

# DATA SCIENCE WITH PYTHON: EXPLORATORY ANALYSIS OF MOVIE RECOMMENDATION

BING MIU , RONG ZHANG , MEILAN CHEN , AND JINCHAO FENG

## **Abstract.**

In this project, we explore the Netflix dataset with two linear models (with/without regularization) and two non-parametric models (Random Forest and Matrix Factorization) to predict the movie rating and build accurate movie recommendation system. Comparing the results based on Mean Square Error (MSE), we find that Matrix Factorization has the best performance.

## **1. Introduction.**

Movie recommendation systems offer users with personalized suggestions for movie choice based on their preferences. It is extensively used by Netflix to suggest movies to the users and to provide users with information to help them decide which movie to watch [2]. In 2006, Netflix released over 100 million customer generated movie ratings on nearly 18,000 movies as part of the Netflix Prize competition [1]. Contestants were rated based on their algorithm's error in the predictions of users' ratings of movies [3]. We analyze the dataset in Section 2.1, and build the two linear models in Section 2.2; then we attempt Random Forest and Matrix Factorization methods on Section 2.3; finally, we include the compared results in Section 3 and summary on Section 4.

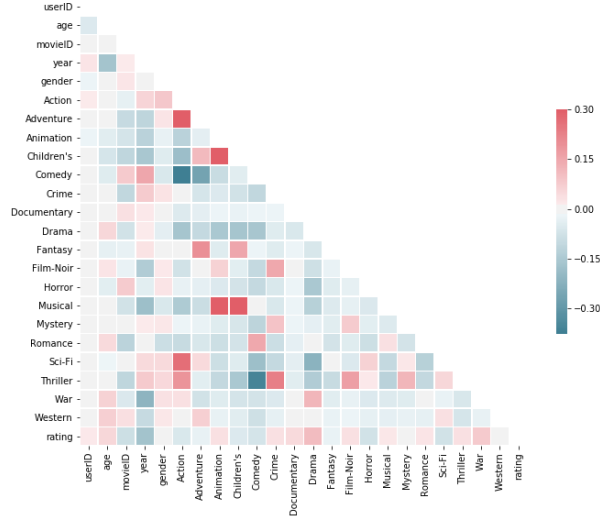
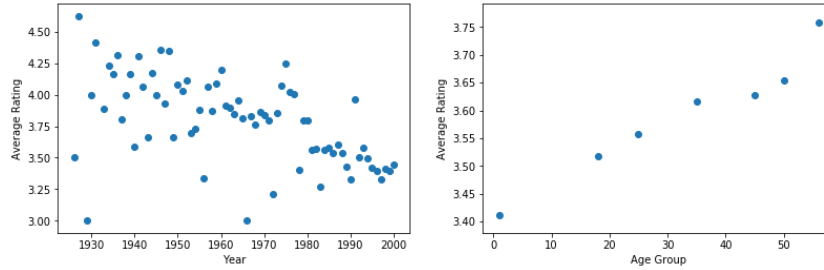
## 22      **2. Methods.**

23      To explore the best model for future movie rating prediction, we try  
24 not only linear regression, ridge and lasso regression but also non-parametric  
25 models including random forest and matrix factorization. We split the dataset  
26 into train set and test set with size 0.2. And we choose the best model based  
27 on the mean square error (MSE) of the validation set.

### 28      **2.1. Data Analysis.**

29      The Netflix dataset contains 31620 observations and 10 variables. The  
30 response variable (rating) is on an integral scale from 1 to 5. The level of  
31 1, 2, 3, 4 and 5 each received a total of 1847, 3350, 8158, 11300 and 6965  
32 ratings respectively. The other key variables include userID, movieID, age  
33 group, gender, year of the movie and genres of the movie. There are 1465  
34 movieID for movies from the year 1926-2000 and 2353 userID for users across  
35 five age groups. Females contributed to 23.58% of the overall ratings. There  
36 are missing data in the genres, but after we transformed data into dummy  
37 variable, we solve the missing data issue, and as a result, we have 25 variables.

38      The correlation plot (FIG.1) shows a relatively weak to moderate cor-  
39 relation between each of the variable. The highest we see is the positive  
40 correlation between musical and comedy.

FIG. 1. *Correlation Map*FIG. 2. *Left: Scatterplot Average Rating vs Year; Right: Scatterplot Average Rating vs Age Group;*

41 The average rating group by year plot (FIG.2) shows a negative linear  
 42 relationship between rating and the year of the movie. The average rating by  
 43 age group plot (FIG. 2) shows a slightly positive linear relationship between  
 44 average rating and age group.

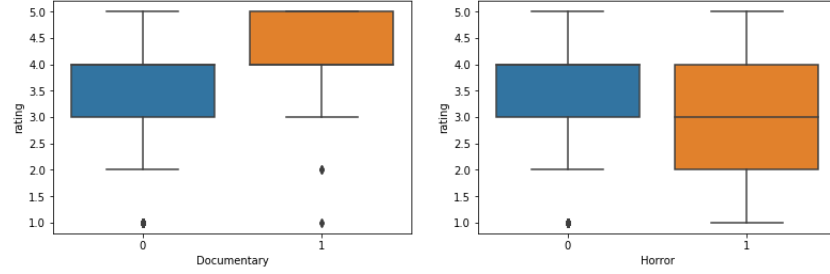


FIG. 3. *Left: Boxplot of Documentary Movie Rating; Right: Boxplot of Horror Movie Rating*

Movies in documentary genera have higher average rating than movies that are not in documentary genera. Movies with genera listed as horror have average ratings lower than movies not in the horror genera.

## 2.2. Original Linear Regression.

$$(2.1) \quad Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Linear regression models the relationship between outcomes and predictors by fitting a linear equation to the observed data. The Linear regression formula is:

$$(2.2) \quad Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_{p-1} X_{i,p-1} + \epsilon_i$$

The intercept  $\beta_0$  is the constant term. The model parameters  $\beta_1$  to  $\beta_{p-1}$  are unknown constants relating the p-1 explanatory variables to the vari-

ables of interest. The  $\epsilon_i$  are normally distributed, independent random error components [6, 10].

The advantages of linear regression models are easy to implement and straightforward interpretations. However, multiple linear regression analysis makes several fundamental assumptions: linear relationship between the outcome and independent variables, uncorrelated errors normally distributed with mean zero and variance  $\sigma^2$ , constant variance, and no or little multicollinearity [10].

Stepwise selection methods are widely utilized to identify predictors for inclusion in regression models [8]. We attempt the backward stepwise selection, which involves setting the significance level at 0.05, starting with a model with all predictors, and remove the variable with the highest non-significant p-value for the coefficient. We repeat this process until no further variables can be deleted.

After five rounds of selection, predictors for musical (p-value 0.950), mystery (p-value 0.864), gender (p-value 0.437), adventure (p-value 0.418) and age (p-value 0.166) are removed from the model sequentially. The final model with MSE of 1.1502 is as follow:

$$\begin{aligned} \hat{Y}_{rating} = & 29.815 - 0.013x_{year} - 0.143x_{Act} + 0.467x_{Ani} - 0.467x_{Child} - 0.087x_{Comedy} \\ & + 0.175x_{Crime} + 0.687x_{Doc} + 0.202x_{Dra} + 0.090x_{Fan} - 0.138x_{Film} - 0.377x_{Hor} \\ & + 0.088x_{Rom} - 0.080x_{Sci} + 0.167x_{Thr} + 0.145x_{War} - 0.108x_{West} \end{aligned}$$

However, variable screening based on statistical significance and stepwise vari-

able selection may lead to unreliable models with biased regression coefficients that need shrinkage[5, 7]. In the next section, we apply Ridge and Lasso regression to address this issue.

### 2.3. Ridge and Lasso Regression.

We apply Ridge and Lasso regularization techniques to the Ordinary Linear Squares(OLS), which shrink the coefficients of the linear model and prevent overfitting of the model[10]. Ridge regression adds the squared magnitude of coefficients as penalty terms to the function and minimizes the function to find the parameters:

$$(2.3) \quad \min(\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \lambda \sum_{j=1}^p \beta^2)$$

The parameter  $\lambda$  controls how much we want to shrink the coefficients. When  $\lambda = 0$ , we restore the equation to OLS. When  $\lambda$  increases, we add more weight to the penalty term and shrink the  $\beta$ 's to a more considerable extent. Lasso regression adds absolute value of magnitude of coefficients as penalty terms and shrinks the less essential features' coefficients to zero and thus useful in removing some features when we have a large number of features [10].

$$(2.4) \quad \min(\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \alpha \sum_{j=1}^p |\beta|)$$

The following graphs show how the coefficients of the predictors change with  $\lambda$  in Ridge and  $\alpha$  in Lasso regression for our data, with Ridge on the left and Lasso on the right:

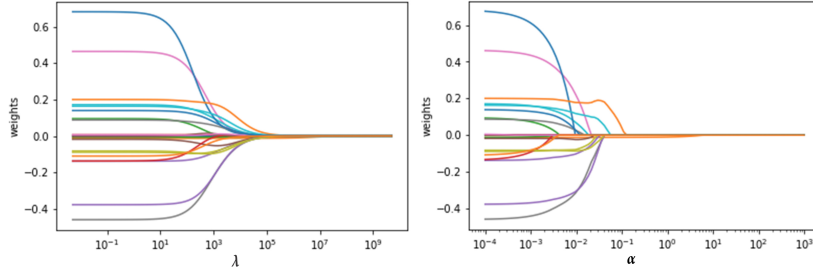


FIG. 4. Left:Regularization of Ridge Regression;Right:Regularization of Lasso Regression

Using *RidgeCV* and *LassoCV* in scikit-learn in python, we find the optimal parameters to be  $\lambda = 9.7$  and  $\alpha = 0.0003$ , which give the smallest MSE of 1.1500 and 1.1499 respectively.

As  $\lambda = 9.7$ , the Ridge regression reduces the MSE in OLS by a number smaller than 0.0001 and shrinks the coefficients of the variables on a small scale as shown in *Figure 4*. As  $\alpha = 0.0003$ , Lasso regression shrinks the coefficients of the variables of movie genres of Musical and Mystery to zero and reduces the MSE in OLS by 0.001. Since Lasso regression performs better than Ridge regression, we display the final result of Lasso regression:

$$\begin{aligned} \hat{Y}_{rating} = & 29.560 + 0.0008x_{age} - 0.013x_{year} - 0.011x_{gender} - 0.135x_{Act} - 0.018x_{adv} \\ & + 0.446x_{Ani} - 0.449x_{Child} - 0.088x_{Comedy} + 0.165x_{Crime} + 0.639x_{Doc} \\ & + 0.199x_{Dra} + 0.082x_{Fan} - 0.110x_{Film} - 0.372x_{Hor} + 0.084x_{Rom} \\ & - 0.082x_{Sci} + 0.160x_{Thr} + 0.133x_{War} - 0.093x_{West} \end{aligned}$$

## 2.4. Random Forest.

A random forest is an ensemble of trees. Each tree consists of split nodes and leaf nodes. The basic idea is that we bootstrap the data for each iteration

114 first, then build a tree on those bootstrap samples. For each split in the tree, we  
 115 randomly select a subset of features and use the best predictors to split. In other  
 116 words, we only consider a subset of the variables for each potential split, and this  
 117 makes it possible to build a diverse set of possible trees. After growing a large  
 118 number of trees, we average those trees to get the prediction for a new outcome.

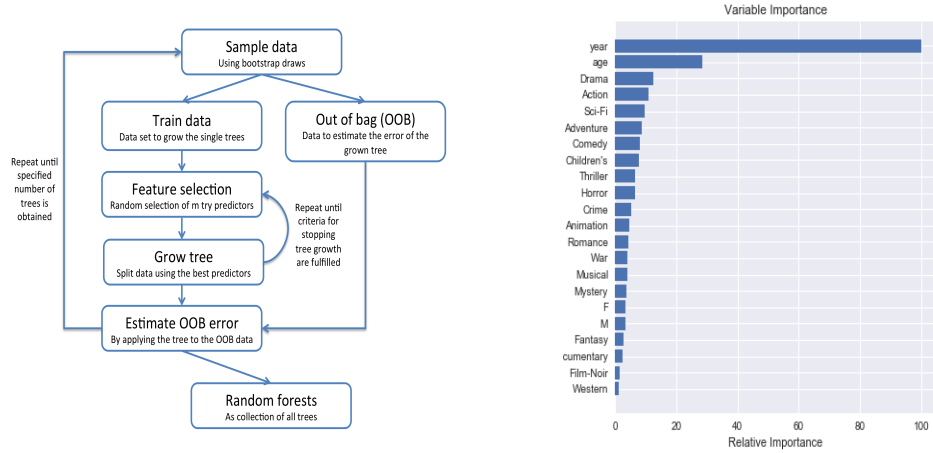


FIG. 5. *Left: The process of building random forest; Right: Diagram of variable importance*

119 This approach is one of the most widely used and highly accurate methods.  
 120 However, this approach can be hard to interpret since we build a large number  
 121 of trees that are averaged together to represent bootstrap samples with bootstrap  
 122 nodes. It is also computational expansive and can lead to other issues such as  
 123 overfitting.

124 To optimize the model performance in the random forest, we apply hyperpa-  
 125 rameter tunning. By using Randomized Search Cross Validation and Grid Search  
 126 with Cross Validation, we finalize the optimization hyperparameter and reduce the



127 MSE from 1.08 to 0.98.

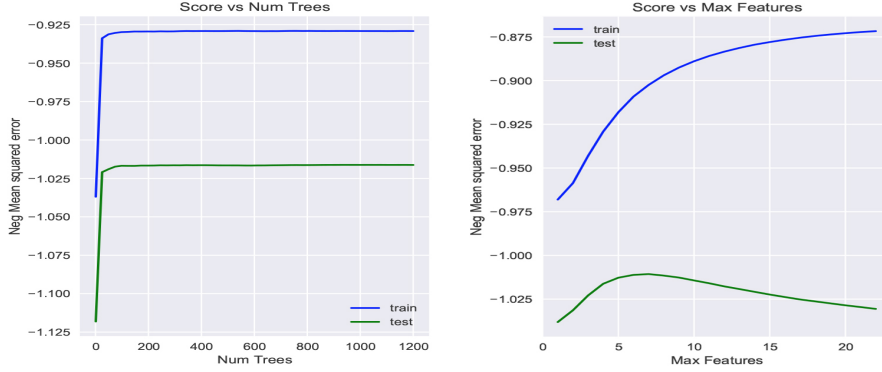


FIG. 6. Left: relationship of number of trees in the forest with negative mean squared error; Right: relationship of maximum features in the forest with negative mean squared error;

## 128 2.5. Matrix Factorization.

129 The idea of Matrix Factorization is approximating the rating matrix  $R$  as the  
 130 product of two matrices:  $R \approx PQ$ , where  $P$  is an  $N \times K$  and  $Q$  is a  $K \times M$  matrix.  
 131 This factorization gives a low dimensional numerical representation of both users  
 132 and movies.

133 In our case, we have lots of unknown ratings of  $R$  ( $\sim 99\%$ ) which cannot be  
 134 represented by zero. Thus, we formulate our problem by minimizing

$$135 \quad (2.5) \quad MSE = \frac{1}{|\mathcal{T}|} \sum_{(u,i) \in \mathcal{T}} e_{ui}^2$$

136 where  $e_{ui} = r_{ui} - \hat{r}_{ui}$  for  $(u, i) \in \mathcal{T}$  and  $\hat{r} = \sum_{k=1}^K p_{uk}q_{ki}$ . And we apply non-  
 137 negative matrix factorization method, which restricts elements of  $P, Q$  to be non-

negative since each matrix element  $r \in [1, 5]$ , and solving the optimization problem

$$(2.6) \quad (P^*, Q^*) = \underset{(P \geq 0, Q \geq 0)}{\operatorname{argmin}} \operatorname{MSE}.$$

Furthermore, we want to introduce regularization to avoid overfitting and this is done by adding a parameter  $\beta$  and modify the squared error as following:

$$(2.7) \quad e_{ui}^2 = (r_{ui} - \hat{r}_{ui})^2 + \frac{\beta}{2} \sum_{k=1}^K (p_{uk}^2 + q_{ki}^2).$$

Then we apply the gradient decent method to find a local minimum of  $MSE$ , i.e. update  $P, Q$  with

$$(2.8) \quad p'_{uk} = p_{uk} + \alpha \frac{\partial}{\partial p_{uk}} e_{ui}^2 = p_{uk} + \alpha(2e_{ui}q_{ki} - \beta p_{uk})$$

$$(2.9) \quad q'_{ki} = q_{ki} + \alpha \frac{\partial}{\partial q_{ki}} e_{ui}^2 = q_{ki} + \alpha(2e_{ui}p_{uk} - \beta q_{ki})$$

Since we want to optimize the  $MSE$  for the validation set, not for the training set, we separate the training set  $\mathcal{T}$  into two part:  $\mathcal{T}_1, \mathcal{T}_2$ , then train with  $\mathcal{T}_1$  and terminate when the  $MSE$  for  $\mathcal{T}_2$  does not decrease during two epochs. And we also select the dimensional parameter  $K \in \{2, 5, 10, 15, 20\}$ , and regularization parameter  $\beta \in (0, 0.3)$  base on the MSE of  $\mathcal{T}_2$ .

With this method, we can achieve more precise prediction without detailed information from users and movies. But the results are not interpretable and need higher computational cost than the previous models. Moreover, there are some upgraded methods by tuning the regularization parameter and learning rate, or pre-process data [9].

### 3. Main Results.

The following table shows the MSE of the validation set by the four different models. We can see that with regularization, the result display only marginal improvement to the linear regression. The two non-parametric models have a more extensive reduction of MSE compare to the linear models. Matrix factorization is the most accurate model with the lowest MSE of 0.9334. However, all models are inadequate with high MSE, and we think it is due to the sparsity of the data and the data with high variance. In addition, all models have more accurate predictions for higher ratings than lower ratings. One explanation is that we have less data for lower ratings, e.g. about 5% of the total ratings is for rating of 1.

Model	OLR	Ridge/Lasso	RF	MF
MSE	1.1502	1.1499	0.9819	0.9334

TABLE 1  
*Result of MSE of the four models*

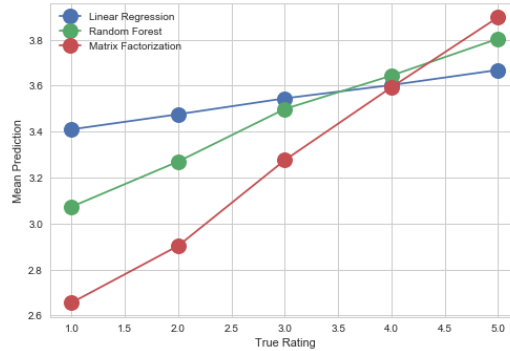


FIG. 7. *prediction vs true value*

### 4. Summary.

In this project, we design and implement parametric and non-parametric statistical models on a real-world data set and compare their performance. Upon

analysis of the movie recommendation data using Python, we explore Python packages such as pandas to deal with data, seaborn and matplotlib to perform data visualization, statsmodels and scikit-learn to build the model, etc. After exploration, we find matrix factorization method, with the largest parameter size, is superior to linear regression and random forest techniques for producing movie recommendation. In future, we will continue to explore other models that may fit the data better to improve the accuracy of the recommendation system further, and practice better data processing strategies such as missing data imputation[4].

## Appendix A. Data Preprocessing.

```

179
180 # Movie Recommendation Lab
181 import numpy as np
182 import pandas as pd
183 import matplotlib.pyplot as plt
184 import matplotlib.cbook as cbook
185 import seaborn as sns
186 from sklearn import datasets, linear_model
187 from sklearn.model_selection import train_test_split
188 from sklearn.model_selection import KFold # import KFold
189 from sklearn import feature_selection
190 from sklearn.linear_model import LinearRegression
191 from sklearn.metrics import mean_squared_error, r2_score
192 from sklearn.linear_model import Ridge, Lasso
193 from sklearn.linear_model import LassoCV
194 import statsmodels.formula.api as smf
195 import statsmodels.api as sm
196
197 #upload data
198 data_users = pd.read_csv("users.csv")
199 data_movies = pd.read_csv("movies.tsv", sep='\t')
200 data_ratings = pd.read_csv("ratings.csv")
201 data_all = pd.read_csv("allData.tsv", sep='\t')
202
203 #view data
204 data_all.head()
205
206 #view data
207 data_all.head()
208 data_all.columns
209 data_all.info()
210
211 #analyze data

```

```

212
213 data_all['age'].describe()
214 data_all['gender'].describe()
215
216 #missing data
217 data_all.count()
218 data_all.isnull().sum()
219
220 # Create a set of dummy variables from the gender and genre variable
221 df_gender = pd.get_dummies(data_all['gender'])
222 df_genre1 = pd.get_dummies(data_all['genre1'])
223 df_genre2 = pd.get_dummies(data_all['genre2'])
224 df_genre3 = pd.get_dummies(data_all['genre3'])
225
226 df_genre = pd.concat([df_genre1, df_genre2, df_genre3]).groupby(level=0).any().astype(int)
227
228 data_new = pd.concat([data_all[['userID', 'age', 'movieID', 'name', 'year']],
229                      df_gender, df_genre, data_all['rating']], axis=1)
230
231 data_X = pd.concat([data_all[['age', 'year']], df_gender, df_genre], axis=1)
232 data_y = data_all['rating']
233
234 # create training and testing vars
235 X_train, X_test, y_train, y_test = train_test_split(data_X, data_y, test_size=0.2, random_state=1234)
236 print (X_train.shape, y_train.shape)
237 print (X_test.shape, y_test.shape)
238
239 # create Cross-Validation K-Fold
240 kf = KFold(n_splits=10, random_state=1234, shuffle=True)
241 for train_index, test_index in kf.split(X_train):
242     X_CV_train, X_CV_test = X_train.iloc[train_index], X_train.iloc[test_index]
243     y_CV_train, y_CV_test = y_train.iloc[train_index], y_train.iloc[test_index]
244
245 #Drop M (only include F in the model)
246 X_train = X_train.drop('M', 1)
247 X_test = X_test.drop('M', 1)

```

## 248 Appendix B. Data visualization.

```

249
250 # Correlation matrix
251 d = data_new
252 # Compute the correlation matrix
253 corr = d.corr()
254 # Generate a mask for the upper triangle
255 mask = np.zeros_like(corr, dtype=np.bool)
256 mask[np.triu_indices_from(mask)] = True
257
258 # Set up the matplotlib figure
259 f, ax = plt.subplots(figsize=(22, 19))
260
261 # Generate a custom diverging colormap
262 cmap = sns.diverging_palette(220, 10, as_cmap=True)
263
264 # Draw the heatmap with the mask and correct aspect ratio

```

```

265 sns.heatmap(corr, mask=mask, cmap=cmap, vmax=.3, center=0,annot=True,
266             square=True, linewidths=.5, cbar_kws={"shrink": .5})
267
268 plt.savefig('fig/correlation_v4.png')
269
270
271
272 #Plot year against average rating by year
273 rating_year=data_all.groupby('year')
274 means=rating_year.mean()
275 plt.scatter(means.index, means['rating'])
276 plt.ylabel('Average Rating')
277 plt.xlabel('Year')
278 plt.show()
279
280 #Plot age against average rating by age
281 rating_year=data_all.groupby('age')
282 means=rating_year.mean()
283 plt.scatter(means.index, means['rating'])
284 plt.ylabel('Average Rating')
285 plt.xlabel('Age Group')
286 plt.show()
287
288

```

## 289 Appendix C. Linear Regression.

```

290 #Linear regression
291 # Train the model using the training sets
292 X = sm.add_constant(X_train)
293 olsmod = sm.OLS(y_train, X).fit()
294 print(olsmod.summary())
295
296 # Test the model using the test sets
297 X2 = sm.add_constant(X_test)
298 ypred = olsmod.predict(X2)
299
300 # Test linear regression model using MSE
301 np.mean((ypred - y_test)**2)
302
303 #Stepwise Selection:
304
305 #Drop Musical
306 Xa = X.drop('Musical', 1)
307 X2a= X2.drop('Musical', 1)
308 olsmod = sm.OLS(y_train, Xa).fit()
309 print(olsmod.summary())
310 ypred = olsmod.predict(X2a)
311 MSE = np.mean((ypred - y_test)**2)
312 print(MSE)
313
314 #Drop Musical+ Mystery
315 Xb = Xa.drop('Mystery', 1)
316 X2b= X2a.drop('Mystery', 1)
317 olsmod = sm.OLS(y_train, Xb).fit()

```

```

318 print(olsmod.summary())
319 ypred = olsmod.predict(X2b)
320 MSE = np.mean((ypred - y_test)**2)
321 print(MSE)
322
323 #Drop Musical+ Mystery+Adventure
324 Xc = Xb.drop('Adventure', 1)
325 X2c= X2b.drop('Adventure', 1)
326 olsmod = sm.OLS(y_train, Xc).fit()
327 print(olsmod.summary())
328 ypred = olsmod.predict(X2c)
329 MSE = np.mean((ypred - y_test)**2)
330 print(MSE)
331
332 #Drop Musical+ Mystery+Adventure+F
333 Xd = Xc.drop('F', 1)
334 X2d= X2c.drop('F', 1)
335 olsmod = sm.OLS(y_train, Xd).fit()
336 print(olsmod.summary())
337 ypred = olsmod.predict(X2d)
338 MSE = np.mean((ypred - y_test)**2)
339 print(MSE)
340
341 #Drop Musical+ Mystery+Adventure+F+age
342 Xe = Xd.drop('age', 1)
343 X2e= X2d.drop('age', 1)
344 olsmod = sm.OLS(y_train, Xe).fit()
345 print(olsmod.summary())
346 ypred = olsmod.predict(X2e)
347 MSE = np.mean((ypred - y_test)**2)
348 print(MSE)
349
350
351 #Linear model using forward_selection using adj R^2 as selection criterion
352 def forward_selected(data):
353     """
354     model: an "optimal" fitted statsmodels linear model
355           with an intercept
356           selected by forward selection
357           evaluated by adjusted R-squared
358     """
359     remaining = set(data.columns)
360     remaining.remove('rating')
361     selected = []
362     current_score, best_new_score = 0.0, 0.0
363     while remaining and current_score == best_new_score:
364         scores_with_candidates = []
365         for candidate in remaining:
366             formula = "{} ~ {} + 1".format('rating',
367                                           ' + '.join(selected + [candidate]))
368             score = smf.ols(formula, data).fit().rsquared_adj
369             scores_with_candidates.append((score, candidate))
370         scores_with_candidates.sort()
371         best_new_score, best_candidate = scores_with_candidates.pop()
372         if current_score < best_new_score:
373             remaining.remove(best_candidate)
374             selected.append(best_candidate)

```

```

375         print(selected)
376         current_score = best_new_score
377         formula = "{} ~ {} + 1".format('rating',
378                                     ' + '.join(selected))
379         model = smf.ols(formula, data).fit()
380         return model
381
382 model = forward_selected(data)
383 print (model.model.formula)
384 print (model.rsquared_adj)
385

```

## 386 Appendix D. Ridge and Lasso.

```

387 #Ridge and Lasso Regression
388 from sklearn.linear_model import Ridge, Lasso
389 from sklearn.linear_model import RidgeCV, LassoCV
390
391 #Find the best alpha in Lasso
392 alphas = np.logspace(-5, -2, 1000)
393 lassocv = linear_model.LassoCV(alphas=alphas, cv=10, random_state=1234)
394 lassocv.fit(data_X, data_y)
395 lassocv_score = lassocv.score(data_X, data_y)
396 lassocv_alpha = lassocv.alpha_
397 lassocv_alpha
398
399
400 #Find the best alpha in Ridge
401 alphas = np.linspace(0,20,200)
402 ridgecv = linear_model.RidgeCV(alphas=alphas, scoring=None, cv=10)
403 ridgecv.fit(data_X, data_y)
404 ridgecv_score = ridgecv.score(data_X, data_y)
405 ridgecv_alpha = ridgecv.alpha_
406 ridgecv_alpha
407
408 #Fit Lasso Regression
409 lasso=Lasso(alpha=0.0003)
410 lasso.fit(X_train, y_train)
411 y_est=lasso.predict(X_test)
412 mse=np.mean(np.square(y_test-y_est))
413 print(mse)
414 print(lasso.coef_)
415 lassodf=pd.DataFrame(lasso.coef_, index=data_X.columns)
416 lasso.intercept_
417
418 #Fit Ridge Regression
419 ridge=Ridge(alpha=9.7)
420 ridge.fit(X_train, y_train)
421 y_est=ridge.predict(X_test)
422 mse=np.mean(np.square(y_test-y_est))
423 print(mse)
424 print(ridge.coef_)
425 ridgedf=pd.DataFrame(ridge.coef_, index=data_X.columns)
426
427 #Plot alpha against coefficients in Ridge

```



```

428 alphas = 10**np.linspace(10,-2,100)*0.5
429 coefs = []
430 for a in alphas:
431     ridge.set_params(alpha=a)
432     ridge.fit(X_train, y_train)
433     coefs.append(ridge.coef_)
434 ax = plt.gca()
435 ax.plot(alphas, coefs)
436 ax.set_xscale("log")
437 plt.axis("tight")
438 plt.xlabel("alpha")
439 plt.ylabel("weights")
440
441 #Plot alpha against coefficients in Lasso
442 alphas = 10**np.linspace(3,-4,100)*0.5
443 coefs = []
444 for a in alphas:
445     lasso.set_params(alpha=a)
446     lasso.fit(X_train, y_train)
447     coefs.append(lasso.coef_)
448 ax = plt.gca()
449 ax.plot(alphas*2, coefs)
450 ax.set_xscale("log")
451 plt.axis("tight")
452 plt.xlabel("alpha")
453 plt.ylabel("weights")

```

## 454 Appendix E. Random Forest.

```

455
456
457 from sklearn.ensemble import RandomForestRegressor
458 from sklearn.model_selection import RandomizedSearchCV
459 from sklearn.model_selection import GridSearchCV
460 from sklearn.metrics import mean_squared_error as MSE
461 import pydot
462 from sklearn.tree import export_graphviz# default parameter
463 rf = RandomForestRegressor(random_state = 1234)
464 rf.fit(X_train, y_train)
465
466 MSE(y_test, rf.predict(X_test))
467 #1.0776258432691597
468
469 # Number of trees in random forest
470 n_estimators = [int(x) for x in np.linspace(start = 200, stop = 2000, num = 10)]
471 # Number of features to consider at every split
472 max_features = ['auto', 'sqrt']
473 # Maximum number of levels in tree
474 max_depth = [int(x) for x in np.linspace(10, 110, num = 11)]
475 max_depth.append(None)
476 # Minimum number of samples required to split a node
477 min_samples_split = [2, 5, 10]
478 # Minimum number of samples required at each leaf node
479 min_samples_leaf = [1, 2, 4]
480 # Method of selecting samples for training each tree

```

```

481 bootstrap = [True, False]
482
483 # Create the random grid
484 random_grid = {'n_estimators': n_estimators,
485               'max_features': max_features,
486               'max_depth': max_depth,
487               'min_samples_split': min_samples_split,
488               'min_samples_leaf': min_samples_leaf,
489               'bootstrap': bootstrap}
490
491
492 # Use the random grid to search for best hyperparameters
493 # First create the base model to tune
494 rf = RandomForestRegressor(random_state = 1234)
495 # Random search of parameters, using 3 fold cross validation,
496 # search across 100 different combinations, and use all available cores
497 rf_random = RandomizedSearchCV(estimator=rf, param_distributions=random_grid,
498                               n_iter = 100, scoring='neg_mean_squared_error',
499                               cv = 10, verbose=2, random_state=1234, n_jobs=-1,
500                               return_train_score=True)
501
502 # Fit the random search model
503 rf_random.fit(X_train, y_train)
504 rf_random.best_params_
505 #{'bootstrap': True,
506  'max_depth': 90,
507  'max_features': 'sqrt',
508  'min_samples_leaf': 4,
509  'min_samples_split': 2,
510  'n_estimators': 1000}
511 rf_random.cv_results_
512
513 best_random = rf_random.best_estimator_
514 random_mse = MSE(best_random, X_test, y_test)
515 print(random_mse)
516 #0.98936
517
518 # Fit the grid search model
519 base_model = RandomForestRegressor(n_estimators = 10, random_state = 1234)
520 base_model.fit(X_train, y_train)
521 MSE(y_test, base_model.predict(X_test))
522 #1.0776
523
524
525 # Create the parameter grid based on the results of random search
526 param_grid = {
527     'bootstrap': [True],
528     'max_depth': [50, 60, 70, 80, 90],
529     'max_features': ['sqrt', 'auto'],
530     'min_samples_leaf': [3, 4, 5],
531     'min_samples_split': [2, 3, 4],
532     'n_estimators': [500, 800, 1000, 1200]
533 }
534
535 # Create a base model
536 rf = RandomForestRegressor(random_state = 1234)
537

```

```

538 # Instantiate the grid search model
539 grid_search = GridSearchCV(estimator = rf, param_grid = param_grid,
540                             cv = 10, n_jobs = -1, verbose = 2, return_train_score=True)
541
542
543 # Fit the grid search to the data
544 grid_search.fit(X_train, y_train)
545
546
547 grid_search.best_params_
548 #{'bootstrap': True,
549  # 'max_depth': 50,
550  # 'max_features': 'sqrt',
551  # 'min_samples_leaf': 4,
552  # 'min_samples_split': 2,
553  # 'n_estimators': 1000}
554 best_grid = grid_search.best_estimator_
555 MSE(y_test, best_grid.predict(X_test))
556 #MSE: 0.9894.
557
558
559 #Training Curves
560 #Number of Trees
561
562 # Grid with only the number of trees changed
563 tree_grid = {'n_estimators': [int(x) for x in np.linspace(1, 1200, 50)]}
564
565 # Create the grid search model and fit to the training data
566 tree_grid_search = GridSearchCV(best_grid, param_grid=tree_grid, verbose = 2,
567                                 n_jobs=-1, cv = 10, scoring = 'neg_mean_squared_error')
568 tree_grid_search.fit(X_train, y_train)
569
570 tree_grid_search.cv_results_
571
572 def plot_results(model, param = 'n_estimators', name = 'Num Trees'):
573     param_name = 'param_%s' % param
574
575     # Extract information from the cross validation model
576     train_scores = model.cv_results_['mean_train_score']
577     test_scores = model.cv_results_['mean_test_score']
578     train_time = model.cv_results_['mean_fit_time']
579     param_values = list(model.cv_results_[param_name])
580
581     # Plot the scores over the parameter
582     plt.subplots(1, 2, figsize=(10, 6))
583     plt.subplot(121)
584     plt.plot(param_values, train_scores, 'b', label = 'train')
585     plt.plot(param_values, test_scores, 'g', label = 'test')
586     plt.legend()
587     plt.xlabel(name)
588     plt.ylabel('Neg Mean squared error')
589     plt.title('Score vs %s' % name)
590
591     plt.subplot(122)
592     plt.plot(param_values, train_time, 'r')
593     plt.xlabel(name)
594     plt.ylabel('Train Time (sec)')

```

```

595     plt.title('Training Time vs %s' % name)
596
597 plot_results(tree_grid_search)
598 plt.savefig("fig/Number_of_Trees.pdf")
599
600
601 #Number of Features at Each Split
602 # Define a grid over only the maximum number of features
603 feature_grid = {'max_features': list(range(1, X_train.shape[1] + 1))}
604 feature_grid_search = GridSearchCV(best_grid, param_grid=feature_grid, cv = 10,
605                                   n_jobs=-1, verbose= 2, scoring = 'neg_mean_squared_error')
606 feature_grid_search.fit(X_train, y_train)
607 plot_results(feature_grid_search, param='max_features', name = 'Max Features')
608 plt.savefig("fig/max_features.pdf")
609
610
611
612
613 #Another Round of Grid Search
614 # Create a base model
615 rf = RandomForestRegressor(random_state = 1234)
616
617 param_grid_2 = {
618     'bootstrap': [True],
619     'max_depth': [50,55, 60],
620     'max_features': [5, 7, 9, 11],
621     'min_samples_leaf': [3, 4, 5],
622     'min_samples_split': [2, 3],
623     'n_estimators': [200,400, 800, 1000, 1200]
624 }
625
626
627 # Instantiate the grid search model
628 grid_search_final = GridSearchCV(estimator = rf, param_grid = param_grid_2,
629                                  cv = 10, n_jobs = -1, verbose = 2, return_train_score=True)
630
631 grid_search_final.fit(X_train, y_train)
632 grid_search_final.best_params_
633 best_grid_final = grid_search_final.best_estimator_
634 MSE(y_test, grid_search_final.predict(X_test))
635 #0.9819552658313127
636 final_model = grid_search.best_estimator_
637
638 #using re_substitution best_estimator
639 rf = RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=50,
640                           max_features='sqrt', max_leaf_nodes=None,
641                           min_impurity_decrease=0.0, min_impurity_split=None,
642                           min_samples_leaf=4, min_samples_split=2,
643                           min_weight_fraction_leaf=0.0, n_estimators=1000, n_jobs=1,
644                           oob_score=False, random_state=1234, verbose=0, warm_start=False)
645
646 rf.fit(X_train, y_train)
647 pred = []
648 pred = rf.predict(X_test)
649
650 print(pred)
651 len(pred)

```

```

652
653
654 d = {'test': y_test, 'predict': pred}
655 df_res = pd.DataFrame(data=d)
656 df_res.groupby('test').mean()
657
658 residual=np.subtract(y_test,pred)
659 print(residual)
660 sns.stripplot(y_test, residual, jitter=True, alpha=.85)
661 sns.despine()
662
663 visual_tree = final_model.estimators_[12]
664 feature_names = list(X_train)
665 export_graphviz(visual_tree, out_file = 'fig/best_tree.dot', feature_names = feature_names,
666                 precision = 2, filled = True, rounded = True, max_depth = None)
667
668 # Use pydot for converting to an image file
669 # Import the dot file to a graph and then convert to a png
670 (graph, ) = pydot.graph_from_dot_file('fig/best_tree.dot')
671 graph.write_png('fig/best_tree.png')
672
673 # Plot feature importance
674 feature_importance = final_model.feature_importances_
675 # make importances relative to max importance
676 feature_importance = 100.0 * (feature_importance / feature_importance.max())
677 sorted_idx = np.argsort(feature_importance)
678 pos = np.arange(sorted_idx.shape[0]) + .5
679 plt.figure(figsize=(6, 6))
680 plt.barh(pos, feature_importance[sorted_idx], align='center')
681 feature_names= np.array(feature_names)
682 plt.yticks(pos, feature_names[sorted_idx])
683 plt.xlabel('Relative Importance')
684 plt.title('Variable Importance')
685 #plt.show()
686 plt.savefig('fig/Variable_Importance_v2.png')
687

```

## 688 Appendix F. Matrix Factorization.

```

689 def matrix_factorization(data, df, R, P, Q, K, alpha, beta, steps=100, ):
690     indices = range(data.shape[0])
691     indices_train, indices_test = train_test_split(indices, test_size=0.1, random_state=123)
692     Q = Q.T
693     e = 100*len(indices_test)
694     P_temp = P.copy()
695     Q_temp = Q.copy()
696     for step in range(steps):
697         for k in indices_train:
698             uid = data.iloc[k,0]
699             mid = data.iloc[k,3]
700             i = df.index.get_loc(uid)
701             j = df.columns.get_loc(mid)
702             eij = R[i][j] - np.dot(P_temp[i,:],Q_temp[:,j])
703             for k in range(K):
704                 P_temp[i][k] = max(P_temp[i][k] + alpha * (2 * eij * Q_temp[k][j] - beta * P_temp[i][k]),0)

```

```

705         Q_temp[k][j] = max(Q_temp[k][j] + alpha * (2 * eij * P_temp[i][k] - beta * Q_temp[k][j]), 0)
706
707
708     #eR = np.dot(P,Q)
709     print(step)
710
711     e_temp = 0
712     for k in indices_test:
713         uid = data.iloc[k,0]
714         mid = data.iloc[k,3]
715         i = df.index.get_loc(uid)
716         j = df.columns.get_loc(mid)
717         e_temp = e_temp + pow(R[i][j] - np.dot(P_temp[i,:],Q_temp[:,j]), 2)
718     #         for k in range(K):
719     #             e = e + (beta/2) * ( pow(P[i][k],2) + pow(Q[k][j],2) )
720     print(e_temp/len(indices_test))
721     #     if (e_temp/31620) < 0.1:
722     #         break
723     if e_temp < e:
724         if e-e_temp < 0.00001:
725             break
726         alpha = alpha*1.05
727         P = P_temp.copy()
728         Q = Q_temp.copy()
729         e = e_temp
730         temp = 0
731     else:
732         alpha = alpha*0.5
733         P_temp = P.copy()
734         Q_temp = Q.copy()
735         temp += 1
736     if temp == 3:
737         break
738     return P, Q.T
739
740 R = df.values
741
742 N = len(R)
743 M = len(R[0])
744 K = 15 #{2, 5, 10, 15, 20}
745 alpha = 0.01
746 beta = 0.25 #[0,0.3]
747
748 np.random.seed(123)
749 P = np.random.rand(N,K)
750 Q = np.random.rand(M,K)
751
752 nP, nQ = matrix_factorization(data_all.iloc[indices_train], df, R, P, Q, K, alpha, beta)
753 nR = np.dot(nP, nQ.T)
754
755 # compute MSE on validation set with given P,Q
756 e = 0
757 r_test = []
758 r_pred = []
759 for k in indices_test:
760     uid = data_all.iloc[k,0]
761     mid = data_all.iloc[k,3]

```

```

762     i = df.index.get_loc(uid)
763     j = df.columns.get_loc(mid)
764     e += (R[i,j]-nR[i,j])**2
765     r_test.append(R[i,j])
766     r_pred.append(nR[i,j])
767
768
769 mse = e/len(indices_test)
770

```

## 771 REFERENCES

- 772 [1] R. M. BELL AND Y. KOREN, *Lessons from the netflix prize challenge*, Acm Sigkdd  
773 Explorations Newsletter, 9 (2007), pp. 75–79.
- 774 [2] J. BENNETT, S. LANNING, ET AL., *The netflix prize*, in Proceedings of KDD cup and  
775 workshop, vol. 2007, New York, NY, USA, 2007, p. 35.
- 776 [3] H. BYSTRÖM, *Movie recommendations from user ratings*, 2013.
- 777 [4] C. CHRISTAKOU, S. VRETTOS, AND A. STAFYLOPATIS, *A hybrid movie recommender*  
778 *system based on neural networks*, International Journal on Artificial Intelligence  
779 Tools, 16 (2007), pp. 771–792.
- 780 [5] S. DERKSEN AND H. J. KESELMAN, *Backward, forward and stepwise automated subset*  
781 *selection algorithms: Frequency of obtaining authentic and noise variables*, British  
782 Journal of Mathematical and Statistical Psychology, 45 (1992), pp. 265–282.
- 783 [6] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *The elements of statistical learning*,  
784 vol. 1, Springer series in statistics New York, 2001.
- 785 [7] M. S. LEWIS-BECK, *Stepwise regression: A caution*, Political Methodology, (1978),  
786 pp. 213–240.
- 787 [8] E. W. STEYERBERG, M. J. EIJKEMANS, AND J. D. F. HABBEMA, *Stepwise selection*  
788 *in small data sets: a simulation study of bias in logistic regression analysis*, Journal  
789 of clinical epidemiology, 52 (1999), pp. 935–942.

- 790 [9] G. TAKÁCS, I. PILÁSZY, B. NÉMETH, AND D. TIKK, *Matrix factorization and neigh-*  
791 *bor based algorithms for the netflix prize problem*, in Proceedings of the 2008 ACM  
792 conference on Recommender systems, ACM, 2008, pp. 267–274.
- 793 [10] S. WEISBERG, *Applied linear regression*, vol. 528, John Wiley & Sons, 2005.