
Though there are 57 regions in total in the data set, more than 80% of the data is collected from the United State. Therefore, the main focus is the film industries in the United State.

Linear Regression

In order to see the relationships between ratings of movies and variables that may affect their scores furthermore, a simple linear regression is fitted to the original data.

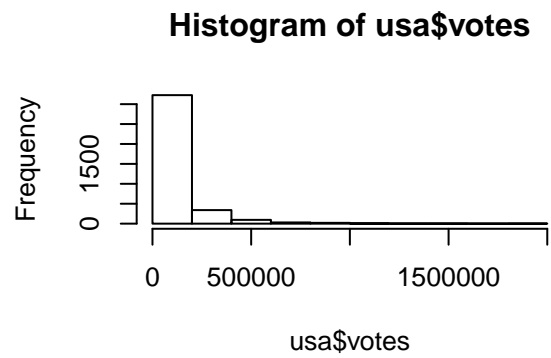
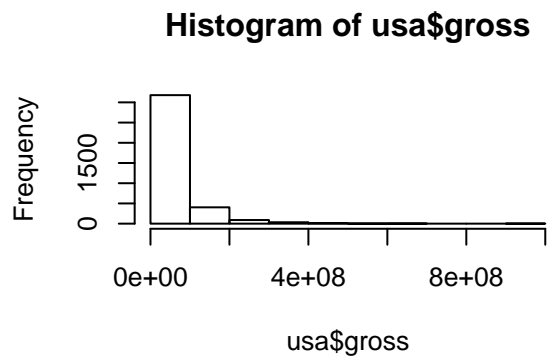
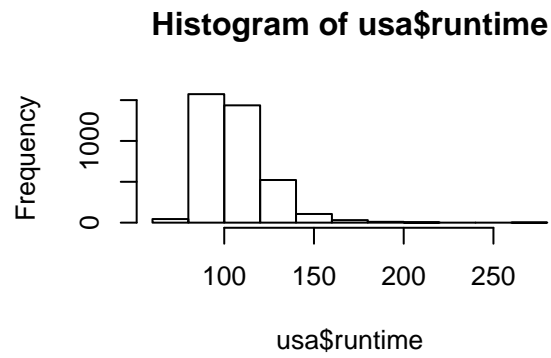
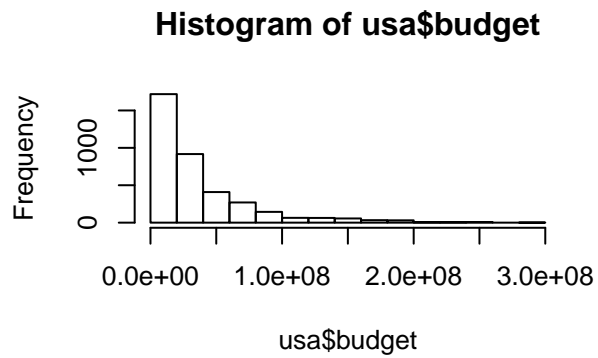
```
##
## Call:
## lm(formula = score ~ budget + runtime + gross + votes + factor(company),
##     data = usa)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.2589 -0.4012  0.0807  0.5327  2.2510
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.295e+00  1.036e-01  41.455 < 2e-16
## budget        -6.143e-09  4.673e-10 -13.145 < 2e-16
## runtime        1.746e-02  8.555e-04  20.405 < 2e-16
## gross          7.732e-10  3.088e-10   2.504 0.012326
## votes          2.986e-06  1.158e-07  25.789 < 2e-16
## factor(company)Dimension Films -4.906e-01  1.297e-01  -3.782 0.000158
## factor(company)DreamWorks      1.493e-01  1.098e-01   1.360 0.173984
## factor(company)Fox 2000 Pictures  1.126e-02  1.261e-01   0.089 0.928813
## factor(company)Lionsgate       -2.433e-01  1.297e-01  -1.876 0.060738
## factor(company)Metro-Goldwyn-Mayer (MGM) -2.297e-02  1.105e-01  -0.208 0.835287
## factor(company)Miramax          1.841e-01  1.277e-01   1.442 0.149443
## factor(company)New Line Cinema  -1.657e-01  9.002e-02  -1.840 0.065797
## factor(company)Others           4.698e-02  5.755e-02   0.816 0.414394
## factor(company)Paramount Pictures  2.747e-02  7.824e-02   0.351 0.725489
## factor(company)Touchstone Pictures 1.493e-01  1.019e-01   1.465 0.142894
## factor(company)TriStar Pictures  1.835e-01  1.234e-01   1.486 0.137294
## factor(company)Twentieth Century Fox -1.146e-01  8.237e-02  -1.391 0.164220
## factor(company)Universal Pictures  2.891e-02  7.656e-02   0.378 0.705715
## factor(company)Walt Disney       5.200e-01  9.778e-02   5.318 1.11e-07
## factor(company)Warner Bros.      -3.997e-02  7.546e-02  -0.530 0.596324
##
## (Intercept) ***
## budget      ***
## runtime     ***
## gross       *
## votes       ***
## factor(company)Dimension Films ***
## factor(company)DreamWorks
## factor(company)Fox 2000 Pictures
## factor(company)Lionsgate .
## factor(company)Metro-Goldwyn-Mayer (MGM)
## factor(company)Miramax
## factor(company)New Line Cinema .
## factor(company)Others
## factor(company)Paramount Pictures
```

```
## factor(company)Touchstone Pictures
## factor(company)TriStar Pictures
## factor(company)Twentieth Century Fox
## factor(company)Universal Pictures
## factor(company)Walt Disney          ***
## factor(company)Warner Bros.
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.806 on 3706 degrees of freedom
## Multiple R-squared:  0.3527, Adjusted R-squared:  0.3494
## F-statistic: 106.3 on 19 and 3706 DF,  p-value: < 2.2e-16
```

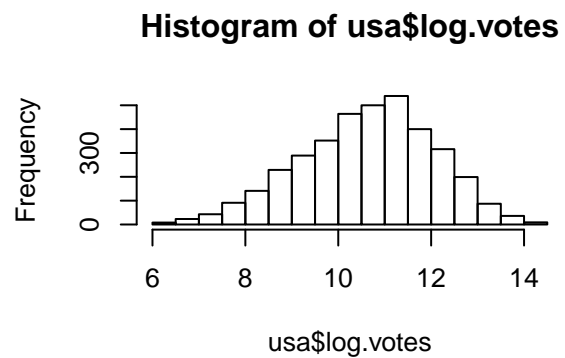
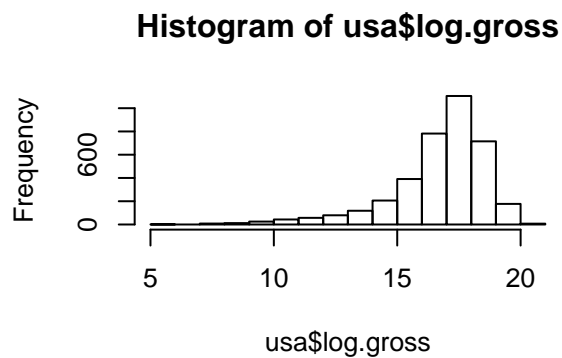
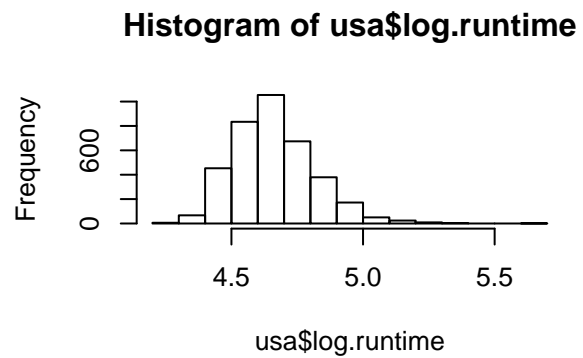
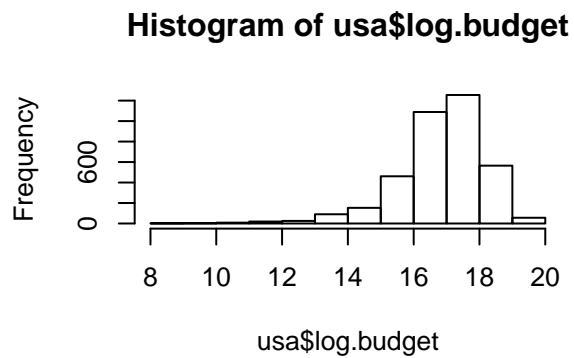


The residual plot generating from the simple linear regression shows both **Heteroscedasticity** and **Nonlinear** issue and the R-square of approximately 0.35 is not very promising. This pattern indicates that transformation to some of the variables is needed.

The histogram of variables with very large skewness also indicates that a log transformation is needed.



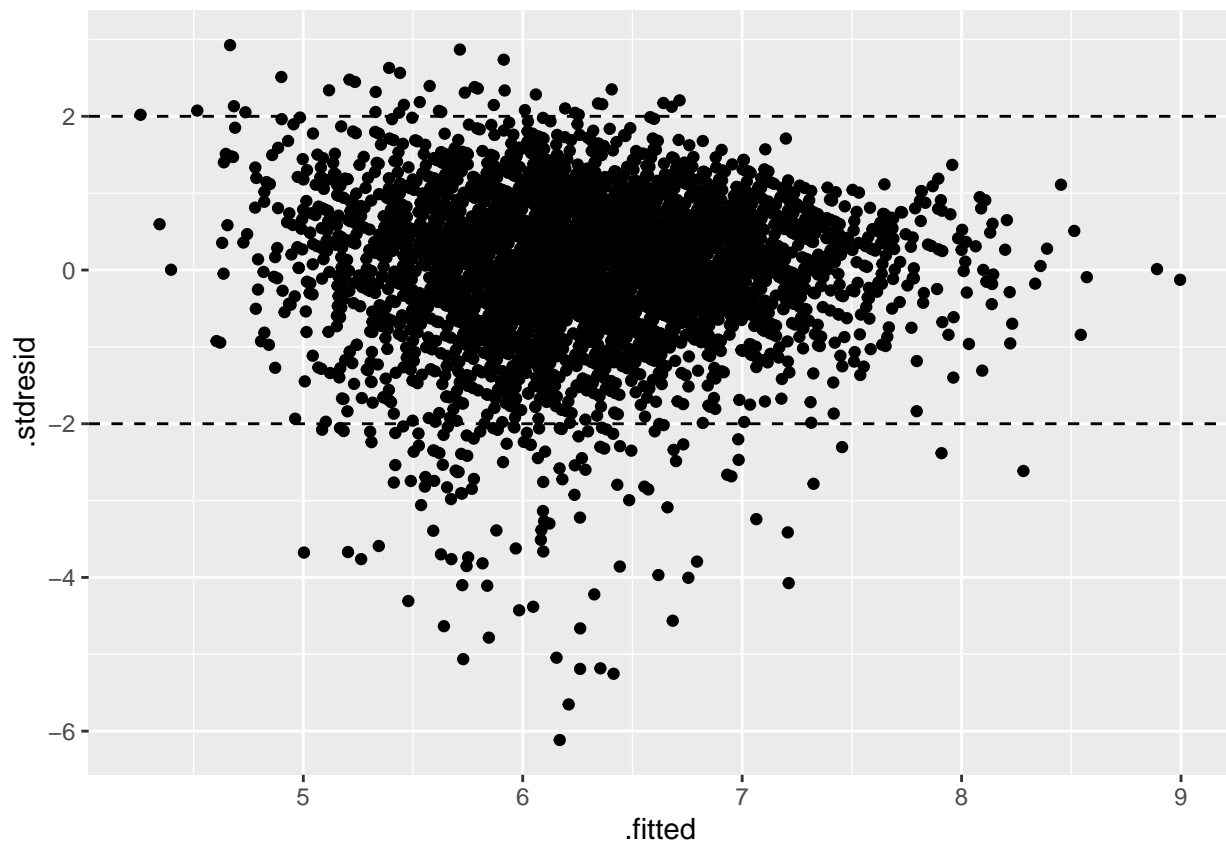
Thus, log transformations are made to budget, runtime, gross and votes and fit the model again.



```
##
## Call:
## lm(formula = score ~ log.budget + log.runtime + log.gross + log.votes +
##     factor(company), data = usa)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.6686 -0.3902  0.0777  0.4949  2.2338
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   -4.239175    0.404985 -10.467  < 2e-16
## log.budget                     -0.257291    0.012944 -19.877  < 2e-16
## log.runtime                     2.325874    0.092142  25.242  < 2e-16
## log.gross                      -0.015984    0.009302  -1.718  0.085808
## log.votes                      0.396036    0.011591  34.167  < 2e-16
## factor(company)Dimension Films -0.471645    0.122996  -3.835  0.000128
## factor(company)DreamWorks      0.188629    0.104137   1.811  0.070166
## factor(company)Fox 2000 Pictures 0.114303    0.119533   0.956  0.339009
## factor(company)Lionsgate       -0.292396    0.122952  -2.378  0.017452
## factor(company)Metro-Goldwyn-Mayer (MGM) 0.155117    0.104926   1.478  0.139399
## factor(company)Miramax         0.284653    0.121070   2.351  0.018768
## factor(company)New Line Cinema -0.127715    0.085264  -1.498  0.134250
## factor(company)Others          0.105126    0.055286   1.901  0.057316
## factor(company)Paramount Pictures 0.114210    0.074196   1.539  0.123813
## factor(company)Touchstone Pictures 0.206587    0.096579   2.139  0.032498
## factor(company)TriStar Pictures 0.360482    0.117079   3.079  0.002092
## factor(company)Twentieth Century Fox -0.046065    0.078109  -0.590  0.555389
## factor(company)Universal Pictures 0.055603    0.072595   0.766  0.443759
## factor(company)Walt Disney      0.528495    0.091842   5.754  9.4e-09
## factor(company)Warner Bros.     0.017742    0.071510   0.248  0.804069
##
## (Intercept) ***
## log.budget ***
## log.runtime ***
## log.gross .
## log.votes ***
## factor(company)Dimension Films ***
## factor(company)DreamWorks .
## factor(company)Fox 2000 Pictures
## factor(company)Lionsgate *
## factor(company)Metro-Goldwyn-Mayer (MGM)
## factor(company)Miramax *
## factor(company)New Line Cinema
## factor(company)Others .
## factor(company)Paramount Pictures
## factor(company)Touchstone Pictures *
## factor(company)TriStar Pictures **
## factor(company)Twentieth Century Fox
## factor(company)Universal Pictures
## factor(company)Walt Disney ***
## factor(company)Warner Bros.
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.7646 on 3706 degrees of freedom
## Multiple R-squared:  0.4174, Adjusted R-squared:  0.4144
## F-statistic: 139.8 on 19 and 3706 DF,  p-value: < 2.2e-16
```

The new model after log transformation is slightly better than the one one without log transformation by looking at the residual plot, though it shows a large variance in the middle portion of the residual plot, and the R-square goes up to 0.41. However, improving of the model fitting is still needed.



Multilevel Model

Random Intercept Model

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: score ~ log.budget + (1 | company)
## Data: usa
##
## REML criterion at convergence: 10559
##
## Scaled residuals:
##   Min      1Q  Median      3Q      Max
## -4.7222 -0.5807  0.0627  0.6831  3.0431
##
## Random effects:
##   Groups   Name      Variance Std.Dev.
##   company  (Intercept) 0.03185  0.1785
```

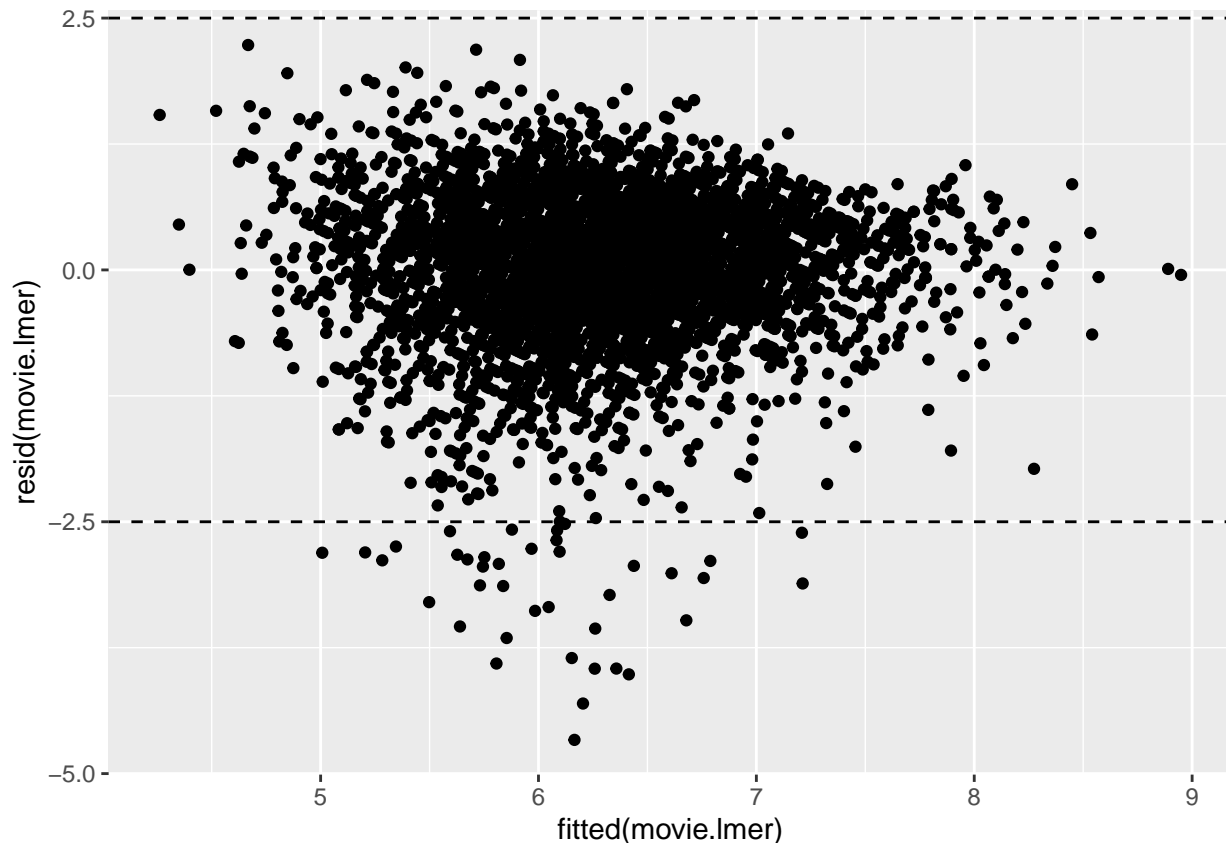


```
## Residual          0.98675  0.9934
## Number of obs: 3726, groups:  company, 16
##
## Fixed effects:
##           Estimate Std. Error t value
## (Intercept)  5.93165    0.23668  25.062
## log.budget   0.02207    0.01343   1.643
##
## Correlation of Fixed Effects:
##           (Intr)
## log.budget -0.976

## Analysis of Deviance Table (Type II Wald chisquare tests)
##
## Response: score
##           Chisq Df Pr(>Chisq)
## log.budget 2.7002  1    0.1003
```

A simple random intercept model is fitted with log.budget coefficient as well as its t value very small, that is budget(in log scale) does not have many effects on movie scores. However, the plot shows that there is some relationship between. A very large AIC, which indicates that this is not a good model. Thus, other variables are added to the model to see whether will improve the fit or not.

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: score ~ log.budget + log.gross + log.runtime + log.votes + (1 |
##           company)
##           Data: usa
##
## REML criterion at convergence: 8636.1
##
## Scaled residuals:
##           Min       1Q   Median       3Q      Max
## -6.1006 -0.5100  0.1014  0.6445  2.9198
##
## Random effects:
## Groups   Name                Variance Std.Dev.
## company (Intercept) 0.04513  0.2124
## Residual              0.58476  0.7647
## Number of obs: 3726, groups:  company, 16
##
## Fixed effects:
##           Estimate Std. Error t value
## (Intercept) -4.188854    0.403124 -10.391
## log.budget   -0.256453    0.012927 -19.839
## log.gross    -0.015298    0.009294  -1.646
## log.runtime   2.329151    0.092008  25.315
## log.votes     0.394604    0.011577  34.085
##
## Correlation of Fixed Effects:
##           (Intr) lg.bdg lg.grs lg.rnt
## log.budget  -0.098
## log.gross    -0.073 -0.404
## log.runtime -0.882 -0.245  0.037
## log.votes     0.138 -0.101 -0.453 -0.202
```



```
##              2.5 %      97.5 %
## .sig01      0.13393619  0.319571415
## .sigma      0.74723017  0.782032231
## (Intercept) -4.97975357 -3.400856761
## log.budget  -0.28172472 -0.231055920
## log.gross   -0.03346486  0.002967442
## log.runtime  2.14905545  2.509602340
## log.votes   0.37179554  0.417208341
```

The 95% confidence interval shows that every fixed effect estimates are significant besides log.gross. This shows that budget of the movie(in log scale), runtime of the movie(in log scale) and number of votes(in log scale) will affect the score of a movie. Among all, it is interesting that runtime have a positive relationship with the score of the movie. It is also interesting that with every unit increase in runtime(in log scale), the score of a movie will be expected to increase by 2.3, which is a significant amount of increase. The variance of the random effect company is close to 0, which means that there are no huge differences in scores among different film production companies. The residual plot still shows a large variance in the middle portion.

Notice that genre can be another group factor, so can be added to the model.

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: score ~ log.budget + log.gross + log.runtime + log.votes + (1 |
##      company) + (1 | genre)
##      Data: usa
##
## REML criterion at convergence: 8148.9
##
## Scaled residuals:
##      Min      1Q  Median      3Q      Max
```

```
## -6.5269 -0.4972  0.1019  0.6310  2.6257
##
## Random effects:
##   Groups   Name                Variance Std.Dev.
##   company   (Intercept) 0.01698  0.1303
##   genre     (Intercept) 0.13569  0.3684
##   Residual                   0.50906  0.7135
## Number of obs: 3726, groups:  company, 16; genre, 15
##
## Fixed effects:
##               Estimate Std. Error t value
## (Intercept) -3.685804   0.421810  -8.738
## log.budget  -0.269802   0.012770 -21.128
## log.gross   -0.010879   0.008793  -1.237
## log.runtime  2.267491   0.097238  23.319
## log.votes    0.393985   0.010915  36.097
##
## Correlation of Fixed Effects:
##              (Intr) lg.bdg lg.grs lg.rnt
## log.budget  -0.062
## log.gross   -0.011 -0.367
## log.runtime -0.865 -0.283 -0.026
## log.votes    0.126 -0.077 -0.447 -0.191
##
##              2.5 %      97.5 %
## .sig01        0.06978558  0.220318405
## .sig02        0.24599042  0.559667935
## .sigma        0.69713829  0.729669620
## (Intercept) -4.51241103 -2.861387821
## log.budget  -0.29474868 -0.244655973
## log.gross   -0.02812394  0.006351319
## log.runtime  2.07680344  2.457991011
## log.votes    0.37253400  0.415395863
```

An estimate of the variance of 0.14 explained by the random effect genre indicates that scores varies slightly between different genres. The model still gives the same significant fixed effect estimates: budget(in log scale), runtime(in log scale),votes(in log scale).

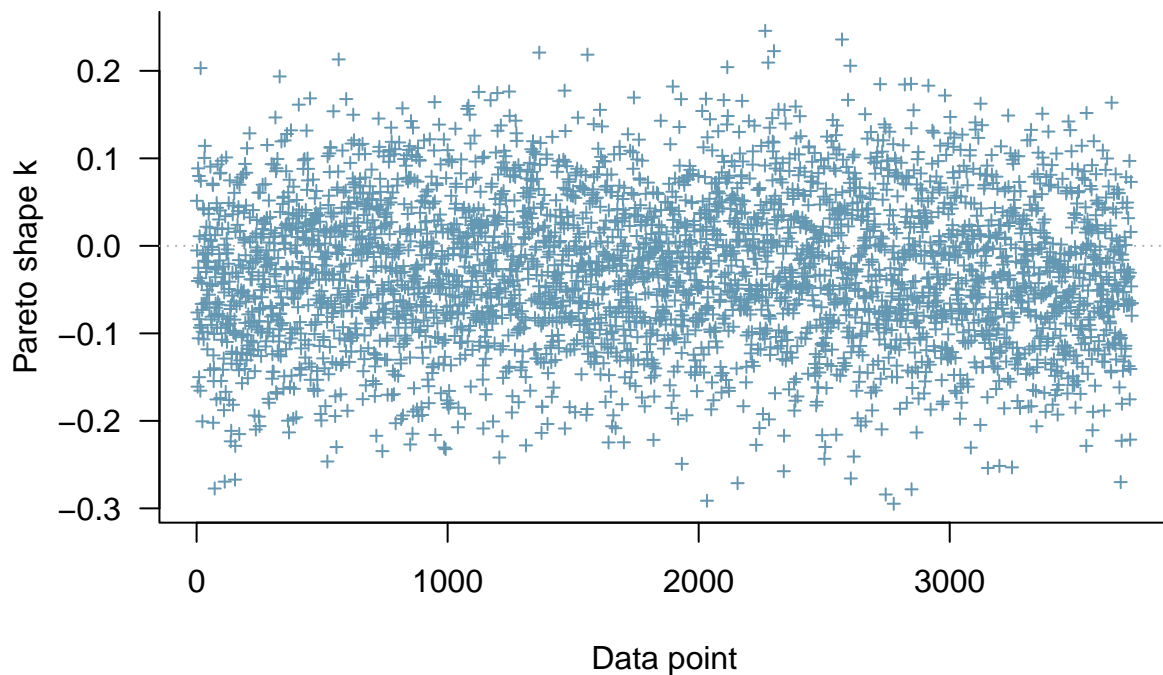
Stan to Fit Multilevel Model

```
##
##               5%          95%
## (Intercept)   -4.885898653 -3.5234461367
## log.budget    -0.278308146 -0.2350202854
## log.gross     -0.030743264 -0.0001811552
## log.runtime    2.180064920  2.4858353919
## log.votes      0.375670642  0.4136054038
## b[(Intercept) company:Columbia_Pictures] -0.193238982  0.0571164564
## b[(Intercept) company:Dimension_Films]   -0.626637986 -0.2436995030
## b[(Intercept) company:DreamWorks]        -0.058407868  0.2613052266
## b[(Intercept) company:Fox_2000_Pictures] -0.140841754  0.2083785128
## b[(Intercept) company:Lionsgate]         -0.477574726 -0.0994943226
## b[(Intercept) company:Metro-Goldwyn-Mayer_(MGM)] -0.097971316  0.2372436043
## b[(Intercept) company:Miramax]           -0.008225794  0.3476797543
## b[(Intercept) company:New_Line_Cinema]    -0.323089819 -0.0425002939
```

```
## b[(Intercept) company:Others] -0.069074806 0.1368039253
## b[(Intercept) company:Paramount_Pictures] -0.087416674 0.1699721867
## b[(Intercept) company:Touchstone_Pictures] -0.032807954 0.2722822969
## b[(Intercept) company:TriStar_Pictures] 0.057247308 0.4193454547
## b[(Intercept) company:Twentieth_Century_Fox] -0.243385088 0.0240642295
## b[(Intercept) company:Universal_Pictures] -0.140388142 0.1094747822
## b[(Intercept) company:Walt_Disney] 0.252809753 0.5632802792
## b[(Intercept) company:Warner_Bros.] -0.174721606 0.0738011769
## sigma 0.750276378 0.7793443945
## Sigma[company:(Intercept),(Intercept)] 0.023289448 0.1044169200

##
## Computed from 4000 by 3726 log-likelihood matrix
##
## Estimate SE
## elpd_loo -4298.1 68.6
## p_loo 19.8 1.1
## looic 8596.2 137.2
## -----
## Monte Carlo SE of elpd_loo is 0.1.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.
```

PSIS diagnostic plot



Stan is fitted to the model and loo package is used to check goodness of fit of the model. The table above shows a summary of Pareto k diagnostic with a Pareto k estimates getting from loo function less than 0.5 and a Monte Carlo SE of elpd_loo 0.1, which kind of gives the sence that this model is a good fit. Also the plot of Pareto k diagnostic shows that all of the values are below 0.7. Moreover, in this case `p_loo` estimate of 19.7 also indicate the model is a good fit because the effective number of parameters (`p_loo`) is similar to the total number of parameters in the model.

Summary

Ideally, there should be a relationship between the scores of movies and the budget of the movie varying by different film production companies and genre. However, none of the model gives a very good fit to support this hypothesis. Though the results given by stan indicates a good fit of the model, Pareto k estimates sometimes can be unreliable because of the sample size of the data. Therefore, more types of model need to be fitted in order to see the relationship between movie scores and budgets as well as make predictions based on the model.

Appendix

```
# make some changes of the names
# movie.new$country <- ifelse(movie.new$country == "Hong Kong", "Hong Kong China", as.character(movie.new$country))
# movie.new$country <- ifelse(movie.new$country == "China", "Mainland China", as.character(movie.new$country))

#p1 <- ggplot(movie.new) + geom_point(aes(y=score, x=budget, group=factor(year), color = factor(year))) + theme(legend.position = "none")

#log.budget <- log(movie.new$budget)
# distribution of the budget in log scale of movies within each country
#p2 <- plotly_build(ggplot(movie.new) + aes(x= log.budget, color = country) + geom_density() + theme(legend.position = "none"))
# distribution of the score of movies within each country
#p3 <- plotly_build(ggplot(movie.new) + aes(x=score, color = country) + geom_density() + theme(legend.position = "none"))

#p4 <- ggplot(usa) + geom_point(aes(x=score, y=budget, group=factor(year), color = factor(year))) + theme(legend.position = "none")

#p5 <- ggplot(usa) + stat_summary(aes(x=genre, y=score, color = factor(genre)), fun.ymin = min, fun.ymax = max, fun.mean = mean)

#p6 <- ggplot(usa) + geom_bar(aes(x=score, color = factor(company)), fill = NA)

# set year to be the group
#p7 <- plotly_build(ggplot(usa) +
  # geom_point(aes(x=budget, y=score, color = factor(year), alpha=0.8)) +
  # geom_smooth(aes(x=budget, y=score, color = factor(year), alpha=0.2), method = "lm", se = FALSE) +
  # theme(legend.position = "none"))

# set genre to be the group
#p8 <- plotly_build(ggplot(movie.new) +
  # geom_point(aes(x=budget, y=score, color = factor(genre), alpha=0.8)) +
  # geom_smooth(aes(x=budget, y=score, color = factor(company), alpha = 0.8), method = "lm", se = FALSE))

#movie.3.coef <- lmer(formula = score ~ budget + (1|genre), data = usa)
#movie.4.coef <- lmer(formula = score ~ log.budget + (1|genre), data = usa)
#movie.5.coef <- lmer(formula = score ~ log.budget + runtime + votes + (1|genre), data = usa)
#movie.6.coef <- lmer(formula = score ~ factor(log.budget) + factor(log.runtime) + factor(log.votes) + (1|genre), data = usa)
#movie.7.coef <- lmer(formula = score ~ log.budget + (1|genre) + (1|year), data = usa)
#movie.8.coef <- lmer(formula = score ~ log.budget + log.runtime + log.votes + (1|genre) + (1|year), data = usa)
#ggplot() + aes(fitted(movie.8.coef), resid(movie.8.coef)) + geom_point()
#plot(movie.8.coef, resid(., scale = TRUE) ~ fitted(.))
#hist(usa$log.budget)
# ggplot(usa) + aes(x = log.runtime, y = log(score)) + geom_point()
#movie.10.coef <- lmer(formula = log(score) ~ log.budget + 1/log.runtime + log.votes + (1|genre) + (1|year), data = usa)
#plot(movie.10.coef, resid(., scale = TRUE) ~ fitted(.))

#movie.2 <- lmer(formula = log(score) ~ log.budget + log(gross) + log.runtime + log.votes + (1|genre) + (1|year), data = usa)
#movie.3 <- lmer(formula = log(score) ~ log.budget + log(gross) + log.runtime + log.votes + (1|rating), data = usa)

# usa %<>% mutate(cube.gross = gross^(1/3))
# movie.cube <- lmer(log(score) ~ log.budget + cube.gross + log.runtime + log.votes + (1|company), data = usa)

#score.c <- (usa$score - mean(usa$score))/sd(usa$score)
#movie.fit1 <- lmer(score.c ~ log.budget + log.gross + log.runtime + log.votes + (1|company), data = usa)

#budget.c <- (usa$budget - mean(usa$budget))/sd(usa$budget)
```

```
#gross.c <- (usa$gross - mean(usa$gross))/sd(usa$gross)
#runtime.c <- (usa$runtime - mean(usa$runtime))/sd(usa$runtime)
#votes.c <- (usa$votes - mean(usa$votes))/sd(usa$votes)

#movie.fit2 <- lmer(score.c ~ budget.c + gross.c + runtime.c + votes.c + (1|company), data = usa)
```