# Edx MovieLens Project

*June 6, 2019*

## Introduction

The goal of this project is to create a model that accurately predicts moving ratings for users. The data used in this project is from the MovieLens data set. The performance metric used was RMSE.

Key steps in the project were:

1. Data exploration and cleansing

2. Feature engineering

3. Model training

4. Model comparison and test estimation.

(Additional note: This data set large. This experiment was performed on a workstation with 16gb of ram. Workstations with less than 16gb of ram may not be able to process the data set.)

## Anlysis

### Data Load and libraries

First step was to load the data and libraries.

```r
#Load MovieLens dataset (Code copied from Edx project instructions)

#load packages

if(!require(tidyverse)) install.packages("tidyverse", repos = "http://cran.us.r-project.org")
if(!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")
```

```
## Loading required package: caret

## Loading required package: lattice

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
##     lift
```

```r
# MovieLens 10M dataset:
# https://grouplens.org/datasets/movielens/10m/
# http://files.grouplens.org/datasets/movielens/ml-10m.zip

dl <- tempfile()
download.file("http://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)

ratings <- read.table(text = gsub("::", "\t", readLines(unzip(dl, "ml-10M100K/ratings.dat"))),
                      col.names = c("userId", "movieId", "rating", "timestamp"))

movies <- str_split_fixed(readLines(unzip(dl, "ml-10M100K/movies.dat")), "\\::", 3)
colnames(movies) <- c("movieId", "title", "genres")
movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(levels(movieId))[movieId],
                                           title = as.character(title),
                                           genres = as.character(genres))

movielens <- left_join(ratings, movies, by = "movieId")

# Validation set will be 10% of MovieLens data

set.seed(1) # if using R 3.6.0: set.seed(1, sample.kind = "Rounding")
test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)
edx <- movielens[-test_index,]
temp <- movielens[test_index,]

# Make sure userId and movieId in validation set are also in edx set

validation <- temp %>%
  semi_join(edx, by = "movieId") %>%
  semi_join(edx, by = "userId")

# Add rows removed from validation set back into edx set

removed <- anti_join(temp, validation)
```

```
## Joining, by = c("userId", "movieId", "rating", "timestamp", "title", "genres")
```

```r
edx <- rbind(edx, removed)

rm(dl, ratings, movies, test_index, temp, movielens, removed)
```

```r
#load libraries

library(tidyverse)
library(splitstackshape)
```

## Data cleansing and error checking

Next we checked the data for errors.

```r
#check for NA's

sum(is.na(edx))
```

```
## [1] 0
```

```r
#view summaries of each field for 0's, outliers, and anomalies

summary(edx$userId)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       1   18124   35738   35870   53607   71567
```

```r
summary(edx$rating)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   0.500   3.000   4.000   3.512   4.000   5.000
```

```r
summary(edx$movieId)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       1     648    1834    4122    3626   65133
```

```r
data.frame(table(edx$genres))
```

```
##                                                            Var1   Freq
## 1                                            (no genres listed)      7
## 2                                                        Action  24482
## 3                                               Action|Adventure  68688
## 4               Action|Adventure|Animation|Children|Comedy        7467
## 5        Action|Adventure|Animation|Children|Comedy|Fantasy        187
## 6        Action|Adventure|Animation|Children|Comedy|IMAX           66
## 7        Action|Adventure|Animation|Children|Comedy|Sci-Fi         600
## 8               Action|Adventure|Animation|Children|Fantasy        737
## 9               Action|Adventure|Animation|Children|Sci-Fi          50
## 10                Action|Adventure|Animation|Comedy|Drama         1902
## 11                Action|Adventure|Animation|Comedy|Sci-Fi           3
## 12                Action|Adventure|Animation|Drama|Fantasy        4333
## 13        Action|Adventure|Animation|Drama|Fantasy|Sci-Fi          239
## 14                    Action|Adventure|Animation|Fantasy          868
## 15                Action|Adventure|Animation|Fantasy|Sci-Fi        451
## 16                Action|Adventure|Animation|Horror|Sci-Fi        4087
## 17                Action|Adventure|Animation|Mystery|Romance       317
## 18                    Action|Adventure|Animation|Sci-Fi          6430
## 19                Action|Adventure|Animation|Sci-Fi|Thriller      4738
## 20                              Action|Adventure|Children          824
## 21                      Action|Adventure|Children|Comedy          4785
## 22                Action|Adventure|Children|Comedy|Crime            44
## 23        Action|Adventure|Children|Comedy|Fantasy|Sci-Fi         2832
## 24                Action|Adventure|Children|Comedy|Mystery         305
```

```
## 25              Action|Adventure|Children|Crime|Mystery|Thriller       62
## 26                         Action|Adventure|Children|Fantasy     3837
## 27                                 Action|Adventure|Comedy    45118
## 28                           Action|Adventure|Comedy|Crime    27714
## 29                     Action|Adventure|Comedy|Crime|Drama      159
## 30             Action|Adventure|Comedy|Crime|Drama|Romance     6288
## 31            Action|Adventure|Comedy|Crime|Horror|Thriller      520
## 32           Action|Adventure|Comedy|Crime|Romance|Thriller       62
## 33                   Action|Adventure|Comedy|Crime|Sci-Fi       23
## 34                 Action|Adventure|Comedy|Crime|Thriller     2246
## 35                            Action|Adventure|Comedy|Drama       35
## 36 Action|Adventure|Comedy|Drama|Fantasy|Horror|Sci-Fi|Thriller      256
## 37                   Action|Adventure|Comedy|Drama|Romance      198
## 38          Action|Adventure|Comedy|Drama|Romance|Thriller      373
## 39                  Action|Adventure|Comedy|Drama|Thriller       72
## 40                       Action|Adventure|Comedy|Drama|War     5920
## 41                          Action|Adventure|Comedy|Fantasy    23562
## 42                  Action|Adventure|Comedy|Fantasy|Horror     6066
## 43                 Action|Adventure|Comedy|Fantasy|Mystery     2047
## 44                 Action|Adventure|Comedy|Fantasy|Romance    14809
## 45                  Action|Adventure|Comedy|Fantasy|Sci-Fi     1112
## 46          Action|Adventure|Comedy|Fantasy|Sci-Fi|Western     5301
## 47                         Action|Adventure|Comedy|Musical     9879
## 48                         Action|Adventure|Comedy|Romance    10835
## 49                 Action|Adventure|Comedy|Romance|Thriller    26390
## 50                          Action|Adventure|Comedy|Sci-Fi    13878
## 51                        Action|Adventure|Comedy|Thriller     2159
## 52                             Action|Adventure|Comedy|War     2204
## 53                         Action|Adventure|Comedy|Western     3485
## 54                                Action|Adventure|Crime       32
## 55                          Action|Adventure|Crime|Drama     4466
## 56           Action|Adventure|Crime|Drama|Mystery|Thriller     1539
## 57           Action|Adventure|Crime|Drama|Romance|Thriller      634
## 58                  Action|Adventure|Crime|Drama|Western      336
## 59                 Action|Adventure|Crime|Mystery|Thriller        9
## 60                  Action|Adventure|Crime|Sci-Fi|Thriller      345
## 61                        Action|Adventure|Crime|Thriller    11482
## 62                         Action|Adventure|Documentary       19
## 63                             Action|Adventure|Drama    31166
## 64                       Action|Adventure|Drama|Fantasy    11941
## 65                Action|Adventure|Drama|Fantasy|Sci-Fi       57
## 66                Action|Adventure|Drama|Fantasy|Thriller     3561
## 67                Action|Adventure|Drama|Mystery|Thriller     6657
## 68                         Action|Adventure|Drama|Romance     5021
## 69         Action|Adventure|Drama|Romance|Thriller|Western       44
## 70                Action|Adventure|Drama|Romance|War       83
## 71                         Action|Adventure|Drama|Sci-Fi    20366
## 72                Action|Adventure|Drama|Sci-Fi|Thriller     4575
## 73                        Action|Adventure|Drama|Thriller     4931
## 74                Action|Adventure|Drama|Thriller|War       86
## 75                Action|Adventure|Drama|Thriller|Western      170
## 76                            Action|Adventure|Drama|War    19652
## 77                Action|Adventure|Drama|War|Western       55
## 78                        Action|Adventure|Drama|Western     5122
```

```
## 79                              Action|Adventure|Fantasy  71176
## 80                        Action|Adventure|Fantasy|Horror   1425
## 81                Action|Adventure|Fantasy|Horror|Romance      6
## 82               Action|Adventure|Fantasy|Horror|Thriller  11460
## 83                      Action|Adventure|Fantasy|Mystery    2211
## 84                       Action|Adventure|Fantasy|Sci-Fi   16494
## 85               Action|Adventure|Fantasy|Sci-Fi|Thriller   1054
## 86                      Action|Adventure|Fantasy|Thriller    4420
## 87                           Action|Adventure|Fantasy|War     47
## 88                               Action|Adventure|Horror    4761
## 89               Action|Adventure|Horror|Sci-Fi|Thriller    7461
## 90                       Action|Adventure|Horror|Thriller    3976
## 91                              Action|Adventure|Mystery     1446
## 92                       Action|Adventure|Mystery|Sci-Fi     8304
## 93                     Action|Adventure|Mystery|Thriller    21369
## 94                              Action|Adventure|Romance    10517
## 95              Action|Adventure|Romance|Sci-Fi|Thriller     2202
## 96                      Action|Adventure|Romance|Thriller   16934
## 97                          Action|Adventure|Romance|War       6
## 98                               Action|Adventure|Sci-Fi  219938
## 99                       Action|Adventure|Sci-Fi|Thriller  105144
## 100                           Action|Adventure|Sci-Fi|War   23449
## 101                       Action|Adventure|Sci-Fi|Western    2203
## 102                             Action|Adventure|Thriller  149091
## 103                                  Action|Adventure|War    2869
## 104                          Action|Adventure|War|Western     130
## 105                             Action|Adventure|Western    10237
## 106                     Action|Animation|Children|Comedy      896
## 107             Action|Animation|Children|Comedy|Musical       30
## 108                    Action|Animation|Children|Fantasy      221
## 109                      Action|Animation|Comedy|Horror         2
## 110                        Action|Animation|Drama|Sci-Fi      883
## 111                Action|Animation|Drama|Sci-Fi|Thriller     453
## 112                             Action|Animation|Fantasy      220
## 113                      Action|Animation|Fantasy|Sci-Fi      248
## 114             Action|Animation|Film-Noir|Sci-Fi|Thriller   1103
## 115                              Action|Animation|Sci-Fi      132
## 116                                      Action|Children    3922
## 117                               Action|Children|Comedy      518
## 118                Action|Children|Comedy|Fantasy|Sci-Fi     2408
## 119                              Action|Children|Fantasy     2318
## 120                                        Action|Comedy   51289
## 121                                  Action|Comedy|Crime     9434
## 122                            Action|Comedy|Crime|Drama    27781
## 123         Action|Comedy|Crime|Drama|Horror|Thriller        396
## 124                   Action|Comedy|Crime|Drama|Thriller     8325
## 125                          Action|Comedy|Crime|Fantasy    17043
## 126                  Action|Comedy|Crime|Horror|Thriller      564
## 127                          Action|Comedy|Crime|Romance     3310
## 128                         Action|Comedy|Crime|Thriller    32549
## 129                          Action|Comedy|Crime|Western       98
## 130                                  Action|Comedy|Drama    16595
## 131                          Action|Comedy|Drama|Fantasy      548
## 132                          Action|Comedy|Drama|Romance       21
```

```
## 133                        Action|Comedy|Drama|Thriller     111
## 134                             Action|Comedy|Fantasy    2735
## 135                      Action|Comedy|Fantasy|Horror    3546
## 136               Action|Comedy|Fantasy|Horror|Sci-Fi    3757
## 137                      Action|Comedy|Fantasy|Sci-Fi   13914
## 138                             Action|Comedy|Horror    9829
## 139                      Action|Comedy|Horror|Sci-Fi     777
## 140                     Action|Comedy|Horror|Thriller    7125
## 141                            Action|Comedy|Musical    1305
## 142                     Action|Comedy|Musical|Sci-Fi    3464
## 143                    Action|Comedy|Mystery|Thriller    1529
## 144                             Action|Comedy|Romance    8069
## 145                    Action|Comedy|Romance|Thriller    1714
## 146                         Action|Comedy|Romance|War    1307
## 147                     Action|Comedy|Romance|Western      14
## 148                              Action|Comedy|Sci-Fi   30196
## 149                      Action|Comedy|Sci-Fi|Thriller      42
## 150                             Action|Comedy|Thriller    1081
## 151                     Action|Comedy|Thriller|Western      27
## 152                                 Action|Comedy|War    8313
## 153                             Action|Comedy|Western   11430
## 154                                       Action|Crime   21886
## 155                                 Action|Crime|Drama   35597
## 156            Action|Crime|Drama|Film-Noir|Mystery    1103
## 157            Action|Crime|Drama|Film-Noir|Thriller      31
## 158              Action|Crime|Drama|Horror|Thriller     930
## 159                          Action|Crime|Drama|IMAX    2353
## 160                       Action|Crime|Drama|Musical      68
## 161                       Action|Crime|Drama|Mystery     164
## 162     Action|Crime|Drama|Mystery|Sci-Fi|Thriller    5014
## 163            Action|Crime|Drama|Mystery|Thriller    9480
## 164                       Action|Crime|Drama|Romance      18
## 165             Action|Crime|Drama|Romance|Thriller      33
## 166     Action|Crime|Drama|Romance|Thriller|War    1678
## 167             Action|Crime|Drama|Romance|Western      27
## 168                        Action|Crime|Drama|Sci-Fi    6669
## 169                      Action|Crime|Drama|Thriller   65183
## 170                          Action|Crime|Drama|War       9
## 171                      Action|Crime|Drama|Western    1452
## 172                              Action|Crime|Fantasy    2693
## 173     Action|Crime|Fantasy|Mystery|Romance|Thriller    6360
## 174                    Action|Crime|Fantasy|Thriller    8807
## 175                            Action|Crime|Film-Noir      17
## 176              Action|Crime|Film-Noir|Thriller    5242
## 177                               Action|Crime|Horror    1128
## 178     Action|Crime|Horror|Mystery|Thriller     262
## 179              Action|Crime|Horror|Sci-Fi|Thriller    1601
## 180              Action|Crime|Mystery|Romance|Thriller     793
## 181              Action|Crime|Mystery|Sci-Fi|Thriller   15798
## 182                     Action|Crime|Mystery|Thriller    4511
## 183                              Action|Crime|Romance   16090
## 184                     Action|Crime|Romance|Thriller    6580
## 185                               Action|Crime|Sci-Fi   17331
## 186                      Action|Crime|Sci-Fi|Thriller   24318
```

```
## 187                              Action|Crime|Thriller 102259
## 188                          Action|Crime|Thriller|War     60
## 189                      Action|Crime|Thriller|Western   2155
## 190                              Action|Crime|Western     26
## 191                                     Action|Drama  38748
## 192                             Action|Drama|Fantasy     41
## 193                     Action|Drama|Fantasy|Romance    561
## 194                      Action|Drama|Fantasy|Sci-Fi     91
## 195         Action|Drama|Fantasy|Sci-Fi|Thriller|War     52
## 196                              Action|Drama|Horror   1432
## 197                      Action|Drama|Horror|Sci-Fi      4
## 198              Action|Drama|Horror|Sci-Fi|Thriller   4198
## 199                     Action|Drama|Horror|Thriller   1689
## 200                     Action|Drama|Musical|Romance     13
## 201                          Action|Drama|Musical|War     12
## 202                              Action|Drama|Mystery   2155
## 203            Action|Drama|Mystery|Romance|Thriller   1479
## 204                              Action|Drama|Romance  20327
## 205                     Action|Drama|Romance|Sci-Fi   3822
## 206                    Action|Drama|Romance|Thriller   1459
## 207                         Action|Drama|Romance|War  14629
## 208                     Action|Drama|Romance|Western    398
## 209                              Action|Drama|Sci-Fi   6950
## 210                     Action|Drama|Sci-Fi|Thriller  18444
## 211                             Action|Drama|Thriller  45246
## 212                         Action|Drama|Thriller|War    480
## 213                     Action|Drama|Thriller|Western   1879
## 214                                 Action|Drama|War  99183
## 215                         Action|Drama|War|Western    330
## 216                              Action|Drama|Western   9999
## 217                                   Action|Fantasy     29
## 218                             Action|Fantasy|Horror   2345
## 219         Action|Fantasy|Horror|Mystery|Sci-Fi|Thriller    409
## 220                     Action|Fantasy|Horror|Romance    130
## 221                     Action|Fantasy|Horror|Thriller   1226
## 222                     Action|Fantasy|Mystery|Thriller   1405
## 223                             Action|Fantasy|Romance      9
## 224                              Action|Fantasy|Sci-Fi    568
## 225                     Action|Fantasy|Sci-Fi|Thriller    297
## 226                             Action|Fantasy|Thriller    459
## 227                                 Action|Fantasy|War   2522
## 228                                     Action|Horror  13312
## 229                             Action|Horror|Mystery     30
## 230                     Action|Horror|Mystery|Sci-Fi     21
## 231                     Action|Horror|Mystery|Thriller    327
## 232                              Action|Horror|Sci-Fi  10364
## 233                     Action|Horror|Sci-Fi|Thriller  47724
## 234                             Action|Horror|Thriller  13262
## 235                                     Action|Mystery   6072
## 236                              Action|Mystery|Sci-Fi   5789
## 237                     Action|Mystery|Sci-Fi|Thriller    806
## 238                             Action|Mystery|Thriller   4349
## 239                                     Action|Romance  10527
## 240                     Action|Romance|Sci-Fi|Thriller  10663
```

```
## 241                                             Action|Romance|Thriller 41025
## 242                                        Action|Romance|War|Western  6383
## 243                                         Action|Romance|Western     7
## 244                                                    Action|Sci-Fi 49733
## 245                                            Action|Sci-Fi|Thriller 95280
## 246                                    Action|Sci-Fi|Thriller|Western  2260
## 247                                                Action|Sci-Fi|War  8788
## 248                                            Action|Sci-Fi|Western    58
## 249                                                  Action|Thriller 96535
## 250                                              Action|Thriller|War  2202
## 251                                                       Action|War  7567
## 252                                               Action|War|Western     2
## 253                                                   Action|Western  5397
## 254                                                        Adventure  2276
## 255                                              Adventure|Animation   163
## 256                                     Adventure|Animation|Children  1206
## 257                             Adventure|Animation|Children|Comedy 31404
## 258     Adventure|Animation|Children|Comedy|Crime|Fantasy|Mystery 10975
## 259                       Adventure|Animation|Children|Comedy|Drama  3375
## 260     Adventure|Animation|Children|Comedy|Drama|Fantasy|Mystery   355
## 261                     Adventure|Animation|Children|Comedy|Fantasy 43527
## 262              Adventure|Animation|Children|Comedy|Fantasy|Musical   836
## 263  Adventure|Animation|Children|Comedy|Fantasy|Musical|Romance   515
## 264             Adventure|Animation|Children|Comedy|Fantasy|Romance 13063
## 265               Adventure|Animation|Children|Comedy|Fantasy|Sci-Fi  5297
## 266                 Adventure|Animation|Children|Comedy|Fantasy|War   112
## 267                     Adventure|Animation|Children|Comedy|Musical 26874
## 268             Adventure|Animation|Children|Comedy|Musical|Romance  5900
## 269              Adventure|Animation|Children|Comedy|Romance|Sci-Fi  1254
## 270                      Adventure|Animation|Children|Comedy|Sci-Fi  3529
## 271             Adventure|Animation|Children|Crime|Drama|Fantasy  1354
## 272                        Adventure|Animation|Children|Drama  2165
## 273                Adventure|Animation|Children|Drama|Fantasy  1578
## 274                 Adventure|Animation|Children|Drama|Musical 19309
## 275                  Adventure|Animation|Children|Drama|Sci-Fi  3745
## 276                      Adventure|Animation|Children|Fantasy  8922
## 277                 Adventure|Animation|Children|Fantasy|IMAX   550
## 278               Adventure|Animation|Children|Fantasy|Sci-Fi   691
## 279                       Adventure|Animation|Children|Musical  6032
## 280                        Adventure|Animation|Children|Sci-Fi  1994
## 281                               Adventure|Animation|Comedy  5324
## 282                         Adventure|Animation|Comedy|Crime  4200
## 283                       Adventure|Animation|Comedy|Fantasy    19
## 284                 Adventure|Animation|Comedy|Fantasy|Musical  2141
## 285                 Adventure|Animation|Comedy|Fantasy|Romance   530
## 286                                Adventure|Animation|Drama   210
## 287                 Adventure|Animation|Drama|Fantasy|Sci-Fi   870
## 288                 Adventure|Animation|Fantasy|Romance|Sci-Fi  1333
## 289                      Adventure|Animation|Fantasy|Sci-Fi  3026
## 290                       Adventure|Animation|Musical|Sci-Fi     4
## 291                                         Adventure|Children 20260
## 292                                  Adventure|Children|Comedy  9473
## 293                        Adventure|Children|Comedy|Drama  1884
## 294                      Adventure|Children|Comedy|Fantasy 13823
```

8

```
## 295                    Adventure|Children|Comedy|Fantasy|Musical     928
## 296                    Adventure|Children|Comedy|Fantasy|Romance     931
## 297                     Adventure|Children|Comedy|Fantasy|Sci-Fi    8289
## 298                           Adventure|Children|Comedy|Musical    3070
## 299                                   Adventure|Children|Drama   10144
## 300                           Adventure|Children|Drama|Fantasy     528
## 301                    Adventure|Children|Drama|Fantasy|Sci-Fi     172
## 302                       Adventure|Children|Drama|Fantasy|War      68
## 303                           Adventure|Children|Drama|Romance     228
## 304                                 Adventure|Children|Fantasy   50123
## 305                         Adventure|Children|Fantasy|Musical   11784
## 306                          Adventure|Children|Fantasy|Sci-Fi    3485
## 307                         Adventure|Children|Fantasy|Western     108
## 308                                 Adventure|Children|Musical    2073
## 309                                 Adventure|Children|Romance    3619
## 310                                  Adventure|Children|Sci-Fi    2706
## 311                                          Adventure|Comedy   32003
## 312                                    Adventure|Comedy|Crime    9128
## 313                       Adventure|Comedy|Crime|Drama|Fantasy    3017
## 314                       Adventure|Comedy|Crime|Drama|Romance     205
## 315                            Adventure|Comedy|Crime|Mystery     236
## 316                            Adventure|Comedy|Crime|Romance    1984
## 317                           Adventure|Comedy|Crime|Thriller     251
## 318                          Adventure|Comedy|Documentary      38
## 319                                   Adventure|Comedy|Drama   15835
## 320                           Adventure|Comedy|Drama|Fantasy      28
## 321           Adventure|Comedy|Drama|Fantasy|Mystery|Sci-Fi       7
## 322                           Adventure|Comedy|Drama|Musical    8667
## 323                           Adventure|Comedy|Drama|Romance    6488
## 324                               Adventure|Comedy|Drama|War      32
## 325                           Adventure|Comedy|Drama|Western    2140
## 326                                 Adventure|Comedy|Fantasy   18464
## 327                         Adventure|Comedy|Fantasy|Musical     166
## 328                         Adventure|Comedy|Fantasy|Romance       6
## 329                          Adventure|Comedy|Fantasy|Sci-Fi      25
## 330                             Adventure|Comedy|Horror     251
## 331                      Adventure|Comedy|Horror|Romance       6
## 332                             Adventure|Comedy|Musical    2657
## 333                     Adventure|Comedy|Musical|Romance      63
## 334                             Adventure|Comedy|Romance    4391
## 335                         Adventure|Comedy|Romance|War    5223
## 336                              Adventure|Comedy|Sci-Fi   43658
## 337                            Adventure|Comedy|Thriller     320
## 338                                 Adventure|Comedy|War     207
## 339                             Adventure|Comedy|Western   19899
## 340                               Adventure|Crime|Drama     234
## 341                      Adventure|Crime|Drama|Thriller     144
## 342                  Adventure|Crime|Drama|Thriller|War    1626
## 343                      Adventure|Crime|Horror|Thriller       7
## 344                       Adventure|Crime|Sci-Fi|Thriller     711
## 345                            Adventure|Crime|Thriller     562
## 346                             Adventure|Crime|Western     784
## 347                             Adventure|Documentary     798
## 348                         Adventure|Documentary|IMAX      62
```

9

```
## 349                              Adventure|Drama  51227
## 350                      Adventure|Drama|Fantasy    198
## 351                 Adventure|Drama|Fantasy|IMAX   1385
## 352       Adventure|Drama|Fantasy|Mystery|Sci-Fi   3948
## 353              Adventure|Drama|Fantasy|Romance   1896
## 354             Adventure|Drama|Fantasy|Thriller     70
## 355      Adventure|Drama|Film-Noir|Sci-Fi|Thriller 13957
## 356      Adventure|Drama|Horror|Mystery|Thriller     64
## 357       Adventure|Drama|Horror|Sci-Fi|Thriller    217
## 358              Adventure|Drama|Horror|Thriller   1320
## 359             Adventure|Drama|Mystery|Thriller   1104
## 360                      Adventure|Drama|Romance   6226
## 361               Adventure|Drama|Romance|Sci-Fi   3849
## 362       Adventure|Drama|Romance|Sci-Fi|Thriller     44
## 363             Adventure|Drama|Romance|Thriller     33
## 364                  Adventure|Drama|Romance|War    497
## 365              Adventure|Drama|Romance|Western     16
## 366                       Adventure|Drama|Sci-Fi  13520
## 367               Adventure|Drama|Sci-Fi|Thriller   1067
## 368                    Adventure|Drama|Sci-Fi|War    262
## 369                     Adventure|Drama|Thriller   4581
## 370                 Adventure|Drama|Thriller|War     10
## 371                          Adventure|Drama|War  14137
## 372                  Adventure|Drama|War|Western     35
## 373                      Adventure|Drama|Western  23960
## 374                            Adventure|Fantasy   4264
## 375      Adventure|Fantasy|Film-Noir|Mystery|Sci-Fi      2
## 376      Adventure|Fantasy|Film-Noir|Sci-Fi|Thriller   5429
## 377                    Adventure|Fantasy|Mystery     16
## 378                    Adventure|Fantasy|Romance   4618
## 379                     Adventure|Fantasy|Sci-Fi   1002
## 380                   Adventure|Fantasy|Thriller   3080
## 381                             Adventure|Horror     34
## 382      Adventure|Horror|Mystery|Thriller     47
## 383             Adventure|Horror|Romance|Sci-Fi      5
## 384                      Adventure|Horror|Sci-Fi    577
## 385                    Adventure|Horror|Thriller     47
## 386                     Adventure|IMAX|Romance     43
## 387                            Adventure|Musical   1244
## 388                    Adventure|Musical|Romance     86
## 389                            Adventure|Mystery      2
## 390                   Adventure|Mystery|Thriller  14712
## 391                            Adventure|Romance   1944
## 392              Adventure|Romance|War|Western     90
## 393                  Adventure|Romance|Western     11
## 394                             Adventure|Sci-Fi  26893
## 395                    Adventure|Sci-Fi|Thriller    893
## 396                           Adventure|Thriller    647
## 397                               Adventure|War    631
## 398                       Adventure|War|Western    156
## 399                         Adventure|Western   4877
## 400                                    Animation    329
## 401                           Animation|Children  23723
## 402                    Animation|Children|Comedy  22914
```

```
## 403                          Animation|Children|Comedy|Crime     7167
## 404                        Animation|Children|Comedy|Fantasy    12671
## 405                        Animation|Children|Comedy|Musical     4927
## 406                Animation|Children|Comedy|Musical|Romance    11681
## 407                 Animation|Children|Comedy|Musical|Western      114
## 408                        Animation|Children|Comedy|Romance     1514
## 409                               Animation|Children|Drama     5211
## 410                 Animation|Children|Drama|Fantasy|Musical     9308
## 411                         Animation|Children|Drama|Musical     4180
## 412                              Animation|Children|Fantasy     3418
## 413                      Animation|Children|Fantasy|Musical    33089
## 414              Animation|Children|Fantasy|Musical|Romance    22050
## 415                          Animation|Children|Fantasy|War     1848
## 416                         Animation|Children|IMAX|Musical     1934
## 417                              Animation|Children|Musical    13794
## 418                      Animation|Children|Musical|Romance     6547
## 419                              Animation|Children|Western      320
## 420                                    Animation|Comedy     4206
## 421                        Animation|Comedy|Drama|Fantasy       15
## 422                 Animation|Comedy|Drama|Fantasy|Sci-Fi       51
## 423                 Animation|Comedy|Drama|Romance|Sci-Fi       36
## 424                              Animation|Comedy|Fantasy     1316
## 425                      Animation|Comedy|Fantasy|Musical      243
## 426              Animation|Comedy|Fantasy|Musical|Romance     1429
## 427                       Animation|Comedy|Fantasy|Sci-Fi      151
## 428                              Animation|Comedy|Musical     9333
## 429                               Animation|Comedy|Sci-Fi      529
## 430                             Animation|Comedy|Thriller      201
## 431                                 Animation|Comedy|War      227
## 432                             Animation|Crime|Mystery      329
## 433                               Animation|Documentary       79
## 434                           Animation|Documentary|War        4
## 435                                   Animation|Drama      392
## 436                             Animation|Drama|Fantasy     1227
## 437              Animation|Drama|Mystery|Sci-Fi|Thriller      903
## 438                             Animation|Drama|Romance      226
## 439                           Animation|Drama|Sci-Fi|War      127
## 440                                Animation|Drama|War     1039
## 441                                  Animation|Fantasy       72
## 442                           Animation|Fantasy|Horror      623
## 443                         Animation|Fantasy|Sci-Fi|War       44
## 444                         Animation|Fantasy|Thriller      228
## 445                              Animation|Horror|IMAX       10
## 446                             Animation|IMAX|Sci-Fi        7
## 447                                  Animation|Musical     3206
## 448                           Animation|Mystery|Sci-Fi      225
## 449                       Animation|Mystery|Thriller      459
## 450                                  Animation|Sci-Fi     4447
## 451                                         Children      745
## 452                                  Children|Comedy    63483
## 453                        Children|Comedy|Crime|Drama     1057
## 454                Children|Comedy|Crime|Drama|Fantasy      432
## 455                       Children|Comedy|Crime|Musical      139
## 456                             Children|Comedy|Drama     1777
```

11

```
## 457                         Children|Comedy|Drama|Fantasy  17146
## 458                 Children|Comedy|Drama|Musical|Romance     53
## 459                       Children|Comedy|Drama|Mystery     21
## 460                         Children|Comedy|Fantasy  18784
## 461                   Children|Comedy|Fantasy|Horror    257
## 462                  Children|Comedy|Fantasy|Musical  21023
## 463                  Children|Comedy|Fantasy|Romance    479
## 464                   Children|Comedy|Fantasy|Sci-Fi    438
## 465                         Children|Comedy|Musical   4648
## 466                  Children|Comedy|Musical|Romance    292
## 467                         Children|Comedy|Mystery   1150
## 468                         Children|Comedy|Romance   4881
## 469                  Children|Comedy|Romance|Sci-Fi    451
## 470                          Children|Comedy|Sci-Fi   2940
## 471                         Children|Comedy|Western   1477
## 472                            Children|Documentary    274
## 473                                  Children|Drama  13503
## 474                          Children|Drama|Fantasy    721
## 475                  Children|Drama|Fantasy|Mystery   2355
## 476                          Children|Drama|Musical    185
## 477                           Children|Drama|Sci-Fi  15851
## 478                                Children|Fantasy    688
## 479                  Children|Fantasy|Horror|Mystery    796
## 480                        Children|Fantasy|Musical   2351
## 481                Children|Fantasy|Musical|Romance     23
## 482                         Children|Fantasy|Sci-Fi     59
## 483                                 Children|Horror    113
## 484                                Children|Musical   1813
## 485                        Children|Musical|Romance     27
## 486                                Children|Mystery     48
## 487                                 Children|Sci-Fi    685
## 488                                Children|Western     52
## 489                                          Comedy 700889
## 490                                    Comedy|Crime  73286
## 491                              Comedy|Crime|Drama  59071
## 492                    Comedy|Crime|Drama|Film-Noir   2012
## 493            Comedy|Crime|Drama|Film-Noir|Musical     60
## 494              Comedy|Crime|Drama|Horror|Mystery     17
## 495    Comedy|Crime|Drama|Horror|Mystery|Thriller   3553
## 496                      Comedy|Crime|Drama|Musical   3842
## 497      Comedy|Crime|Drama|Musical|Mystery|Romance     79
## 498                      Comedy|Crime|Drama|Mystery   2626
## 499              Comedy|Crime|Drama|Mystery|Romance    140
## 500             Comedy|Crime|Drama|Mystery|Thriller   1301
## 501                      Comedy|Crime|Drama|Romance    507
## 502             Comedy|Crime|Drama|Romance|Thriller   6591
## 503                     Comedy|Crime|Drama|Thriller  24341
## 504                          Comedy|Crime|Drama|War    332
## 505                            Comedy|Crime|Fantasy     61
## 506                             Comedy|Crime|Horror   3305
## 507                            Comedy|Crime|Musical   3736
## 508                    Comedy|Crime|Musical|Mystery    388
## 509                            Comedy|Crime|Mystery   3936
## 510                    Comedy|Crime|Mystery|Romance    391
```

```
## 511                 Comedy|Crime|Mystery|Romance|Thriller    1925
## 512                       Comedy|Crime|Mystery|Thriller    9868
## 513                             Comedy|Crime|Romance    1909
## 514                   Comedy|Crime|Romance|Thriller    8066
## 515                         Comedy|Crime|Thriller   19099
## 516                          Comedy|Crime|Western      23
## 517                           Comedy|Documentary   10023
## 518                     Comedy|Documentary|Drama     276
## 519                   Comedy|Documentary|Musical     326
## 520                   Comedy|Documentary|Romance      65
## 521                                 Comedy|Drama 323637
## 522                         Comedy|Drama|Fantasy   30247
## 523           Comedy|Drama|Fantasy|Horror|Thriller    4132
## 524           Comedy|Drama|Fantasy|Musical|Romance      77
## 525                 Comedy|Drama|Fantasy|Mystery     114
## 526                 Comedy|Drama|Fantasy|Romance   42975
## 527         Comedy|Drama|Fantasy|Romance|Thriller   15488
## 528                   Comedy|Drama|Fantasy|Sci-Fi      12
## 529                       Comedy|Drama|Film-Noir     264
## 530                         Comedy|Drama|Horror     961
## 531                   Comedy|Drama|Horror|Sci-Fi       5
## 532           Comedy|Drama|Horror|Sci-Fi|Thriller      37
## 533                         Comedy|Drama|Musical    9861
## 534                 Comedy|Drama|Musical|Mystery      78
## 535                 Comedy|Drama|Musical|Romance    5647
## 536                         Comedy|Drama|Mystery    3259
## 537                 Comedy|Drama|Mystery|Romance     606
## 538                 Comedy|Drama|Mystery|Thriller      83
## 539                         Comedy|Drama|Romance 261425
## 540                 Comedy|Drama|Romance|Sci-Fi    7593
## 541                 Comedy|Drama|Romance|Thriller    3044
## 542                   Comedy|Drama|Romance|War   41762
## 543                         Comedy|Drama|Sci-Fi      90
## 544                   Comedy|Drama|Sci-Fi|War     290
## 545                       Comedy|Drama|Thriller   14269
## 546                             Comedy|Drama|War   19823
## 547                         Comedy|Drama|Western      52
## 548                               Comedy|Fantasy   19764
## 549                         Comedy|Fantasy|Horror    5004
## 550                 Comedy|Fantasy|Horror|Romance       8
## 551                 Comedy|Fantasy|Horror|Thriller    2091
## 552                       Comedy|Fantasy|Musical    1003
## 553                 Comedy|Fantasy|Musical|Romance      77
## 554                 Comedy|Fantasy|Mystery|Sci-Fi       6
## 555                       Comedy|Fantasy|Romance   28110
## 556                 Comedy|Fantasy|Romance|Sci-Fi   16511
## 557                         Comedy|Fantasy|Sci-Fi    2764
## 558                     Comedy|Film-Noir|Thriller      21
## 559                               Comedy|Horror   37394
## 560                         Comedy|Horror|Musical    4278
## 561                 Comedy|Horror|Musical|Sci-Fi    6771
## 562                         Comedy|Horror|Mystery      72
## 563                 Comedy|Horror|Mystery|Sci-Fi      48
## 564                 Comedy|Horror|Mystery|Thriller     138
```

```
## 565                          Comedy|Horror|Romance    3563
## 566                  Comedy|Horror|Romance|Thriller    1078
## 567                           Comedy|Horror|Sci-Fi    5764
## 568                         Comedy|Horror|Thriller   25700
## 569                                  Comedy|Musical   31055
## 570                          Comedy|Musical|Romance   27682
## 571                  Comedy|Musical|Romance|Western     387
## 572                          Comedy|Musical|Sci-Fi     901
## 573                          Comedy|Musical|Western     353
## 574                                  Comedy|Mystery    7898
## 575                          Comedy|Mystery|Romance     342
## 576                         Comedy|Mystery|Thriller    5240
## 577                                  Comedy|Romance  365468
## 578                          Comedy|Romance|Sci-Fi     953
## 579                         Comedy|Romance|Thriller    6549
## 580                             Comedy|Romance|War     372
## 581                         Comedy|Romance|Western      13
## 582                                   Comedy|Sci-Fi   44599
## 583                          Comedy|Sci-Fi|Thriller     389
## 584                          Comedy|Sci-Fi|Western    9467
## 585                                 Comedy|Thriller   13558
## 586                                      Comedy|War   22586
## 587                              Comedy|War|Western     342
## 588                                  Comedy|Western   11066
## 589                                           Crime    3197
## 590                              Crime|Documentary       4
## 591                        Crime|Documentary|Drama     121
## 592                          Crime|Documentary|War      11
## 593                                     Crime|Drama  137387
## 594                             Crime|Drama|Fantasy    8826
## 595              Crime|Drama|Fantasy|Film-Noir|Horror     543
## 596      Crime|Drama|Fantasy|Film-Noir|Mystery|Romance    1887
## 597             Crime|Drama|Fantasy|Romance|Thriller    4137
## 598                          Crime|Drama|Film-Noir   11249
## 599                  Crime|Drama|Film-Noir|Mystery     610
## 600         Crime|Drama|Film-Noir|Mystery|Thriller    3913
## 601                  Crime|Drama|Film-Noir|Romance       5
## 602         Crime|Drama|Film-Noir|Romance|Thriller     162
## 603                 Crime|Drama|Film-Noir|Thriller    7861
## 604                             Crime|Drama|Horror     465
## 605                     Crime|Drama|Horror|Mystery      17
## 606             Crime|Drama|Horror|Mystery|Thriller    6362
## 607                      Crime|Drama|Horror|Sci-Fi       2
## 608                     Crime|Drama|Horror|Thriller     212
## 609                            Crime|Drama|Musical      42
## 610                            Crime|Drama|Mystery    8799
## 611                    Crime|Drama|Mystery|Romance      46
## 612            Crime|Drama|Mystery|Romance|Thriller    1646
## 613                   Crime|Drama|Mystery|Thriller   28570
## 614               Crime|Drama|Mystery|Thriller|War     180
## 615                            Crime|Drama|Romance   11135
## 616                   Crime|Drama|Romance|Thriller   19470
## 617                             Crime|Drama|Sci-Fi    1244
## 618                    Crime|Drama|Sci-Fi|Thriller   10730
```

```
## 619                             Crime|Drama|Thriller 106101
## 620                             Crime|Drama|War    206
## 621                          Crime|Drama|Western   1211
## 622                             Crime|Film-Noir   4241
## 623                       Crime|Film-Noir|Mystery   4029
## 624              Crime|Film-Noir|Mystery|Thriller  24961
## 625                     Crime|Film-Noir|Romance     30
## 626                     Crime|Film-Noir|Thriller   4844
## 627                             Crime|Horror    143
## 628                        Crime|Horror|Mystery    302
## 629                   Crime|Horror|Mystery|Sci-Fi     29
## 630          Crime|Horror|Mystery|Sci-Fi|Thriller     32
## 631              Crime|Horror|Mystery|Thriller  27240
## 632                         Crime|Horror|Sci-Fi     15
## 633                Crime|Horror|Sci-Fi|Thriller     38
## 634                     Crime|Horror|Thriller  33757
## 635                             Crime|Musical    295
## 636                             Crime|Mystery   1002
## 637              Crime|Mystery|Romance|Thriller   2394
## 638              Crime|Mystery|Sci-Fi|Thriller     61
## 639                     Crime|Mystery|Thriller  26892
## 640                     Crime|Romance|Thriller    180
## 641                     Crime|Sci-Fi|Thriller   2351
## 642                             Crime|Thriller  15739
## 643                         Crime|Thriller|War   4595
## 644                                Documentary  70041
## 645                          Documentary|Drama   1859
## 646                   Documentary|Drama|Musical     95
## 647           Documentary|Drama|Romance|War     55
## 648                   Documentary|Drama|War    105
## 649                       Documentary|Fantasy     27
## 650                        Documentary|Horror    619
## 651                          Documentary|IMAX   1714
## 652                  Documentary|IMAX|Musical     43
## 653                       Documentary|Musical   5100
## 654                       Documentary|Mystery     29
## 655                       Documentary|Romance      2
## 656                      Documentary|Thriller     71
## 657                          Documentary|War   1206
## 658                                      Drama 733296
## 659                              Drama|Fantasy  16554
## 660                       Drama|Fantasy|Horror    694
## 661               Drama|Fantasy|Horror|Mystery     13
## 662      Drama|Fantasy|Horror|Mystery|Thriller   3574
## 663               Drama|Fantasy|Horror|Romance     55
## 664                Drama|Fantasy|Horror|Sci-Fi     33
## 665       Drama|Fantasy|Horror|Sci-Fi|Thriller     59
## 666              Drama|Fantasy|Horror|Thriller  10183
## 667          Drama|Fantasy|Horror|Thriller|War    575
## 668                      Drama|Fantasy|Musical   1351
## 669              Drama|Fantasy|Musical|Romance    156
## 670                      Drama|Fantasy|Mystery   1310
## 671              Drama|Fantasy|Mystery|Romance   1547
## 672     Drama|Fantasy|Mystery|Romance|Thriller   1597
```

```
## 673                                  Drama|Fantasy|Mystery|Sci-Fi    3351
## 674                                Drama|Fantasy|Mystery|Thriller      75
## 675                                         Drama|Fantasy|Romance   17021
## 676                                Drama|Fantasy|Romance|Thriller      54
## 677                                          Drama|Fantasy|Sci-Fi     918
## 678                                         Drama|Fantasy|Thriller    4011
## 679                                              Drama|Fantasy|War     222
## 680                                                Drama|Film-Noir    3010
## 681              Drama|Film-Noir|Horror|Mystery|Thriller             195
## 682                                        Drama|Film-Noir|Mystery     385
## 683                                Drama|Film-Noir|Mystery|Romance     618
## 684                                Drama|Film-Noir|Mystery|Thriller    467
## 685                                        Drama|Film-Noir|Romance    2989
## 686                                        Drama|Film-Noir|Thriller   1490
## 687                                                   Drama|Horror   28860
## 688                                            Drama|Horror|Musical     54
## 689                                   Drama|Horror|Musical|Thriller    825
## 690                                            Drama|Horror|Mystery    102
## 691                        Drama|Horror|Mystery|Sci-Fi|Thriller        2
## 692                                   Drama|Horror|Mystery|Thriller   7223
## 693                                            Drama|Horror|Romance    101
## 694                                   Drama|Horror|Romance|Thriller   3438
## 695                                             Drama|Horror|Sci-Fi   6115
## 696                                    Drama|Horror|Sci-Fi|Thriller   4378
## 697                               Drama|Horror|Sci-Fi|Thriller|War    813
## 698                                           Drama|Horror|Thriller  24171
## 699                                               Drama|Horror|War      65
## 700                                                  Drama|Musical   25646
## 701                                          Drama|Musical|Romance   22344
## 702                                      Drama|Musical|Romance|War     313
## 703                                         Drama|Musical|Thriller       4
## 704                                              Drama|Musical|War    2944
## 705                                                  Drama|Mystery   33390
## 706                                          Drama|Mystery|Romance   14441
## 707                        Drama|Mystery|Romance|Sci-Fi|Thriller      537
## 708                                 Drama|Mystery|Romance|Thriller   17099
## 709                                      Drama|Mystery|Romance|War     854
## 710                                           Drama|Mystery|Sci-Fi    1657
## 711                                  Drama|Mystery|Sci-Fi|Thriller    7711
## 712                                          Drama|Mystery|Thriller  61069
## 713                                              Drama|Mystery|War      28
## 714                                          Drama|Mystery|Western    1575
## 715                                                  Drama|Romance  259355
## 716                                           Drama|Romance|Sci-Fi    3555
## 717                                   Drama|Romance|Sci-Fi|Thriller   3549
## 718                                          Drama|Romance|Thriller  17115
## 719                                  Drama|Romance|Thriller|War       224
## 720                                              Drama|Romance|War   35471
## 721                                  Drama|Romance|War|Western        9011
## 722                                          Drama|Romance|Western    1210
## 723                                                   Drama|Sci-Fi   22926
## 724                                           Drama|Sci-Fi|Thriller  30623
## 725                                               Drama|Sci-Fi|War     369
## 726                                                  Drama|Thriller 145373
```

```
## 727                          Drama|Thriller|War  16959
## 728                      Drama|Thriller|Western    224
## 729                                     Drama|War 111029
## 730                             Drama|War|Western    137
## 731                                 Drama|Western  12976
## 732                                       Fantasy     86
## 733                                Fantasy|Horror   5641
## 734                        Fantasy|Horror|Mystery   1392
## 735                Fantasy|Horror|Mystery|Romance   6431
## 736               Fantasy|Horror|Mystery|Thriller   1815
## 737               Fantasy|Horror|Romance|Thriller   4770
## 738                         Fantasy|Horror|Sci-Fi      6
## 739                Fantasy|Horror|Sci-Fi|Thriller    886
## 740                       Fantasy|Horror|Thriller   1427
## 741                               Fantasy|Musical     56
## 742                       Fantasy|Musical|Romance    464
## 743                   Fantasy|Mystery|Sci-Fi|War      2
## 744                      Fantasy|Mystery|Thriller    198
## 745                       Fantasy|Mystery|Western      5
## 746                               Fantasy|Romance     76
## 747                                Fantasy|Sci-Fi    517
## 748                       Fantasy|Sci-Fi|Thriller   2221
## 749                               Fantasy|Western     87
## 750                                     Film-Noir   1575
## 751                             Film-Noir|Horror     24
## 752                    Film-Noir|Horror|Thriller     14
## 753                             Film-Noir|Mystery   5988
## 754                    Film-Noir|Mystery|Thriller   4011
## 755                    Film-Noir|Romance|Thriller   2453
## 756                            Film-Noir|Thriller   1746
## 757                                        Horror  68738
## 758           Horror|Musical|Mystery|Thriller    686
## 759                  Horror|Musical|Thriller    288
## 760                                 Horror|Mystery   9060
## 761                          Horror|Mystery|Sci-Fi     49
## 762                 Horror|Mystery|Sci-Fi|Thriller   3565
## 763                         Horror|Mystery|Thriller  14789
## 764                                 Horror|Romance   1002
## 765                         Horror|Romance|Thriller      4
## 766                                  Horror|Sci-Fi  31281
## 767                         Horror|Sci-Fi|Thriller  27194
## 768                                Horror|Thriller  75000
## 769                            Horror|War|Western      3
## 770                                Horror|Western   1415
## 771                                          IMAX     14
## 772                                       Musical   3851
## 773                               Musical|Romance   9150
## 774                           Musical|Romance|War   1151
## 775                       Musical|Romance|Western    291
## 776                                 Musical|Sci-Fi     50
## 777                                   Musical|War   1502
## 778                               Musical|Western    269
## 779                                       Mystery    246
## 780                       Mystery|Romance|Thriller   3418
```

17

```
## 781                                          Mystery|Sci-Fi     939
## 782                                   Mystery|Sci-Fi|Thriller     850
## 783                                         Mystery|Thriller   25048
## 784                                          Mystery|Western      35
## 785                                                  Romance    8410
## 786                                           Romance|Sci-Fi     645
## 787                                          Romance|Thriller    1967
## 788                                              Romance|War     862
## 789                                          Romance|Western     849
## 790                                                   Sci-Fi   10125
## 791                                           Sci-Fi|Thriller   40129
## 792                                                 Thriller   94662
## 793                                             Thriller|War      41
## 794                                         Thriller|Western      15
## 795                                                      War    2300
## 796                                              War|Western      14
## 797                                                  Western   15300
```

After reviewing the data there seem to be no data integrity or outline problems. There were 7 movies that did not have a genres listed. Genres also contains multiple entries. These will be broken out into individual dummy variables. Looking at the data, additional features could be useful for the mode. We will engineer features for average movie rating, individuals average movie rating, and individuals average movie ratings for the most popular genres.

## Feature engineering

Next we engineered additional features for the data set.

```
#create dummy variables for genres

genres <- cSplit(edx[6], 'genres', sep="|", type.convert=FALSE)
genres <- unique(genres[,1])

for (i in 1:19) {
  gen <- genres[i]
  index <- grep(gen, edx$genres)
  edx[index,6+i] <-1
  edx[-index,6+i] <-0
}

for (i in 1:19) {
  gen <- genres[i]
  index <- grep(gen, validation$genres)
  validation[index,6+i] <-1
  validation[-index,6+i] <-0
}

#create average movie ratings and user ratings

avgmovierating <- edx %>% group_by(movieId) %>% summarize(avgmovierating = mean(rating))
avguserrating <- edx %>% group_by(userId) %>% summarize(avguserrating = mean(rating))

edx <- merge(edx, avgmovierating, 'movieId', all.x = TRUE)
```

```
edx <- merge(edx, avguserrating, 'userId', all.x = TRUE)

validation <- merge(validation, avgmovierating, 'movieId', all.x = TRUE)
validation <- merge(validation, avguserrating, 'userId', all.x = TRUE)

#create individual average movie ratings for top 5 genreas

avguserrating_comedy <- edx[edx$V7 == 1,] %>% group_by(userId) %>% summarize(avguserrating_comedy = mean
edx <- merge(edx, avguserrating_comedy, 'userId', all.x = TRUE)
validation <- merge(validation, avguserrating_comedy, 'userId', all.x = TRUE)

avguserrating_action <- edx[edx$V8 == 1,] %>% group_by(userId) %>% summarize(avguserrating_action = mean
edx <- merge(edx, avguserrating_action, 'userId', all.x = TRUE)
validation <- merge(validation, avguserrating_action, 'userId', all.x = TRUE)

avguserrating_adventure <- edx[edx$V10 == 1,] %>% group_by(userId) %>% summarize(avguserrating_adventure
edx <- merge(edx, avguserrating_adventure, 'userId', all.x = TRUE)
validation <- merge(validation, avguserrating_adventure, 'userId', all.x = TRUE)

avguserrating_drama <- edx[edx$V12 == 1,] %>% group_by(userId) %>% summarize(avguserrating_drama = mean
edx <- merge(edx, avguserrating_drama, 'userId', all.x = TRUE)
validation <- merge(validation, avguserrating_drama, 'userId', all.x = TRUE)

avguserrating_crime <- edx[edx$V13 == 1,] %>% group_by(userId) %>% summarize(avguserrating_crime = mean
edx <- merge(edx, avguserrating_crime, 'userId', all.x = TRUE)
validation <- merge(validation, avguserrating_crime, 'userId', all.x = TRUE)

#view summary of engineered features
summary(edx[,c(3,26:32)])
```

```
##      rating        avgmovierating   avguserrating    avguserrating_comedy
##  Min.   :0.500   Min.   :0.500    Min.   :0.500    Min.   :0.500
##  1st Qu.:3.000   1st Qu.:3.218    1st Qu.:3.252    1st Qu.:3.153
##  Median :4.000   Median :3.591    Median :3.529    Median :3.455
##  Mean   :3.512   Mean   :3.512    Mean   :3.512    Mean   :3.436
##  3rd Qu.:4.000   3rd Qu.:3.876    3rd Qu.:3.800    3rd Qu.:3.750
##  Max.   :5.000   Max.   :5.000    Max.   :5.000    Max.   :5.000
##                                                    NA's   :301
##  avguserrating_action avguserrating_adventure avguserrating_drama
##  Min.   :0.500        Min.   :0.500           Min.   :0.500
##  1st Qu.:3.087        1st Qu.:3.160           1st Qu.:3.405
##  Median :3.412        Median :3.486           Median :3.675
##  Mean   :3.394        Mean   :3.463           Mean   :3.658
##  3rd Qu.:3.721        3rd Qu.:3.787           3rd Qu.:3.940
##  Max.   :5.000        Max.   :5.000           Max.   :5.000
##  NA's   :7211         NA's   :9054            NA's   :288
##  avguserrating_crime
##  Min.   :0.50
##  1st Qu.:3.35
##  Median :3.66
##  Mean   :3.64
##  3rd Qu.:3.95
##  Max.   :5.00
```

```
##  NA's    :32358
```

```r
#remove extra tables and clear memory
rm(avgmovierating,avguserrating,avguserrating_action,avguserrating_adventure,avguserrating_comedy,avguse
gc()
```

```
##             used    (Mb) gc trigger   (Mb)    max used    (Mb)
## Ncells    2035195   108.7   16153870   862.8    43871432  2343.0
## Vcells 330220313  2519.4 1017595051 7763.7 1267029607  9666.7
```

The engineered features contain NA's that need to be removed because some users have not rated those
genres. The NAs will be replaced with the average genre rating for users who haven't rated that genre.

```r
#replace NA's with means
edx$avguserrating_comedy[is.na(edx$avguserrating_comedy)] <- mean(na.omit(edx$avguserrating_comedy))
validation$avguserrating_comedy[is.na(validation$avguserrating_comedy)] <- mean(na.omit(edx$avguserratin

edx$avguserrating_action[is.na(edx$avguserrating_action)] <- mean(na.omit(edx$avguserrating_action))
validation$avguserrating_action[is.na(validation$avguserrating_action)] <- mean(na.omit(edx$avguserratin

edx$avguserrating_adventure[is.na(edx$avguserrating_adventure)] <- mean(na.omit(edx$avguserrating_advent
validation$avguserrating_adventure[is.na(validation$avguserrating_adventure)] <- mean(na.omit(edx$avguse

edx$avguserrating_drama[is.na(edx$avguserrating_drama)] <- mean(na.omit(edx$avguserrating_drama))
validation$avguserrating_drama[is.na(validation$avguserrating_drama)] <- mean(na.omit(edx$avguserrating_

edx$avguserrating_crime[is.na(edx$avguserrating_crime)] <- mean(na.omit(edx$avguserrating_crime))
validation$avguserrating_crime[is.na(validation$avguserrating_crime)] <- mean(na.omit(edx$avguserrating_

summary(edx[,c(3,26:32)])
```

```
##      rating        avgmovierating    avguserrating    avguserrating_comedy
##  Min.   :0.500   Min.   :0.500   Min.   :0.500   Min.   :0.500
##  1st Qu.:3.000   1st Qu.:3.218   1st Qu.:3.252   1st Qu.:3.153
##  Median :4.000   Median :3.591   Median :3.529   Median :3.455
##  Mean   :3.512   Mean   :3.512   Mean   :3.512   Mean   :3.436
##  3rd Qu.:4.000   3rd Qu.:3.876   3rd Qu.:3.800   3rd Qu.:3.750
##  Max.   :5.000   Max.   :5.000   Max.   :5.000   Max.   :5.000
##  avguserrating_action avguserrating_adventure avguserrating_drama
##  Min.   :0.500        Min.   :0.500           Min.   :0.500
##  1st Qu.:3.088        1st Qu.:3.161           1st Qu.:3.405
##  Median :3.411        Median :3.485           Median :3.675
##  Mean   :3.394        Mean   :3.463           Mean   :3.658
##  3rd Qu.:3.721        3rd Qu.:3.787           3rd Qu.:3.940
##  Max.   :5.000        Max.   :5.000           Max.   :5.000
##  avguserrating_crime
##  Min.   :0.500
##  1st Qu.:3.356
##  Median :3.655
##  Mean   :3.639
##  3rd Qu.:3.954
##  Max.   :5.000
```

NA's have been replaced and features are ready for modeling.

## Modeling

We will attempt to fit a simple linear model.

```
#fit LM model
fit.lm <- lm(rating~., data = edx[,c(3,7:32)])
summary(fit.lm)
```

```
##
## Call:
## lm(formula = rating ~ ., data = edx[, c(3, 7:32)])
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -4.7936 -0.4978  0.0661  0.5823  4.9380
##
## Coefficients:
##                          Estimate Std. Error   t value Pr(>|t|)
## (Intercept)            -2.6488830  0.0033139  -799.335  < 2e-16 ***
## V7                     -0.0051231  0.0007330    -6.989 2.77e-12 ***
## V8                     -0.0093468  0.0008277   -11.292  < 2e-16 ***
## V9                     -0.0203407  0.0015128   -13.446  < 2e-16 ***
## V10                    -0.0095331  0.0008525   -11.182  < 2e-16 ***
## V11                     0.0159105  0.0017705     8.986  < 2e-16 ***
## V12                     0.0088877  0.0007387    12.031  < 2e-16 ***
## V13                     0.0042884  0.0009112     4.706 2.52e-06 ***
## V14                    -0.0169109  0.0009106   -18.571  < 2e-16 ***
## V15                     0.0151429  0.0011847    12.782  < 2e-16 ***
## V16                    -0.0030801  0.0008025    -3.838 0.000124 ***
## V17                    -0.0028225  0.0026390    -1.070 0.284834
## V18                     0.0035303  0.0012742     2.771 0.005597 **
## V19                    -0.0083860  0.0020473    -4.096 4.20e-05 ***
## V20                     0.0386823  0.0029774    12.992  < 2e-16 ***
## V21                    -0.0107734  0.0007855   -13.716  < 2e-16 ***
## V22                     0.0015821  0.0010449     1.514 0.129994
## V23                    -0.0102416  0.0014969    -6.842 7.81e-12 ***
## V24                    -0.0096776  0.0013195    -7.334 2.23e-13 ***
## V25                    -0.0148780  0.0096576    -1.541 0.123425
## avgmovierating          0.8955754  0.0006646  1347.493  < 2e-16 ***
## avguserrating           0.3277623  0.0037542    87.307  < 2e-16 ***
## avguserrating_comedy    0.0874988  0.0018060    48.448  < 2e-16 ***
## avguserrating_action    0.1635695  0.0014466   113.068  < 2e-16 ***
## avguserrating_adventure 0.0899487  0.0014255    63.099  < 2e-16 ***
## avguserrating_drama     0.2246446  0.0019342   116.145  < 2e-16 ***
## avguserrating_crime    -0.0320895  0.0010487   -30.599  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8699 on 9000028 degrees of freedom
## Multiple R-squared:  0.327,  Adjusted R-squared:  0.3269
## F-statistic: 1.682e+05 on 26 and 9000028 DF,  p-value: < 2.2e-16
```

```
#train RMSE
predict.lm <- predict(fit.lm, newdata = edx[,c(3,7:32)])
sqrt(mean((edx$rating - predict.lm)^2))
```

## [1] 0.8698914

The majority of the features in the model are significant and the F-statistic is significant. The RMSE on the training set was .8698, within the target range of $<= 0.87750$. We will remove non-significant(p>.05) variables and refit the model.

```
#remove previous fit to save memory
rm(fit.lm)
rm(predict.lm)

#fit LM model and remove non-signicant variables
fit.lm <- lm(rating~., data = edx[,c(3,7:16,18:21,23:24,26:32)])
summary(fit.lm)
```

```
##
## Call:
## lm(formula = rating ~ ., data = edx[, c(3, 7:16, 18:21, 23:24,
##      26:32)])
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -4.7938 -0.4979  0.0661  0.5823  4.9377
##
## Coefficients:
##                        Estimate Std. Error   t value Pr(>|t|)
## (Intercept)          -2.6488386  0.0033116  -799.872  < 2e-16 ***
## V7                   -0.0050051  0.0007303    -6.853 7.22e-12 ***
## V8                   -0.0091974  0.0008214   -11.197  < 2e-16 ***
## V9                   -0.0197478  0.0014657   -13.473  < 2e-16 ***
## V10                  -0.0093832  0.0008452   -11.102  < 2e-16 ***
## V11                   0.0158718  0.0017700     8.967  < 2e-16 ***
## V12                   0.0089211  0.0007364    12.114  < 2e-16 ***
## V13                   0.0041353  0.0009053     4.568 4.93e-06 ***
## V14                  -0.0169873  0.0009091   -18.685  < 2e-16 ***
## V15                   0.0153728  0.0011783    13.047  < 2e-16 ***
## V16                  -0.0031640  0.0008001    -3.955 7.67e-05 ***
## V18                   0.0034682  0.0012667     2.738  0.00618 **
## V19                  -0.0084680  0.0020456    -4.140 3.48e-05 ***
## V20                   0.0384209  0.0029702    12.935  < 2e-16 ***
## V21                  -0.0106670  0.0007833   -13.618  < 2e-16 ***
## V23                  -0.0102861  0.0014965    -6.873 6.27e-12 ***
## V24                  -0.0097646  0.0013171    -7.414 1.23e-13 ***
## avgmovierating        0.8955571  0.0006618  1353.185  < 2e-16 ***
## avguserrating         0.3276878  0.0037537    87.298  < 2e-16 ***
## avguserrating_comedy  0.0875153  0.0018060    48.458  < 2e-16 ***
## avguserrating_action  0.1635990  0.0014464   113.110  < 2e-16 ***
## avguserrating_adventure  0.0899508  0.0014255   63.100  < 2e-16 ***
## avguserrating_drama   0.2246710  0.0019340   116.168  < 2e-16 ***
```

```
## avguserrating_crime     -0.0320939  0.0010486  -30.606  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8699 on 9000031 degrees of freedom
## Multiple R-squared:  0.327,   Adjusted R-squared:  0.3269
## F-statistic: 1.901e+05 on 23 and 9000031 DF,  p-value: < 2.2e-16
```

```r
#train RMSE
predict.lm <- predict(fit.lm, newdata = edx[,c(3,7:32)])
sqrt(mean((edx$rating - predict.lm)^2))
```

```
## [1] 0.8698917
```

The model and remaining variables are significant. The RMSE on the training data was .8698, within the target range of $<= 0.87750$. This model should be suitable for our final model.

## Results

Lastly we tested the linear model against the validation set to provide a final estimate of performance.

```r
#remove previous predict to save memory
rm(predict.lm)

#final validation RMSE
predict.lm <- predict(fit.lm, newdata = validation[,c(3,7:32)])
sqrt(mean((validation$rating - predict.lm)^2))
```

```
## [1] 0.8772375
```

The final model performance on the validation set was .8772, within the target range of $<= 0.87750$.

## Conclusion

The final linear model, with a RMSE of .8772, was within the target range of $<= 0.87750$. Additional features were engineered that ended up being significant to the model. These included variables for genres, average movie ratings, user's average ratings, and user's average ratings for genres.