
Project-II by Group Rome

Cimen Gokcen

EPFL

`gokcen.cimen@epfl.ch`

Angelopoulos Vasileios

EPFL

`vasileios.angelopoulos@epfl.ch`

Petrescu Viviana

EPFL

`viviana.petrescu@epfl.ch`

Abstract

This report presents experimental results and performance comparison of common machine learning methods on two different tasks. The first task is a computer vision problem, a people detection task. We implemented classifiers that are able to predict to a certain degree if an image contains a standing person or not. The performance was measured using ROC curves. The second task is a music recommender systems challenge. For this problem, we need to predict how many times an user will listen to a certain artist, based on the listening history of this user and the other users. The other subtask of this challenge is to predict the listening counts of new users, for which we do not have music preference information. For this, we are given the friendship information with other users for which we already know the listening counts. The performance of the various method was measured using RMSE.

1 Music Recommendation

1.1 Music recommendation

Collaborative filtering is the part of recommender systems that predicts users' preferences for particular items. The major challenge in predicting users' listening count, as in our task, is that the available user-artist counts are usually too few and sometimes a method's performance relies on a good initial estimation of the unknown entries. These are reasons why we decided to implement besides the common Kmeans, KNN also ALS[cite here], which uses only the known listening counts and avoids dependency on initial estimations.

1.2 Data description

The training data consists in a matrix Y_{train} of size 1774×15082 , corresponding to 1774 users and 15082 artists. Element $Y_{train}(i, j)$ expresses how many times user i has listened to artist j . An entry of 0 means we do not have information for that (user, artist, count) triple. We are also given the friendship graph of the 1774 users stored as an adjacency matrix.

1.3 Exploratory Data Analysis

The listening counts matrix is very sparse with a density of only 0.0026%, corresponding to 69617 (user, artist, count) triples. 1262 artists had 0 listening counts. The variance of the entries is very high, the highest count is 352698 while the average listening count per user and per artist are 5.52 and 5.46 respectively.

A histogram of all the listening counts, shown in Fig1 tells us that the known entries follow a heavy tail distribution. The long tail contains a small number of popular items, the well-known artists, and the rest are located in the heavy tail. One method to transform skewed data such that it becomes more gaussian distributed is to use the Box-Cox transform[cite]

$$\begin{aligned} data(\lambda) &= \log(data), \lambda = 0 \\ &= \frac{data^\lambda - 1}{\lambda}, \lambda \neq 0 \end{aligned}$$

. In our case, we can choose $\lambda = 0$ because our values are very high and positive. This transformation will make the distances between listening counts much smaller and will reduce the influence on RMSE of the (user,artist,count) triples in the long tail.

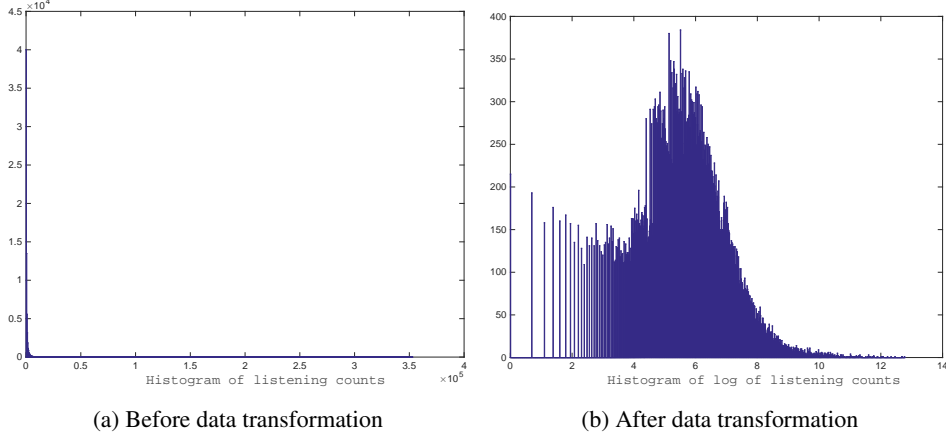


Figure 1: Distribution of all listening counts

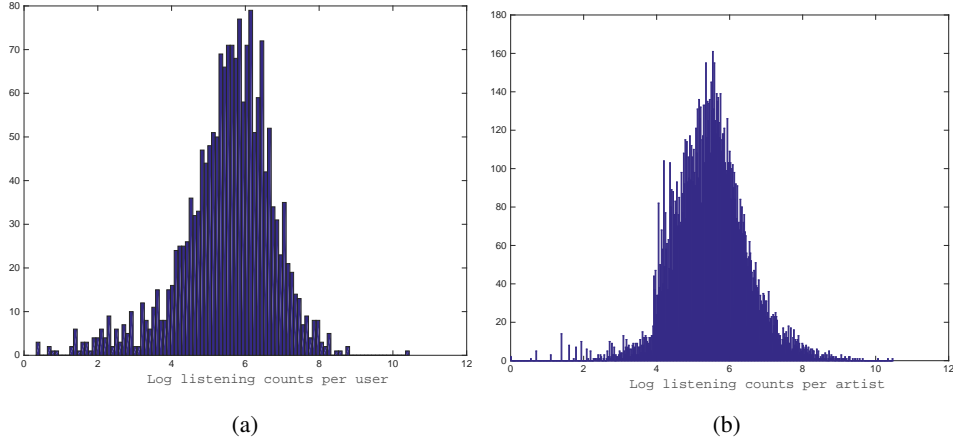


Figure 2: Distribution per user (a) and per artist(b) log of listenining count distribution

1.4 Task 1

In all our experiments we used 10-fold cross validation. For the first task, we randomly omit 10 entries for every artist. Splitting the data was more difficult since we needed to make sure we do not remove all the entries for an artist. If an artist has NA entries with $NA < 10$, then we keep $NA - 1$ entries for testing, to still have one element for training. We do not review here the details of ALS algorithm, the reader can consult the paper[cite] beforehand, since this was not the goal of this report. The methods logALS and logKmeans transformed the training entries using log, but computed the RMSE undoing the transformation, by taking exp of the result.

1.5 Baseline

There are three simple basic predictions one can try: the global average count, the average count per user and the average count per artist prediction. All of these methods give similar RMSE results:

3356, 3293 and 3496 respectively. We note that the final RMSE is composed of various values, small and large. In Fig we can see a plot of the RMSE error terms in log format.

1.5.1 KNN

1.5.2 ALS

1.5.3 logALS

Using 20 features and $\lambda \in [0.01, 0.05, 0.1, 0.5, 1]$ we obtain the results in table 1.5.3 using 10-fold cross validation. The experiments were repeated twice with different seed. We note that we stopped the update steps in the algorithm only after 5 iterations because the update step was very computational intensive.

lambda	RMSE train	RMSE test
0.01	5075 ($\pm 1.0793e+10$)	1.13 (± 0.05)
0.05	3355 (± 895)	1.14 (± 0.02)
0.1	3350 (± 851.5037)	1.58 (± 0.0608)
0.5	3341 (± 851.7858)	13.93 (± 0.8605)
1	3387 (± 842.7349)	25.30 (± 1.48)
1.5	3416 (± 837.1776)	44.25 (± 2.93)

Table 1: Estimated Train and Test RMSE for the two blob models.

We can see that a value of lambda in the middle is a reasonable choice, so we selected lambda to be 0.1 for the next experiments. Our algorithm did a good job training with a RMSE less than 3 but a bad job on unseen data with RMSE ≈ 3000 , meaning it is overfitting the training data and increasing the value of lambda, the regularization parameter did not help.

We notice that the test error across the 10 different splits of the cross validation has values from 2000 to up to 6000.

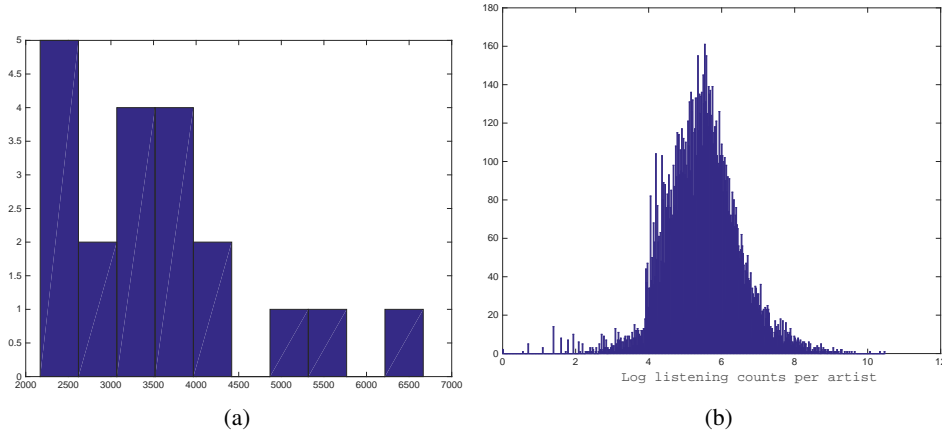


Figure 3: Distribution per user (a) and per artist (b) log of listening count distribution

We varied the number of features from 10, 20, 30, 40, 50, 60 with $\lambda = 0.1$. Giving train and test RMSE of 14 and 3588. So we kept the number of features to 20.

1.5.4 Kmeans

Inspired by the friendship graph information, we tried Kmeans with varying values from 10 to 30 in steps of 5. We plot the mean train and test error using 10 fold cross validation in figure below.

We took log of the nonzeros entries of $Y_{trainnew}$ and before computing RMSE we transformed the data back by taking the exp of the predicted values. We noticed also the fact that sometimes the

training error of Kmeans (reported using exp of data) had very small fluctuations and it was not always decreasing as the algorithm convergence properties would expect. This is due to the fact that Kmeans minimizes squared error but our cost is a little different since we transform the data using exp.

We had a very large variations in the algorithm's output.

One reason for performing Kmeans for clustering over users instead of over artists, although the number of artists is larger than the number of users is to be easier in Task 2.

We can see that our algorithm performs for

1.6 Task 2

1.6.1 Friendship information

Regarding the friendship information, we noticed there are 22 connected components, but many of them contained 1776 nodes and 22904 edges. Using Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Fast unfolding of communities in large networks, in Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P1000 and the Gephi tool to find properties of the graph, it suggested it has about 32 communities, of which 8 had more than 150 members and the others being very small. (22 connected components, of which 20 are very small ; 50 users.) This was run to give us an idea of the number of clusters we could use in Kmeans.

1.7

In this task we are given a set of new users and their friendship information. The task is to predict their listening counts based only on their friendship information with people for which we already have some information about their listening habits. We experiment with three methods. A baseline method in which a (user,artist) count is predicted using the average count of its friends, global average if that information is not present.

The second method uses the 20 clusters from the previous method and predicts the count as an average weight of the listening counts of the clusters of the friends.

The third method uses the 20 clusters from the previous method in a weighted average approach but this time also the friends of the friends information is used.

1.8 Summary

2 People Detection

2.1 Data description and visualization

The original training data *imgs* contains 8545 color images with size 105x43. From every image a HOG descriptor with size 26x10x36 was extracted. All of these descriptors are saved at *feats* and form the training data that we will work on. To process this HOG features we have converted them into a vector of total length 9360.

We notice that there are 7308 images without presence of people (negative images) and 1237 with (positive images). Figure 4 presents a typical example from each category accompanied with the relevant feature descriptor.

Our task is to train various classifiers, so that we will be able to detect the presence of people in unknown cases. For these purposes we evaluate these classifiers, measuring the accuracy of their estimations using Receiver Operating Characteristics (ROC) curves.

2.2 Data preprocessing and used methods

Working on the extracted features, we split the data in order to perform our experiments. In our case we used K-fold validation with $K = 5$ to evaluate the performance of each classifier. Next step is to

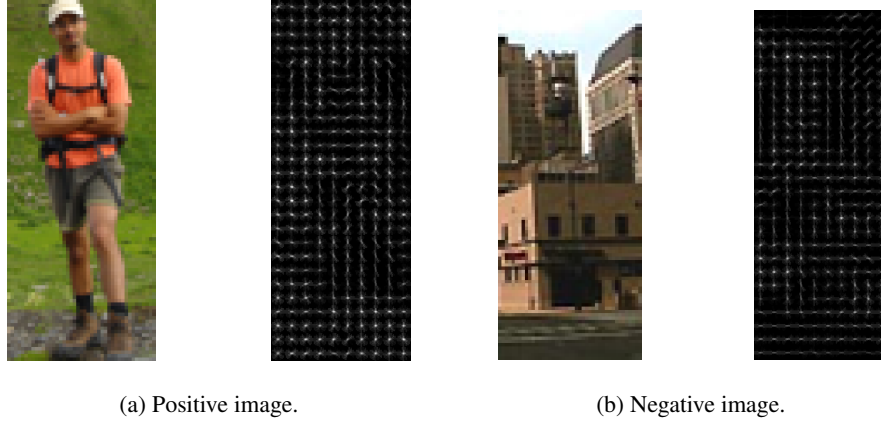


Figure 4: Examples of positive and negative images accompanied with their feature descriptors.

normalize the data so they will be centered with 0 mean and standard deviation 1. Here, we have to mention that we assume that there are no outliers, since the images are manually annotated and they have only two possible values for the classes.

The methods that were tested throughout this experiment are the naive Bayes classifier, the logistic regression and its penalized version, the Support Vector Machine (SVM), the K-NN classifier, and the Neural Network classifiers. All of them are compared also with the random guess of the class in the last subsection of this task.

Moreover, we have to say that at the beginning of the experiments we developed and tried PC Analysis, but, except the huge computational time that was necessary to perform it, it did not give almost any improvements to the dimensionality of our data.

2.3 Naive Bayes classifier

In this part we tried to fit a simple Bayes classifier to our data. The assumptions that we made were that we have data that follow the normal distribution and the estimates of the prior probabilities are empirical and were made from the relative frequencies of the classes in training.

The results from the K-fold validation were not very satisfying, since the average TPR was very low (0.031). The reason that may be hidden behind that are the strong assumptions that we made in order to fit this model. Naive Bayes model assumes that the data that will process every time follow a specific distribution (in our case Gaussian), but also that the value of a particular feature is unrelated to the presence or absence of any other feature, given the variable of the class. In our case, this is possibly wrong, so the classifier cannot fit a appropriate model.

2.4 Logistic and Penalized Logistic Regression

Our purpose in this part was to fit the logistic regression and its penalized version to our data. The value for the parameter α that we used was 10^{-2} , since all of the values close to this performed well enough for our experiments. Moreover, regarding the value for λ in case of the penalized logistic regression, we have to mention that all of the values that were used performed the same. This is the reason that both of the method give the same results, with average TPR equal to 0.754.

During penalization the idea is to avoid the overfitting by imposing penalty on large fluctuations of the model parameters and reduce the influence of the outliers to our model. This is possibly the reason that both models perform equivalently. The data, as we mentioned before, do not have outliers and the fluctuations at the model parameters are avoided by the nature of our data.

2.5 Support Vector Machines

For the certain part of the experiment we used the LIBSVM 3.20 toolbox that includes SVM using various kernels. In our case, we performed analysis and K-fold validation on linear-kernel SVM, polynomial-kernel SVM with degrees 2 and 3, real-basis-kernel SVM and sigmoid-kernel SVM. In all of the cases we used the scores that were given by the program to compute the TPR and the FPR. A comparison of the various cases is given in Figure 5a.

As we can observe radial-basis-kernel SVM have significantly good performance. A fact that can be seen not only by the average TPR, which is equal to 0.893, but also observing the whole range of the FPR, where the classifier gives better results than all of the other SVM cases. The second best in performance is the polynomial SVM with degree 2, which has very similar results with the RBF case for FPR greater than $210e-3$. The rest of the classifier perform well for FPR greater than $210e-2$, with the exception of the linear case, which has quite good results for lower rates, too.

Moreover, in Table 2 we present the average percentage of the successful classifications that were made during the K-fold validation, as they are given by the program.

SVM kernel	Linear	Quadratic	Cubic	RBF	Sigmoid
Average success rate (%)	94.81	97.22	91.45	98.43	84.68

Table 2: Average percentage for the success rate for the various kernels of the SVM classifier, as they are given by the program of the LIBSVM 3.20 toolbox.

2.6 K-NN classifier

Using the K-NN classifier the most important parameter that we have to specify is the number of the neighbouring samples that will be used as a reference for the method so it will decide for any given sample the class that it belongs to. For our experiments we used a ready function, which we modified a lot in order to compute and give in output the scores of the prediction, except of the class labels. We varied the number of neighbours used among 5, 9, 15, 19 and 25. In addition, the euclidean distances were used as the way to specify the distances among the neighbouring samples.

Figure 5b presents the resulting ROC curves that were created using the scores got from the K-fold validation process. There we observe that all of the classifiers show very good performance, especially, for FPR greater than 10^{-3} . Only the 9-NN classifier seems to perform a bit worse for FPR between 10^{-4} and 10^{-3} . These observations can be explained by the distribution of the data, where most probably they are distributed in the N-dimensional ($N = 9360$) space with such a way that their neighbours can specify, with high probability, the class that they belong to.

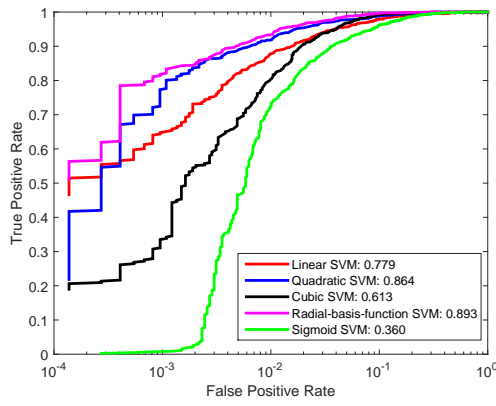
2.7 Neural Network

2.8 Convolutional Neural Network

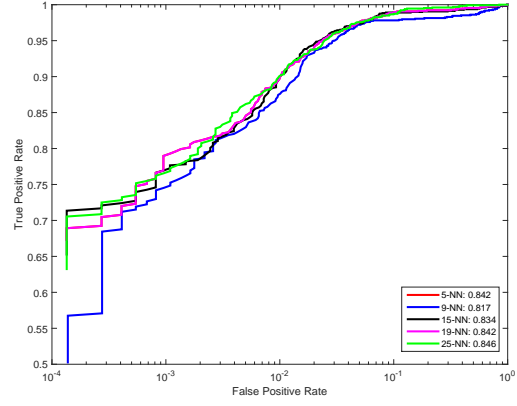
2.9 Comparison

Acknowledgments

We would like to thank the course teaching assistants that helped us a lot during this project preparation, not only during the exercise session, but also at their office hours. Furthermore, the help that our teacher Mohammad Emtiyaz Khan offered us was very important, and we would want to thank him, too. All of the code used and submitted along with this report was written by equally by all of our group members.



(a)



(b)

Figure 5: ROC curves created using: (a) SVM with various kernel types (linear, quadratic, cubic, real-basis-function, sigmoid), (b) K-NN classifier for various number of neighbours for the computations (5, 9, 15, 19, 25).