

MovieLens Project

Adele Taylor

09/11/2020

Introduction

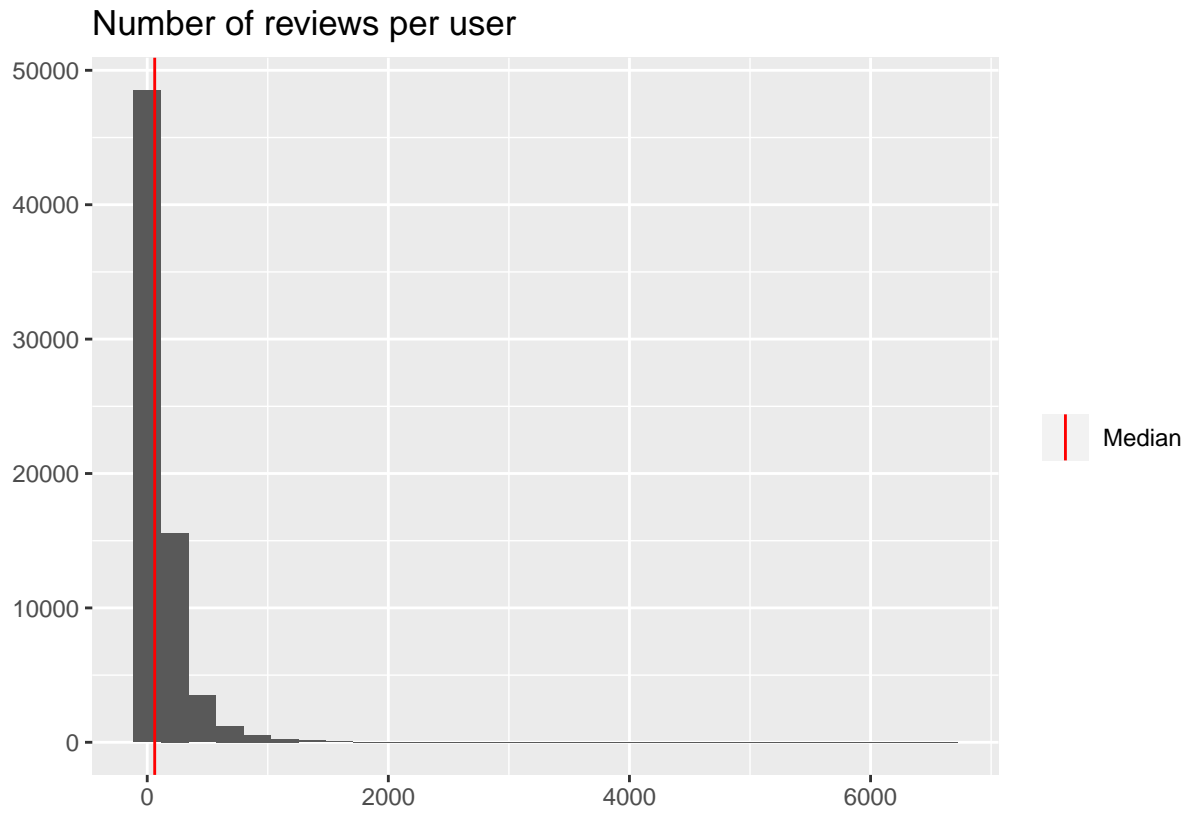
MovieLens (movielens.org) is a non-commercial film rating and recommendation website run by the GroupLens research group at the University of Minnesota. The aim of this project is to use the 10M dataset to design and validate a movie recommendation system.

The dataset was initially divided into training and validation sets using code provided by the edX course “HarvardX PH125.9x Data Science: Capstone”. All investigations and analysis was performed on the training set (“edx”), which comprised over 9 million observations consisting of the variables “userId”, “movieId”, “timestamp”, “title” and “genres”. All further reference to “the dataset” refer to this subset, unless explicitly stated otherwise.

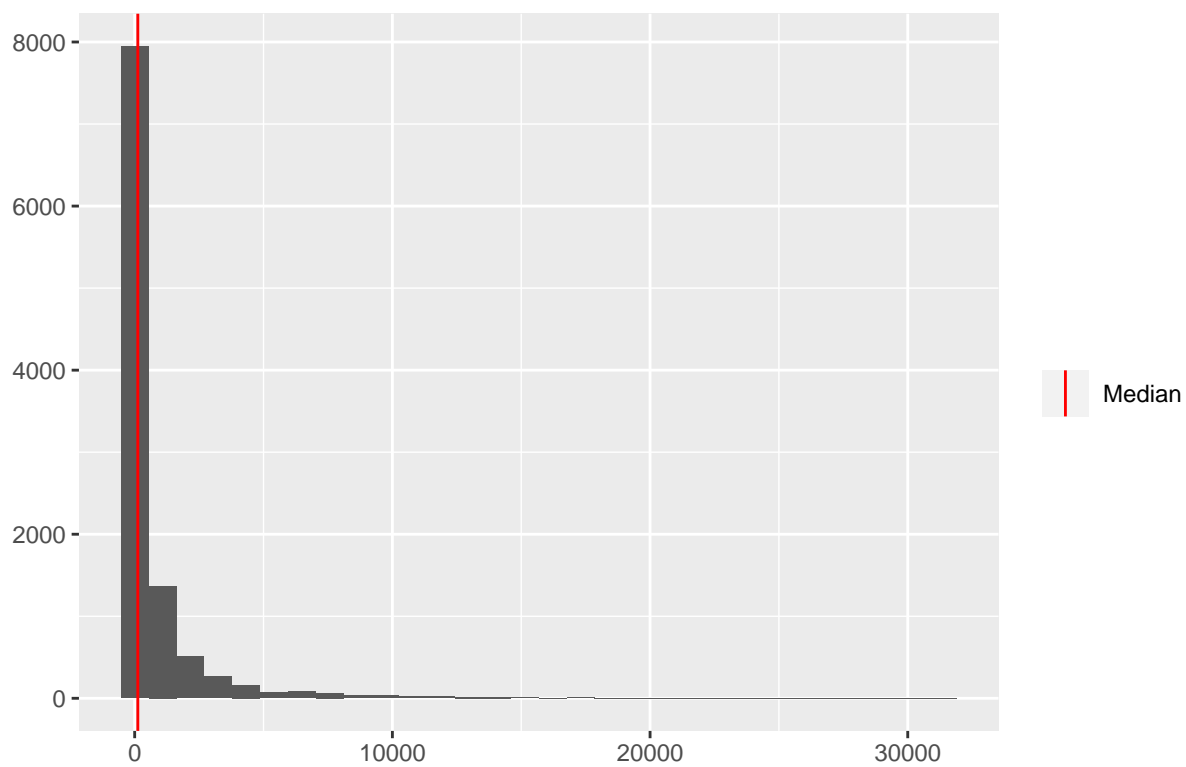
We will briefly examine the dataset to gain insights which may be useful, before trying out a range of models on a smaller subset of the dataset. We will construct a regularised linear model and use this and other potentially well-performing models to create an ensemble which we will use to predict ratings for the validation dataset.

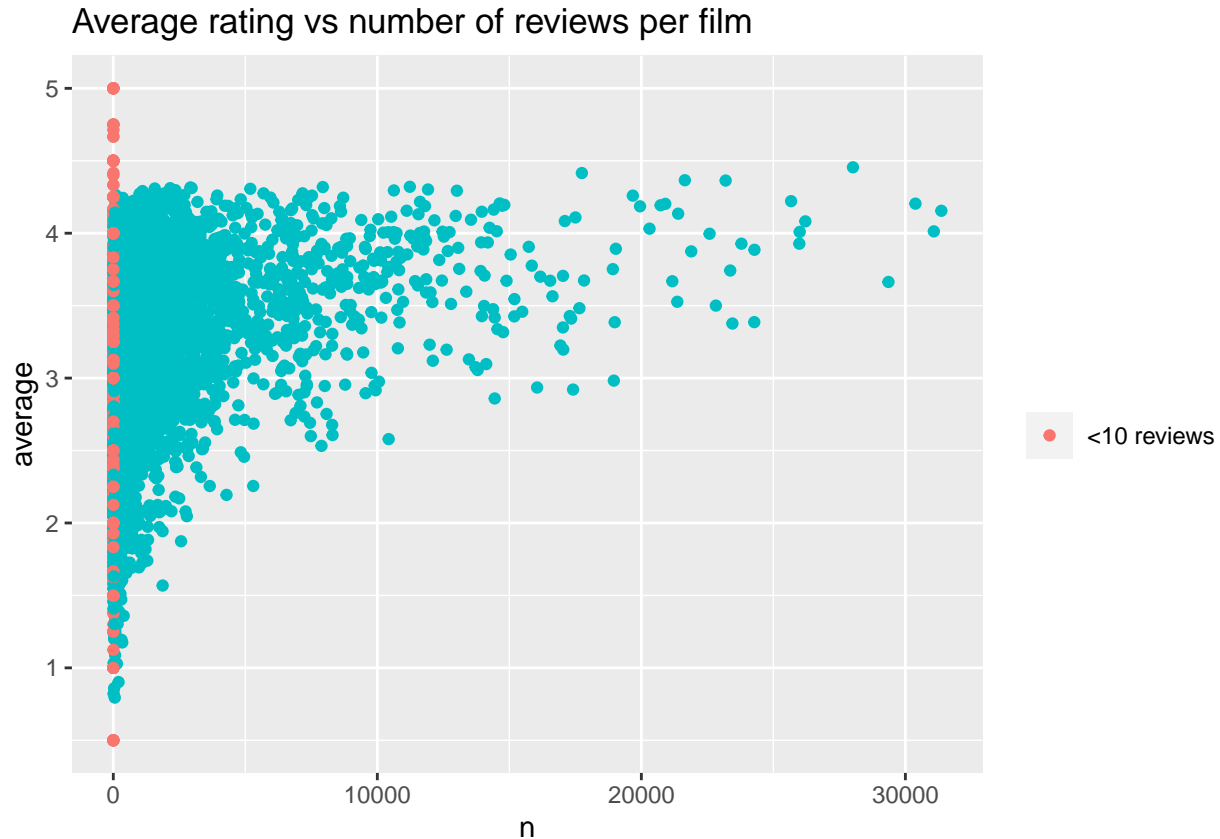
Analysis

The dataset contains ratings for movies by ‘`r length(unique(edx$userId)) %>% print()`’ users. To get an idea of how these ratings are distributed, several graphs were created as below.



Number of reviews per film





It's clear that there is a large variation in the number of reviews a film gets and in the number of ratings a given user makes. We can also see that there is a slight positive correlation ($r = \text{cor}(\text{tempn}, \text{tempaverage})$) between the number of reviews a film gets and its average (mean) rating.

To investigate further which predictors we will train and then test several simple models using randomly generated subsets. The sheer size of the dataset will prevent us from using all of it to train models, but we will use a much larger subset when constructing the final ensemble.

We will use the caret package and train a general linear model ("glm"), a k-nearest neighbours model ("knn") as well as models using random partitioning ("rpart") and the popular Random Forest algorithm ("rf"). First we will only use the `userId`, then only the `movieId`, and finally both.

```
## Warning in set.seed(20, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used
```

```
## Warning in set.seed(10, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used
```

```
## Warning in set.seed(3, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used
```

```
## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range
```

```
## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range
```

[illegible]


```

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in set.seed(1, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used

## Warning in set.seed(8, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used

## Warning in set.seed(4, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:

```


[illegible]

```

## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in set.seed(19, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used

## Warning in set.seed(80, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used

## Warning in set.seed(1, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

```


[illegible]

[illegible]


```
## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range

## Warning in randomForest.default(x, y, mtry = param$mtry, ...): invalid mtry:
## reset to within valid range
```

	glm	knn	rf	rpart
userId only	1.16	1.16	1.16	1.16
movieId only	1.19	1.19	1.19	1.19
both userId and movieId	1.21	1.21	1.21	1.21

Results

Conclusion

Due to time and computer power constraints, we could only perform a subset of the analysis and training of potential models that we would have wished. In particular, we did not fully investigate whether adding some (or all) of the genre tags as predictors would have increased accuracy or provided any additional information beyond what we can glean from user and movie effects. To do this, we would first have to introduce extra variables/columns (one for each genre tag) and then investigate fully how these may interact with the other predictors. I suspect by themselves they would not provide any additional information beyond the movie effect, but the interaction with the user may yield better results.