
Improving Artist Recognition Based on Million Song Dataset

1
2
3 **Shujian Liu**

4 Department of Mechanical and Industrial Engineering
5 University of Massachusetts Amherst
6 Amherst, MA 01003
7 shujian@umass.edu

8 **Abstract**

9 Artist recognition is to predict the artist ID of a song based the other
10 information: analysis sample rate, bars start, beats starts and so on. Machine
11 learning algorithms can learn from the training data and predict the artist ID
12 based on test data. This article is focused on finding the appropriate model
13 and tuning the parameters for high accuracy. The Million Song Dataset is
14 used for the research. It is a large dataset of music consisting of audio
15 features and metadata for one million songs. A variety of methods including
16 supervised learning and unsupervised learning are tested. Good result is
17 generated from LDA method (7.201% compared to 2.103% in KNN with
18 $k=1$).

19 20 **1 Introduction**

21 Columbia University has released a large music dataset named “Million Song Dataset”, for
22 the purposes of encouraging research on machine learning algorithms for commercial scale,
23 providing a reference to evaluate research and help new researchers to get started in the
24 music information retrieval field [1-2].

25 This dataset enables plenty of interesting task such as year recognition (as in homework 3),
26 artist recognition, imputation of missing data, cover song recognition. The original authors
27 have made attempts about artist recognition problem with K Nearest Neighbors algorithm
28 (KNN), however, only low accuracy is achieved.

29 This article is focused on obtaining better accuracy with the machine learning algorithm
30 taught in class. There are significant commercial values in the problem. A successful
31 application of artist/song recognition is Shazam, which is best known for music
32 identification capabilities. Other services include ACRCcloud, Audible Magic and Gracenote.
33

34 **2 Million Song Dataset**

35 The entire Million Song Dataset (MSD) contains 1 million songs and 44,745 unique artists
36 based on their Echo Nest ID, out of that, 18,073 have at least 20 songs (the entire dataset
37 takes 280 GB). Due to limitation of my PC and technical question about accessing it in
38 Amazon AWS, this article used the subset (10,000 songs, 1.8 GB before compression) as
39 provided.

40 The subset can be downloaded from: http://static.echonest.com/millionsongsubset_full.tar.gz

41 In the dataset, each file represents one track with all the related information. The actual file
42 format of each file is HDF5. The wrapper code is provided but user need to install *libhdf5-7*
43 and *python-tables* packages to access data. There are 90 features in MSD, such as analysis
44 sample rate, artist information, bar confidence and start, beats confidence and start.

45 A list is in <http://labrosa.ee.columbia.edu/millionsong/faq> under Field list.

46

47 **3 Related Work**

48 Recognizing the artist from the audio is an interesting task that provides an opportunity for
49 both audio features and machine learning. Artist recognition has been studied for decades,
50 there some early works in [4-7]. However, there was no large scale dataset available to the
51 academics.

52 Realizing this problem, the Laboratory for the Recognition and Organization of Speech and
53 Audio (LabROSA) of Columbia University has built Million Song Dataset and made it
54 available. MSD offers is the 18, 073 artists who have at least 20 songs (compared to 5 artists
55 reported a decade ago).

56 One popular Kaggle completion was made this dataset. In this Million Song Dataset
57 Challenge, participants are working on predicting which songs a user will listen to [8]. It
58 aims at being the best possible offline evaluation of a music recommendation system. Any
59 type of algorithm can be used: collaborative filtering, content-based methods, web crawling.

60 In this artist recognizing problem, they provide two split of the data: 1) in the regular split,
61 each artist has 15 songs in the training set, the rest is in the testing set. 2) in the unbalanced
62 split, each artist has 2/3 of his songs in the training set, the rest in the testing set. This article
63 only considers the unbalanced split which it easier to work with.

64 The authors achieved accuracy of 9.578% for unbalanced split on the entire dataset (each
65 artist has 2/3 of his songs in the training set, the rest in the testing set) with KNN model
66 (K=1). In the present research, using their code on the subset (2805 songs in the test case)
67 leads to an accuracy of 2.2103%. The reason of lower accuracy may be the subset only
68 contains selected tracks with ID starting with A and B. Artist who has tracks from A and B
69 may be more difficult to recognized compare to C-Z.

70 A more recent article [9] uses MLP with backpropagation for this dataset. They boost the
71 accuracy using MLP with backpropagation.

72

73 **4 Proposed Solution**

74 **4.1 Supervised Learning**

75 This classification problem is first solved by supervised learning methods. Supervised
76 learning can analyze the training data and products a model, this model can be used to
77 predict on new data.

78

79 **4.1.1 Basic classifiers**

80 With the starter code from original authors, the hyperparamter of KNN method K is tested.
81 To better resolve this classification problem, a variety of classifiers are going to be used,
82 which include decision tree, Linear Discriminant Analysis (LDA), Support Vector Machines
83 (SVM), Gaussian Naive Bayes (NB), logistic regression (LR) from Scikit-Learn library.
84 Multilayer perceptron (MLP) neural network models such as sknn.mlp, scikit-neuralnetwork
85 and pybrain.

86

87 **4.1.2 Feature Selection**

88 After testing these classifiers with default settings, feature selection is used on the most
89 accurate model. SelectKBest model from Scikit-Learn library. This model can select features
90 according to the k highest scores where k is a user input. The `f_classif` model is used to
91 compute the ANOVA F-value for the provided sample. In feature selection methods,
92 redundant or irrelevant features can be removed without incurring much loss of information.
93 This can both decrease the computation time and increase the predicting accuracy.

94

95 **4.2 Unsupervised Learning**

96 Unsupervised learning methods are also used for this article. Principal Component Analysis

97 (PCA) as a well-used dimensionality reduction model is tested in this part. PCA has the
98 advantages of its low noise sensitivity, the decreased requirements for capacity and memory,
99 and increased efficiency given the processes taking place in a smaller dimension. The
100 hyper-parameter of PCA is the number of components to keep. An exhaustive search over
101 possible to find the best value of hyper-parameter.

102 5 Experiments and Results

104 With the starting KNN code from the Columbia University, a range of K is tested for the
105 unbalanced split of data. It turns out that only k=1 works (59 out of 2805 songs). With k>1,
106 the predicting accuracy is always zero. In KNN methods, everything is far from everything
107 else in high dimensions. Using the whole set of features will make poor prediction since
108 unrelated features will made the songs from one artist far away. In this case, only k=1 can
109 make some good prediction. Trying to take majority vote of several neighbors is unlike to
110 make any correct prediction.

111 Several classifiers are tested and the accuracy is shown in Table 1.

112
113 Table 1: Accuracy of classifiers

114

Classifier Name	Accuracy
KNN (k=1)	2.103%
Decision tree	0.749%
LDA	7.201%
Gaussian NB	2.460%
LR	2.103%
SVM	failed

115

116 The best prediction is given be LDA while decision tree has poor performance. LDA
117 parameters are learned using maximum likelihood, which reduces to using sample estimates.
118 The assumptions LDA makes will rarely be correct for real-world problems. However, the
119 induced linear decision boundaries did perform reasonably well. On the other hand.
120 Decision-tree learners can create over-complex trees that do not generalize the data well.

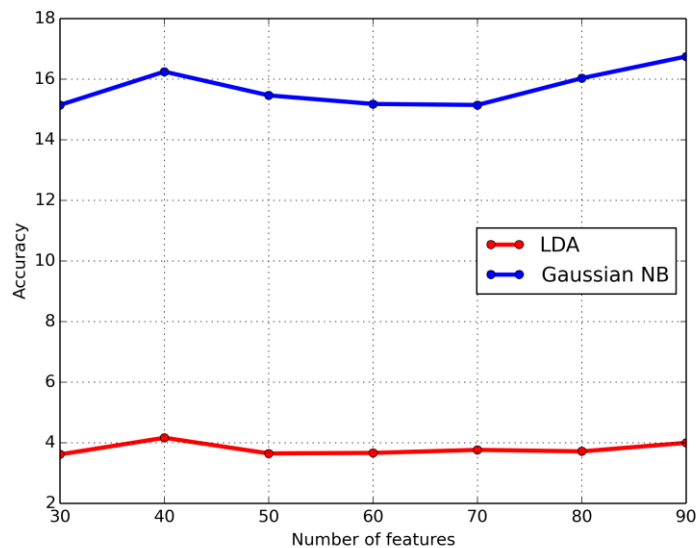
121 In the first and second homework, the well-used ensemble method RandomForestClassifier is
122 also tested here. However, it always return out-of-memory error with a deep tree of large
123 number of forest. The dataset is too large for Random Forest to run in a virtual machine on
124 PC.

125 It is worth mentioned that several neural network models (sknn.mlp, scikit-neuralnetwork
126 and pybrain) has been used in current project, however none of them can work properly or
127 generate reasonable output.

128 By digging into the 90 features. Intuitively, features such as artist_name, the musicbrainz.org
129 ID for this artists, ID of that artist on the service playme.com cannot be used for training. It
130 is cheating. Features such as The Echo Nest song ID, song title are useless or even
131 disturbing. Features such start time of each beat, time of the end of the fade, general loudness
132 of the track are interesting and important for predicting. Feature selection may help increase
133 the accuracy. SelectKBest with f_classif from Scikit-Learn library is used on the LDA and
134 NB since they have best performance in the last part. The number of top features to select k
135 has been searched in range of 40 to 90. The result is shown in Figure 1.

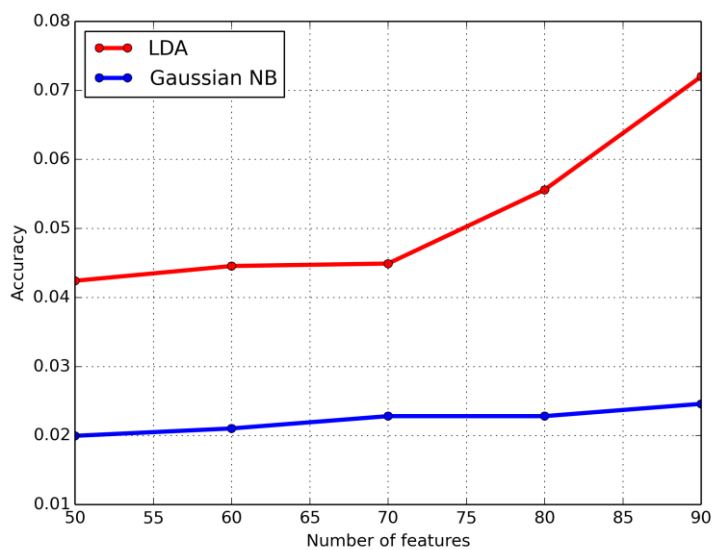
136 Feature selection has negative effects on LDA but slightly improved Gaussian NB. The other
137 way tried to deal with redundant or irrelevant features is dimensional reduction. PCA is used
138 in this part. Number of components to keep has been searched in range of 50 to 90. The
139 results are shown in Figure 2. The negative effect of PCA is even strong for LDA.

140 The main reason that both feature selection and dimensional reduction cannot improve
141 accuracy may be that the irrelevant features are still hard to be find with automatic feature
142 selection and dimensional reduction methods. One potential direct of future work can be
143 manually filtering of these features and reduce them in the dataset to boost the accuracy.
144



145
146
147

Figure 1: Accuracy with feature selection (SelectKBest)



148
149
150
151
152

Figure 2: Accuracy with dimensional reduction (PCA)

5 Discussion and Conclusions

153 The task in real world is much more complicated than in the homework. A lot of effects are
154 done to preprocessing the data and to install packages. The present research is focus on artist
155 recognition problem with machine learning algorithms. It is a large dataset of music

156 consisting of audio features and metadata for one million songs named MSD is used. A
157 variety of methods including supervised learning and unsupervised learning are tested. Both
158 feature selection and dimensional reduction cannot improve the best case with LDA. Best
159 result is generated from LDA method (7.201% compared to 2.103% in KNN with k=1). Tests
160 with deep learning is suggested in the future work. Manually filter the data may be helpful to
161 achieve a higher accuracy and provide more value for commercial use.
162

163 **References**

- 164 [1] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The Million Song
165 Dataset. In Proceedings of the 12th International Society for Music Information Retrieval Conference
166 (ISMIR 2011), 2011.
- 167 [2] Million Song Dataset: <http://labrosa.ee.columbia.edu/millionsong/pages/getting-dataset>
- 168 [3] Bertin-Mahieux, Thierry. Large-Scale Pattern Discovery in Music. Diss. Columbia University,
169 2013.
- 170 [4] Bergstra, James, et al. "Aggregate features and AdaBoost for music classification." Machine
171 learning 65.2-3 (2006): 473-484.
- 172 [5] Whitman, Brian, Gary Flake, and Steve Lawrence. "Artist detection in music with minnowmatch."
173 Neural Networks for Signal Processing XI, 2001. Proceedings of the 2001 IEEE Signal Processing
174 Society Workshop. IEEE, 2001.
- 175 [6] Kaminskas, Marius, and Francesco Ricci. "Contextual music information retrieval and
176 recommendation: State of the art and challenges." Computer Science Review 6.2 (2012): 89-119.
- 177 [7] Henaff, Mikael, et al. "Unsupervised learning of sparse features for scalable audio classification."
178 ISMIR. Vol. 11. 2011.
- 179 [8] Million Song Dataset Challenge. <https://www.kaggle.com/c/msdchallenge>
- 180 [9] Dieleman, Sander, Philémon Brakel, and Benjamin Schrauwen. "Audio-based music classification
181 with a pretrained convolutional network." 12th International Society for Music Information Retrieval
182 Conference (ISMIR-2011). University of Miami, 2011.