



[Course](#) > [Capstone Project: All Learners](#) > [Project Overview: MovieLens](#) >
Create Train and Validation Sets

Audit Access Expires Aug 13, 2019

You lose all access to this course, including your progress, on Aug 13, 2019. Upgrade by Aug 6, 2019 to get unlimited access to the course as long as it exists on the site. [Upgrade now](#)

Create Train and Validation Sets

You will use the following code to generate your datasets. Develop your algorithm using the **edx** set. For a final test of your algorithm, predict movie ratings in the **validation** set as if they were unknown. RMSE will be used to evaluate how close your predictions are to the true values in the **validation** set.

Also remember that by accessing this site, you are agreeing to the terms of the [edX Honor Code](#).

Create test and validation sets

```
#####  
# Create edx set, validation set  
#####  
  
# Note: this process could take a couple of minutes  
  
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")  
if(!require(data.table)) install.packages("data.table", repos =
```

```

"http://cran.us.r-project.org")

# 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 <- fread(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, sample.kind="Rounding")
# if using R 3.5 or earlier, use `set.seed(1)` instead
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)
edx <- rbind(edx, removed)

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

[Learn About Verified Certificates](#)

© All Rights Reserved