
Music Rating Prediction System

Midway report

Team Member	Team Member	Team Member	Team Member
Dhaval Sonawane	Jignesh Darji	Rohit Mandge	Saurabh sakpal
dsonawa@ncsu.edu	jndarji@ncsu.edu	rnmandge@ncsu.edu	ssakpal@ncsu.edu

1 Background and introduction

Music rating prediction system intends to predict how much rating a user would give on a scale of 1-100 for an EMI company music track they are made to hear. This helps the company know how a particular track or an artist would be received upon launch and also helps them to design a music recommendation system by predicting songs for which the user will give a high rating.

This prediction has to be made considering various attributes of users and artists. The data available in the data set about various users and artists will be merged into a single data frame. The resulting data frame will input to various regression models and their outputs will be analyzed. The model parameters for each of the models will be varied and the best suited values will be selected.

The resulting models will be validated against the validation data set. The performance of each of the model will be analyzed using the error rates captured from validation set. Among these models, we will chose the model that is best able to identify the interesting characteristics available in the data set such as correlation between the ratings and the words used, significance of various attributes, how attributes contributed to ratings, etc. The chosen model will be used to predict the values of the test data set and the error rates of the test data set will be captured and the performance of each of the models will be evaluated.

2 Methods

2.1 Data

User data: User data in the users.csv file has information about user specific interests in music. It contains:

- Information about duration for which the users listened to purchased music and how long they listened to other music tracks.
- Demographic information such as age, gender, employment, etc.
- Answers to a set of 19 questions whose values are numeric in the range 0-100.

Artist data: Artist specific information is available in the words.csv file. It contains:

- Information about different users' description for the artist.
- Each row has list of 82 adjectives (0 indicating the adjective was used and 1 indicating the adjective was not used to describe the artist).
- "LIKE_ARTIST" attribute indicates users' rating for the artist on a scale of 0-100.

45 **Training/Test data:** Training/Test data from train.csv/test.csv has following information:

- 46 • Users' rating of the tracks of various artists.
- 47 • 4 attributes: User ID, Artist ID, Track ID and Rating on a scale of 0-100

48

49 **2.2 Data Preprocessing**

50

51 First challenge for data preprocessing is to merge the 3 data sets into a single data frame that
52 consists of the user and artist information. We merge the datasets by using the user ID and
53 artist ID.

54

55 Handling of missing values is done in the following way:

- 56 • For numerical attributes, we replace them with the median values.
- 57 • For categorical attributes, we replace them with the mode for the categorical attribute.

58

59 **2.3 Methods**

60

61 **2.3.1 Linear Regression model**

62

63 Linear regression is an approach for predicting a dependent variable y based on one or more
64 explanatory variables (or independent variables) denoted X . In our problem statement, the
65 ratings for a new track is the dependent variable and all the attributes describing the user and
66 artist are explanatory variables.

67

68 Although the linear regression model is able to identify the correlation between ratings and
69 some of the attributes, the model fails on many occasions. The model is tested against the
70 validation data set and is found to have a high RMSE value of 27.4545. Hence we conclude
71 that the linear regression model is not suitable for the given complex data set.

72

73 **2.3.2 Linear model split by artist**

74

75 As seen in the section 2.3.1, the linear regression model does not work very well for the given
76 complex dataset. To solve this, we have attempted splitting the dataset by various artists and
77 training a linear regression model for each artist. This significantly reduced the complexity of
78 the data for the linear regression model.

79

80 After combining the results from all the models, RMSE was found to be high for this approach
81 at 25.9347. This error rate is slightly better than linear regression on the entire dataset.
82 However this approach is not advisable as it does not account for any unknown artists in the
83 testing set.

84

85 **2.3.3 Gradient Boosting Model**

86

87 Gradient boosting is a machine learning technique which produces a prediction model by
88 combining weak prediction models, typically decision trees. Each successive tree is built for
89 predicting residuals of the preceding tree. It builds the model in a stage-wise fashion, and
90 generalizes them by allowing optimization of an arbitrary differentiable loss function.
91 Shrinkage is used for shrinking the impact of each additional fitted stage. It reduces the size
92 of incremental steps and thus penalizes the importance of each consecutive iteration.

93 For our problem statement, we built a gradient boosting model with shrinkage = 0.08 and
94 interaction depth = 10 and maximum number of trees = 250. Lower shrinkage level means

higher the time to train the model and higher the interaction depth means higher the depth of each of the trees. GBM model could predict the significance of each of the attributes in determining the prediction ratings of songs. Shrinkage value is kept as low as possible to avoid overfitting of the data. GBM's error rate is lower than linear regression with an RMSE of 24.5728.

2.3.4 Random Forest

Random forests are an ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests improve over decision trees by avoiding overfitting to their training set. Random forest models are known to be very good suit for complex data models.

We do not have any results for the random forest approach as it is still under construction.

3 Experiment and Results

In the process of building the regression models, we were able to discover some important characteristics about the data.

3.1 Error rates

RMSE values for prediction of various models against the test data:

Model	RMSE Value
Linear Model	27.4545
Linear Model by Artist	25.9347
GBM	24.5728
Random Forest	Under Construction
Random Forest by Artist	Under Construction

3.2 Artist Based Prediction

We tried to group data based on artist and build separate models for each group. This improved our accuracy in Linear Models and it will improve the accuracy in random forest too. The RMSE values improved by a significant amount. As we can see in the table above, artists based predictions are proved to be useful.

3.3 Track Ratings

Demographic Information based prediction: We tried to use the demographic information about the users and try to come up with some predictions as shown below (under construction)

3.3.1 Track Ratings by Age of User

131 We are trying to divide the age group of users into 6 buckets and plot the average ratings
132 provided by these age groups. Based on these average ratings we will try to predict ratings of
133 new tracks.

134

135 **3.3.2 Track Ratings by User Work Information**

136 Here the data was grouped by User work information and the average ratings will be calculated
137 for each group. Based on these average ratings, new track ratings were predicted.

138

139 **3.3.3 Track ratings by user similarity**

140 We will try to find similar users with the help of Pearson correlation between every pair of
141 users. Values for all pairs of users above a certain threshold will be considered to be similar.
142 To predict the rating of a song by a user we will try to calculate the mean of ratings given by
143 all the similar users to the given user.

144

145 **4 Conclusion**

146

147 We have tried two different techniques so far: Linear regression and Gradient boosting model.
148 Gradient boosting gives the best result among them. Because of the complexity of the data,
149 linear regression does not perform very well. However, our work on Random Forests is still
150 under construction. Random forest should yield better results as it is well suited for complex
151 data sets.