
Group S653 (4 members):Mini Project Report

Anonymous Author(s)

Affiliation

Address

email

Abstract

This document contains various classification methods used in solving the Mini Project problem for the course Statistical Machine Learning, 1RT700. The goal is to classify a test set of 200 songs and predict which songs Andreas Lindholm is going to like. In this report, we analyze four different classification methods with each member implementing one method that was taught in the course and finally implement the best method in production.

1 Introduction

The idea of the project is to train a classification model using the training dataset containing 750 songs alongwith their features extracted from the spotify web-API and predict the songs that Anreas would like. We achieve this by tuning the parameters and implementing the best model in practice.

2 Classification Methods

This section contains the analysis and testing of four different classification methods. The cross-validation method used for all the four different methods is the K-10 fold cross validation method.

2.1 K - Nearest Neighbours Classifier

For this project we have, among other methods, tried the K-NN approach to classify the songs as either liked or disliked. It's an abbreviation for the k-nearest neighbors and works in the following way: given a new unlabeled data point that we want to classify, we observe how the k-nearest neighbors, for k being a positive integer, have been classified and based on a majority vote among them we would classify this new data point. Of course, the majority vote will result in two probabilities, one for like and the other dislike. We choose ourselves a decision boundary, i.e. if the probability is equal to or greater than that boundary, then the song is classified as like, otherwise dislike. K-NN is a non-parametric method and is therefore different from the parametric methods, that have a fixed number of parameters which are estimated using the training data at hand. The more the training data, the better the estimations will be, but lots of training data is not always available, why choosing non-parametric methods may be a wise idea. For this kind of projects, where the task is to classify songs, K-NN is generally considered as being a good method to choose [4], since there is a high probability that you like a song close to what you earlier have liked.

We have for different values of k, plotted the misclassification error versus k. For a plot for k = 1 to 30 using all the features as input see figure 2.

But sometimes one does not need to include all the features. For example, by plotting the different features against each other and observing the scatter-plot matrix or by studying the co-variance

matrix between two features, we can see that some features does indeed have a relationship. Take the scatter-plot matrix between valence and danceability as an example, featured in figure 1.

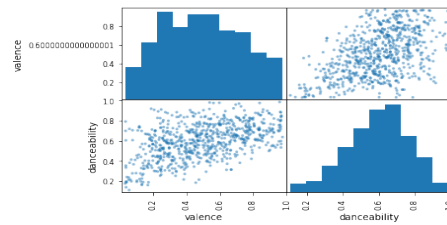


Figure 1: Correlation matrix

This plot is obtained by the following python code: `pd.plotting.scatter-matrix(a)` where a is a vector containing the features valence and danceability from the training set. `pd` is a reference to the python library pandas, and `plotting.scatter-matrix` is simply making a scatter matrix plot of some input data, in this case the vector a .

From this plot we can read that the more a song is danceable, the higher the valence is of that song. Thinking intuitively, this is correct because hos probable is it that someone would want to dance to a depressive song? Therefore, by excluding the valence from the list of features, the result shall not change much, if any at all. By similar arguments, one can exclude even more features from the training set. A second plot is obtained, this time including only the features: acousticness, energy, instrumentalness, key, speechiness. For the resulting plot, see figure 3.

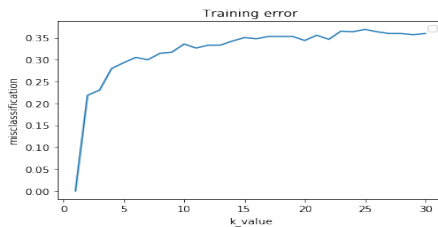


Figure 2: Missclassification error



Figure 3: Missclassification error: selected features

A slightly better result emerges, though not significantly better. The code for this plot and the first one is included in appendix. For both of these missclassification error plots, the missclassification is the error rate for 100 validation data points and the size of the training data set were $750-100 = 650$ data points. Hence for the second plot, we can see that for $k = 10$ we have an error rate of below to 20%, meaning that more that 80% of the songs were classified correctly. The decision boundary is .5 for both cases. Varying this boundary between .4 to .6 could give different results, although not much different that could make an impact.

As a final word it's worth to mention that, although the K-NN method does indeed produce a good result in this case, we choose to use another method in production.

2.2 Discriminant analysis

Within the confines of this course two machine learning methods using discriminant analysis (DA) are discussed: Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis QDA). Although this project will focus on LDA, comparisons to QDA will be made. DA methods don't work by fitting training data to a predefined model, but instead look at each data point and calculate the probability of it falling in a certain class. The method assumes that the training data follows the Gaussian distribution and uses the mean and variance values of this data. Since there are multiple input factors, we're dealing with an average mean ($\hat{\mu}$) and a co-variance matrix ($\hat{\Sigma}$). In LDA this mean is assumed to differ between classes, while the co-variance matrix is constant.

$$p(k|\mathbf{x}_*) = \frac{\hat{\pi}_k \mathcal{N}(\mathbf{x}_* | \hat{\mu}_k, \hat{\Sigma})}{\sum_{j=1}^K \hat{\pi}_j \mathcal{N}(\mathbf{x}_* | \hat{\mu}_j, \hat{\Sigma})} \quad (1)$$

Equation 1 calculates the probability of a new input point \mathbf{x}_* belonging to class k . The numerator in this equation features the occurrence rate of class k in the data ($\hat{\pi}_k$) multiplied by the chance of this new input point existing in a Gaussian/normal distribution with previously mentioned average mean for class k and co-variance matrix. The denominator meanwhile sums the nominator calculations for all classes k . Even though we might not know the actual distribution of data y , we can still learn the parameters in the equation from the training data. We can make a prediction for a certain data point \hat{y}_* by inserting the input values of this data point (\mathbf{x}_*) and choose class k which has the highest probability of occurring.

The results for LDA and QDA were created using the LinearDiscriminantAnalysis and QuadraticDiscriminantAnalysis functions from the sklearn.discriminant_analysis python package. These functions work by creating linear decision boundaries in the domain, which the boundary is located in the input space where the different class predictions for the data points meet.

Figure 5 shows the first attempt at classification and reveals that LDA appears to outperform QDA in this scenario. The average accuracy across all 10 folds in the K-fold validation is about 0.81 for the LDA method and 0.76 for QDA. We're comparing that with an 'always true' prediction, where the outcome of the model is positive for every data point, and which has an accuracy of about 0.6. One way to improve the performance of methods is to focus on a specific number of input factors, instead of including all of them. Figure 4 shows the results of attempting to predict the outcome (like or not like) using only one input factor. The strongest input factors appear to be *accousticness*, *energy*, *loudness* and *speechiness*. Including only those factors in the model slightly elevates the performance of the models (see figure 6 and table 1).

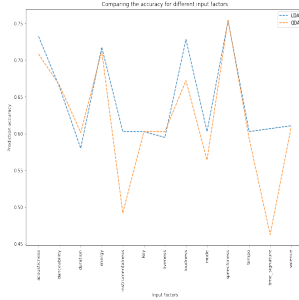


Figure 4: Classification accuracy for models with only one input factor at a time.

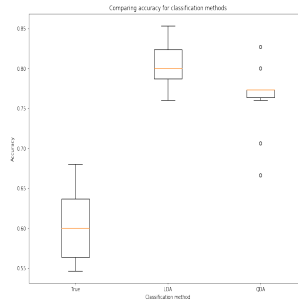


Figure 5: Classification accuracy for methods with all input factors included in model.

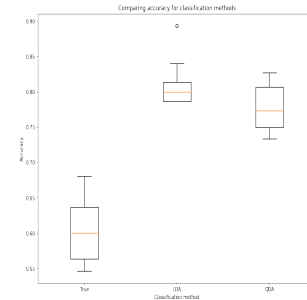


Figure 6: Classification accuracy with four selected input factors included in model.

Another way of optimizing the accuracy using the input factors, is by using an algorithm to compare all the different combinations of input factors and selecting the one with the highest accuracy. The code for this algorithm is featured in the appendices and was used on LDA by comparing all possible combinations of three input factors. Including any more factors took too much time for the algorithm. This resulted in *accousticness*, *loudness* and *speechiness* being selected, with an accuracy of 0.82. This is the highest accuracy we found using any of the DA methods.

Table1: Classification accuracy of methods with either all input factors, or only *accousticness*, *energy*, *loudness* and *speechiness* included.

Method	Accuracy all inputs	Accuracy selected inputs
Always true	0.6027	0.6027
LDA	0.8053	0.8107
QDA	0.7627	0.7773

2.3 Logistic Regression

In the binary classification problem we use logistic regression. We want probability to determine belonging to each class i.e if probability $p > 0.5$ then 'like'. Logistic regression is a parametric model and we are trying to learn a model for probabilities $p(y = 1|x)$ and $p(y = 0|x)$ where:

$$p(y = 1|x) = \frac{e^{\beta^T x}}{e^{\beta^T x} + 1} \quad (2)$$

by maximum likelihood approach to determine β parameters from the training data. We do this by fitting a probability Sigmoid function (2) which has a maximum likelihood to fit the observed data points. Graph of this Sigmoid curve represent dependent variables on the horizontal axis and probability $\in [0, 1]$ on the vertical axis. Likelihood of all data points is equal to the product of observed and unobserved data points probabilities. The algorithm chooses a model for the Sigmoid curve that maximizes the log likelihood of all observed data.

The training data was split 80% and 20% between training and test data i.e (600 and 150). 'Mode' co-variate is changed to categorical variable. Logistic Regression function was used to chose the appropriate Sigmoid curve together with the β parameters. Firstly, all the variables were used to describe the model. Resulting in accuracy ($\frac{TN+TP}{n}$) of 0.6467 and a Receiver Operating Characteristics (ROC) curve as shown in figure 7. We reduce the model to include the most correlated covariates: acoustiness, danceability, energy, loudness, speechiness. Based on the confusion matrix, this reduction in model complexity results in accuracy=0.78 and the ROC curve as in figure 8. Furthermore dropping danceability results in model accuracy=0.8 and ROC as in figure 9, where greater area indicates greater probability of detection.

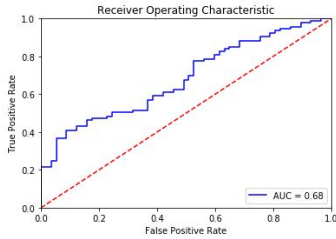


Figure 7: model 1

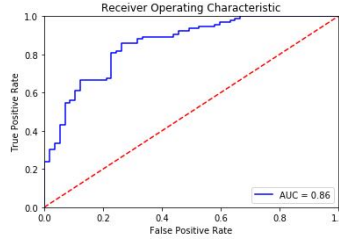


Figure 8: model 2

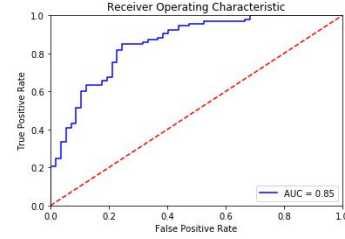


Figure 9: model 3

In order to tune our model we optimize the regularization hyper-parameter C, that is selecting the values of the model that maximize the accuracy, by means of grid search[2]. Where the GridSearchCV function uses C and penalty parameters combinations on a grid, cross-validates the performance of each model and suggests the one obtaining the best results. The combinations can be found in table 2.

Table 2: All combinations of model, penalty parameters and C

model	penalty	C
model 1	11	2.7825594022071245
model 2	12	59.94842503189409
model 3	11	2.7825594022071245

The 10-fold cross validation on the tuned model was carried out. The method was carried out manually, where first we randomize our data points, divide our training data into 10 groups and treat 1 as test data and the remaining as a training data. The process is repeated until each fold has been used once as a test data. Test error is calculated per each fold and average over the 10 folds is obtained as a more accurate estimate of E_{new} . The resulting error rates are as follows: model 1= 0.2, model 2=0.16, model 3=0.16. To help us decide between two last models we use cross-validation score function which again uses 10-fold validation to estimate correct predictions. It allows us to select model2 with accuracy score of 0.8107 winning over model 1 with 0.8035 and model 3 with 0.8085.

2.4 Random Forest Classifier with ADA Boost

Random Forest Classifier comprises of various individual algorithms which combined is termed as an 'ensemble' algorithm. These individual algorithms are nothing but decision trees which in turn are probability spaces borne out of conditional probabilities of data labels given by repetitive partition of the given dataset. With every recursive training of the model the algorithm will calculate the optimum partition of the dataset that would lead to minimum generalization error. Individual trees are weak learners, but Random Forest Classification algorithm takes a majority vote out of many such trees to finally predict the class of the test set. Following is a Validation Score plot of Number of trees Vs Accuracy using Random Forest Classifier.

Boosting is a reinforcement method based on the idea of improving the accuracy by combining many relatively weak learners or classifiers. In the case of Adaptive Boosting (AdaBoost) weak models are attached sequentially and trained using the weighted training data. This goes on till the user-defined number of weak learners have been reached or no refinement on the accuracy can be further made. AdaBoost predicts by calculating the weighted average of all the said weak classifiers. Below is the validation curve plot of Number of estimators Vs Accuracy Score for Cross Validation for 5 splits.

In this section we choose to boost Random Forest Classifier because AdaBoosting works best with non linear models such as Decision Trees which are weak learners. Another reason to choose decision trees is because they are fast learners - the training process is reasonably fast. This is nice as we test in the range of 1 to 600 estimators to get the best possible parameters.

For parameter tuning with AdaBoosting we reduce the learning rate and increase the amount of estimators along with increasing the number of trees for the estimator - Random Forest Classifier.

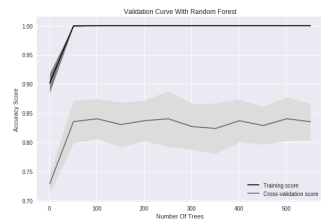


Figure 10: Random Forest Classifier Alone

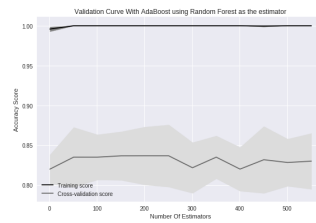


Figure 11: First execution of the model

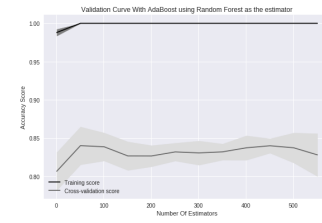


Figure 12: Plot after parameter tuning

As is visible from the accuracy plots, the error for the final tuned model during cross validation which came out to be 0.14 which is by far the best amongst all the models used so far. An accuracy of 81% was achieved in the leaderboard after implementing the solutions of this model.

3 Model Application/Conclusion

We based our optimal model choice on the comparison of 10-fold cross validation obtained from the cross-validation score function. Only the best tuned models from each family were compared and the scores are illustrated in table 3. The Random Forest Classifier with ADA Boosting was implemented on the available test data and achieved a performance of 0.81% correct predictions on the leader board. Therefore we can conclude that Random Forest with ADA Boosting is the most effective way in predicting music preferences taking into account features extracted from songs.

Table 3: Accuracy results of the best tuned version of the models.

Model	Accuracy
K-Nearest Neighbours	0.608
Linear Discriminant Analysis	0.82
Logistic Regression	0.811
Random Forest Classifier with ADA Boost	0.860

4 Reflection

Every machine learning engineer is a consumer of some service and has therefore a dual outlook on ethics of data collection

We would like to discuss the scenario where the music preference data of the user has been used to suggest digital advertisements in the form of popups or online advertisements. More specifically for users with extreme cases of emotional imbalance such as depression who might listen to music of such liking or mood. Song streaming companies can sell such information to pharmaceutical companies who in turn can give out targeted online advertisements to this specific group of unaware users. Within the bracket of said users, the cause of this event can range from redundant expenditures to irrevocable addictions or worse, medication error fatalities. On the ‘upside’, if the music preference data of the aforementioned section of users is given to health authorities or public health organizations, there is a possibility wherein they contact the family/guardian of the said user in a professional manner and provide the necessary attention if required. Although, it is very likely the user will be uncomfortable knowing the data has been used for this particular reason, it is subjective to agree that it has been done so for the benefit of the user and that user alone.

On the other hand, there can be less severe yet still unwarranted outcomes of such data usage. Suggesting google ads for sports shoes/equipments for users who build playlists for working out or jogging for instance is not as harmful. In such scenarios between encouraging unwanted user expenditures or helpful recommendations in some cases and the ‘morality’ of the framework, one can safely agree to process the user data as not only has the privacy policy been agreed upon but quite a few consumers do end up finding the right goods/services.

Furthermore, as an ML engineer working for a company, one has limited options on the course of action to take once it is learnt that the data collection was non-consensual. Thus, to overlook self-integrity for purely economic incentives becomes more of a survival strategy in cases of non sensitive data. It stands to reason that only a freelance or research ML engineer who has a direct influence over a data collection would have a chance to ‘care’ for the consent.

In addition, the morality of non-consensual data collection may be explained by the colloquial ‘aim justifying means’. According to the IEEE Code of Ethics (IEEE 2018): the professional responsibility is to ‘to hold paramount the safety, health, and welfare of the public’. As mentioned earlier, in certain situations, working for the common good might be in conflict with the responsibility to the individual. A perfect example for this is IBM’s machine learning system that analysed patients’ speech patterns to predict psychosis among patients, a serious mental disorder characterized by an impaired relationship with reality. This is again a scenario of personal choice to gauge the damage level for the individual against the benefits for the society. This raises the ideas of a decision boundary - how sensitive is the collected data? In practical everyday life this is not an easily measurable quantity, if at all.

ML engineers in addition respect the humanitarian right of the individual which is entitle to privacy under Article 12 of the 1948 Universal Declaration of Human Rights: ‘No one shall be subjected to arbitrary interference with his privacy, family, home or correspondence, nor to attacks upon his honour and reputation.’

The ethical issues are therefore very limited to legal regulation, business practices and largely dependent on the personal sense of responsibility toward another human being.

202 5 Code

203 This section contains the code of all the four methods implemented in this report.

204 5.1 KNN

```
205
206
207 np.random.seed(1)
208 X_train = trainingSet.drop(columns = ['label'])
209 Y_train = trainingSet['label']
210
211 misclassification = []
212 for k in range(30):
213     model = skl_nb.KNeighborsClassifier(n_neighbors=k+1)
214     beta = model.fit(X_train, Y_train)
215     prediction = beta.predict(X_train)
216     misclassification.append(np.mean(prediction != Y_train))
217
218 K = np.linspace(1,30, 30)
219 plt.plot(K, misclassification)
220 plt.title('Training error')
221 plt.xlabel('k_value')
222 plt.ylabel('misclassification')
223 plt.legend()
224 plt.show()
```

226 **K-NN with selected features**

```
227
228 np.random.seed(1)
229 X_train = trainingSet[['acousticness', 'energy', 'instrumentalness', 'key',
230     'speechiness']]
231 Y_train = trainingSet['label']
232
233 misclassification = []
234 for k in range(30):
235     model = skl_nb.KNeighborsClassifier(n_neighbors=k+1)
236     beta = model.fit(X_train, Y_train)
237     prediction = beta.predict(X_train)
238     misclassification.append(np.mean(prediction != Y_train))
239
240 K = np.linspace(1,30, 30)
241 plt.plot(K, misclassification)
242 plt.title('Training error')
243 plt.xlabel('k_value')
244 plt.ylabel('misclassification')
245 plt.legend()
246 plt.show()
```

248 5.2 LDA

```
249
250 import pandas as pd
251 import numpy as np
252 import matplotlib.pyplot as plt
253 import sklearn.preprocessing as skl_pre
254 import sklearn.linear_model as skl_lm
255 import sklearn.discriminant_analysis as skl_da
256 import sklearn.neighbors as skl_nb
257
258 trainingMusic = pd.read_csv('C:/Users/Erik Jan/Dropbox/MSc Computational
259     Science/UU-61812 Statistical Machine Learning/Project/Data/training_data.csv')
260 Nr_points = len(trainingMusic)
```

```

261 Nr_parts = 10
262 SetDivisor = Nr_points//Nr_parts #Will include all the values. 10 parts of 75 data
263     points each
264
265 acc_LDA = Nr_parts*[0]
266 acc_QDA = Nr_parts*[0]
267 MeanAccLDA = 13*[0]
268 MeanAccQDA = 13*[0]
269
270 outputColumn = 'label'
271 inputColumns = ['acousticness', 'danceability', 'duration', 'energy'
272               , 'instrumentalness' , \
273               'key', 'liveness', 'loudness' , 'mode', 'speechiness', 'tempo', 'time_signature',
274               'valence']
275
276 for factor in range(0,13):
277     inputFactor = [inputColumns[factor]]
278     for i in range(0,Nr_parts):
279         # Setting up data
280         train = trainingMusic[ : SetDivisor*i] #First part of
281             training data
282         train = train.append(trainingMusic[SetDivisor * (i+1) : ]) #Adding second
283             part to the training data
284         validation = trainingMusic[SetDivisor*i : SetDivisor * (i+1)]
285
286         X_train = train[inputFactor]
287         Y_train = train.loc[:, outputColumn]
288         X_validation = validation.loc[:, inputFactor]
289         Y_validation = validation.loc[:, outputColumn]
290
291         # LDA
292         modelLDA = skl_da.LinearDiscriminantAnalysis()
293         modelLDA.fit(X_train, Y_train)
294         prediction = modelLDA.predict(X_validation)
295         acc = np.mean(prediction == Y_validation)
296         acc_LDA[i] = acc
297
298         # QDA
299         modelQDA = skl_da.QuadraticDiscriminantAnalysis()
300         modelQDA.fit(X_train, Y_train)
301         prediction = modelQDA.predict(X_validation)
302         acc = np.mean(prediction == Y_validation)
303         acc_QDA[i] = acc
304
305     MeanAccLDA[factor] = np.mean(acc_LDA)
306     MeanAccQDA[factor] = np.mean(acc_QDA)
307
308
309 # Plotting outcomes
310 plt.plot(inputColumns , MeanAccLDA, linestyle = 'dashed', label = 'LDA')
311 plt.plot(inputColumns , MeanAccQDA, linestyle = 'dashed', label = 'QDA')
312 plt.rcParams["figure.figsize"] = (12,9)
313 plt.legend()
314 plt.title('Comparing the accuracy for different input factors')
315 plt.xlabel('Input factors')
316 plt.xticks(inputColumns, inputColumns, rotation='vertical')
317 plt.ylabel('Prediction accuracy')
318 plt.show()
319
320 #Comparing methods
321 acc_AlwaysTrue = Nr_parts*[0]
322 acc_LDA = Nr_parts*[0]
323 acc_QDA = Nr_parts*[0]
324
325 outputColumn = 'label'

```



```

326 # Here you can select for which combination of input factors you want to test.
327 inputColumns = ['acousticness', 'danceability', 'duration', 'energy'
328                 , 'instrumentalness' , \
329                 'key', 'liveness', 'loudness', 'mode', 'speechiness', 'tempo', 'time_signature',
330                 'valence']
331 #inputColumns = ['acousticness', 'energy', 'loudness', 'speechiness']
332
333 for i in range(0, Nr_parts):
334     # Setting up data
335     train = trainingMusic[ : SetDivisor*i]           #First part of training
336     data
337     train = train.append(trainingMusic[SetDivisor * (i+1) : ]) #Adding second part
338     to the training data
339     validation = trainingMusic[SetDivisor*i : SetDivisor * (i+1)]
340
341     X_train = train.loc[:, inputColumns]
342     Y_train = train.loc[:, outputColumn]
343     X_validation = validation.loc[:, inputColumns]
344     Y_validation = validation.loc[:, outputColumn]
345
346     # Always assume 'like'
347     prediction = len(X_validation)*[1]
348     acc = np.mean(prediction == Y_validation)
349     acc_AlwaysTrue[i] = acc
350
351     # LDA
352     modelLDA = skl_da.LinearDiscriminantAnalysis()
353     modelLDA.fit(X_train, Y_train)
354     prediction = modelLDA.predict(X_validation)
355     acc = np.mean(prediction == Y_validation)
356     #print('Error rate for LDA: ' + str(err))
357     acc_LDA[i] = acc
358
359     # QDA
360     modelQDA = skl_da.QuadraticDiscriminantAnalysis()
361     modelQDA.fit(X_train, Y_train)
362     prediction = modelQDA.predict(X_validation)
363     acc = np.mean(prediction == Y_validation)
364     #print('Error rate for QDA: ' + str(err))
365     acc_QDA[i] = acc
366
367     print("The input columns are: ")
368     print(inputColumns)
369     print
370     # Printing outcomes
371     print("The average accuracy for AlwaysTrue is: " + str(np.mean(acc_AlwaysTrue)))
372     print("The average accuracy for LDA is: " + str(np.mean(acc_LDA)))
373     print("The average accuracy for QDA is: " + str(np.mean(acc_QDA)))
374
375     # Plotting outcomes
376     plt.boxplot([acc_AlwaysTrue, acc_LDA, acc_QDA])
377     plt.xticks([1, 2, 3], ['True', 'LDA', 'QDA'])
378     plt.title("Comparing accuracy for classification methods")
379     plt.xlabel("Classification method")
380     plt.ylabel("Accuracy")
381
382     ### THREE INPUT FACTORS ###
383     # NOTE: Model takes at least 2 minutes to run. You can out-comment the print
384     statements in the code below to keep track of progress.
385     # There are thirteen inout factors to choose from. Let's try to find a combination
386     of three factors which form the best combination.
387
388     MaxAccuracy = 0
389     optimalFactorI = 0
390     optimalFactorJ = 0

```

```

391 optimalFactorK = 0
392
393 outputColumn = 'label'
394 inputColumns = ['acousticness', 'danceability', 'duration', 'energy'
395                 , 'instrumentalness' , \
396                 'key', 'liveness', 'loudness' , 'mode', 'speechiness', 'tempo', 'time_signature',
397                 'valence']
398
399 modelLDA = skl_da.LinearDiscriminantAnalysis()
400
401 for i in range(0,12):
402     for j in range(0,12):
403         for k in range(0,12):
404             if ((i !=j) & (i!=k) & (j!=k)):
405                 #print("The values of i, j and k are: " + str(i)+ " , " + str(j) + "
406                     and " + str(k))
407                 TotalAccuracy = 0
408                 for time in range(0,Nr_parts):
409
410                     # Setting up the data
411                     train = trainingMusic[ : SetDivisor*time]
412                     train = train.append(trainingMusic[SetDivisor * (time+1) : ])
413                     validation = trainingMusic[SetDivisor*time : SetDivisor *
414                         (time+1)]
415
416                     X_train = train.loc[:, [inputColumns[i], inputColumns[j],
417                         inputColumns[k]]]
418                     Y_train = train.loc[:, outputColumn]
419                     X_validation = validation.loc[:, [inputColumns[i],
420                         inputColumns[j], inputColumns[k]]]
421                     Y_validation = validation.loc[:, outputColumn]
422
423                     # Using the model
424                     modelLDA.fit(X_train, Y_train)
425                     prediction = modelLDA.predict(X_validation)
426                     acc = np.mean(prediction == Y_validation)
427                     TotalAccuracy = TotalAccuracy + acc
428
429                     MeanAccuracy = TotalAccuracy/Nr_parts
430                     #print(MeanAccuracy)
431
432                     if (MeanAccuracy > MaxAccuracy):
433                         MaxAccuracy = MeanAccuracy
434                         optimalFactorI = i
435                         optimalFactorJ = j
436                         optimalFactorK = k
437
438 print
439 print("The maximum accuracy is: " + str(MaxAccuracy))
440 print("Optimal factor 1: " + str(inputColumns[optimalFactorI]))
441 print("Optimal factor 2: " + str(inputColumns[optimalFactorJ]))
442 print("Optimal factor 3: " + str(inputColumns[optimalFactorK]))

```

444 5.3 Logistic Regression

```

445 np.random.seed(1)
446 songs=pd.read_csv('training_data.csv')
447 songs.columns.values
448
449 songs["mode"] = songs["mode"].astype('category')
450 #model 1 start:
451 trainI=np.random.choice(songs.shape[0], size=600, replace=False)
452 trainIndex=songs.index.isin(trainI)

```

```

454 train=songs.iloc[trainIndex]
455 test=songs.iloc[~trainIndex]
456
457 model=skl_lm.LogisticRegression()
458
459 X_train=train[['acousticness', 'danceability', 'duration', 'energy',
460               'instrumentalness', 'key', 'liveness', 'loudness', 'mode',
461               'speechiness', 'tempo', 'time_signature', 'valence']]
462 Y_train=train['label']
463
464 X_test=test[['acousticness', 'danceability', 'duration', 'energy',
465              'instrumentalness', 'key', 'liveness', 'loudness', 'mode',
466              'speechiness', 'tempo', 'time_signature', 'valence']]
467 Y_test=test['label']
468
469 model.fit(X_train, Y_train)
470 print('Model summary:')
471 print(model)
472 predict_prob=model.predict_proba(X_test)
473 print('The class order in the model:')
474 print(model.classes_)
475 print('Examples of predicted probabilities for the above classes:')
476 predict_prob[0:5]# inspect the first 5 predictions
477
478 prediction=np.empty(len(X_test), dtype=object)
479 prediction=np.where(predict_prob[:,0]>=0.5,'0','1')
480 prediction[0:5]
481 print(pd.crosstab(prediction, Y_test))
482
483 p = prediction.astype(int) #need to convert since error otherwise
484 y=Y_test.astype(int)
485 np.mean(p==y)
486
487 #Cross Validation
488
489 from sklearn.model_selection import cross_val_score
490 accuracies = cross_val_score(estimator = model, X = X_train, y = Y_train, cv = 10)
491 accuracies.mean()
492
493 # try the ROC curve:
494 import sklearn.metrics as metrics
495 # calculate the fpr and tpr for all thresholds of the classification
496 probs = model.predict_proba(X_test)
497 preds = probs[:,1]
498 fpr, tpr, threshold = metrics.roc_curve(Y_test, preds)
499 roc_auc = metrics.auc(fpr, tpr)
500
501 # method 1: plt
502 import matplotlib.pyplot as plt
503 plt.title('Receiver Operating Characteristic')
504 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
505 plt.legend(loc = 'lower right')
506 plt.plot([0, 1], [0, 1], 'r--')
507 plt.xlim([0, 1])
508 plt.ylim([0, 1])
509 plt.ylabel('True Positive Rate')
510 plt.xlabel('False Positive Rate')
511 plt.show()
512
513 #model 2 start:
514 #Reduce model complexity by selecting only the most correlated covarites: ie
515 acousticness, energy, loudness, speechiness, dacebility
516 X_train2=train[['acousticness', 'danceability', 'energy', 'loudness', 'speechiness']]
517 Y_train2=train['label']
518

```

```

519 X_test2=test[['acousticness', 'danceability','energy','loudness','speechiness']]
520 Y_test2=test['label']
521
522 model2=model.fit(X_train2, Y_train2)
523 print('Model summary:')
524 print(model2)
525
526 predict_prob2=model2.predict_proba(X_test2)
527 print('Examples of predicted probabilities for the above classes:')
528 predict_prob2[0:5]# inspect the first 5 predictions
529
530 prediction2=np.empty(len(X_test2), dtype=object)
531 prediction2=np.where(predict_prob2[:,0]>=0.5,'0','1')
532 prediction2[0:5]
533
534 print(pd.crosstab(prediction2, Y_test2))
535
536 p2 = prediction2.astype(int) #need to convert since error otherwise
537 y2=Y_test2.astype(int)
538 np.mean(p2==y2)
539
540 #Cross Validation
541
542 from sklearn.model_selection import cross_val_score
543 accuracies = cross_val_score(estimator = model2, X = X_train2, y = Y_train2, cv =
544     10)
545 accuracies.mean()
546
547 import sklearn.metrics as metrics
548 # calculate the fpr and tpr for all thresholds of the classification
549 probs = model2.predict_proba(X_test2)
550 preds = probs[:,1]
551 fpr, tpr, threshold = metrics.roc_curve(Y_test2, preds)
552 roc_auc = metrics.auc(fpr, tpr)
553
554 # method 1: plt
555 import matplotlib.pyplot as plt
556 plt.title('Receiver Operating Characteristic')
557 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
558 plt.legend(loc = 'lower right')
559 plt.plot([0, 1], [0, 1], 'r--')
560 plt.xlim([0, 1])
561 plt.ylim([0, 1])
562 plt.ylabel('True Positive Rate')
563 plt.xlabel('False Positive Rate')
564 plt.show()
565
566 #model 3 start:
567 #reduce model again:
568 X_train3=train[['acousticness','energy','loudness','speechiness']]
569 Y_train3=train['label']
570
571 X_test3=test[['acousticness','energy','loudness','speechiness']]
572 Y_test3=test['label']
573
574 model3=model.fit(X_train3, Y_train3)
575 print('Model summary:')
576 print(model3)
577
578 predict_prob3=model3.predict_proba(X_test3)
579 print('Examples of predicted probabilities for the above classes:')
580 predict_prob3[0:5]# inspect the first 5 predictions
581
582 prediction3=np.empty(len(X_test3), dtype=object)
583 prediction3=np.where(predict_prob3[:,0]>=0.5,'0','1')

```

```

584 prediction3[0:5]
585
586 print(pd.crosstab(prediction3, Y_test3))
587
588 p3 = prediction3.astype(int) #need to convert since error otherwise
589 y3=Y_test3.astype(int)
590 np.mean(p3==y3)
591
592 #Cross Validation
593
594 from sklearn.model_selection import cross_val_score
595 accuracies = cross_val_score(estimator = model3, X = X_train3, y = Y_train3, cv =
596     10)
597 accuracies.mean()
598
599 # try the ROC curve:
600 import sklearn.metrics as metrics
601 # calculate the fpr and tpr for all thresholds of the classification
602 probs = model3.predict_proba(X_test3)
603 preds = probs[:,1]
604 fpr, tpr, threshold = metrics.roc_curve(Y_test3, preds)
605 roc_auc = metrics.auc(fpr, tpr)
606
607 # method 1: plt
608 import matplotlib.pyplot as plt
609 plt.title('Receiver Operating Characteristic')
610 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
611 plt.legend(loc = 'lower right')
612 plt.plot([0, 1], [0, 1], 'r--')
613 plt.xlim([0, 1])
614 plt.ylim([0, 1])
615 plt.ylabel('True Positive Rate')
616 plt.xlabel('False Positive Rate')
617 plt.show()
618
619 #tune the parameter - hyperparameter optimisation using grid search: for every
620     available model; then ccompare again:
621 # try to find optimal hyper paramter:
622 # Load libraries
623 import numpy as np
624 from sklearn import linear_model, datasets
625 from sklearn.model_selection import GridSearchCV
626
627 # Create logistic regression
628 logistic = linear_model.LogisticRegression()
629
630 # Create regularization penalty space
631 penalty = ['l1', 'l2']
632
633 # Create regularization hyperparameter space
634 C = np.logspace(0, 4, 10)
635
636 # Create hyperparameter options
637 hyperparameters = dict(C=C, penalty=penalty)
638
639 # Create grid search using 5-fold cross validation
640 clf = GridSearchCV(logistic, hyperparameters, cv=5, verbose=0)
641
642 # Fit grid search
643 best_model1 = clf.fit(X_test,Y_test)
644
645 # View best hyperparameters
646 print('Best Penalty:', best_model1.best_estimator_.get_params()['penalty'])
647 print('Best C:', best_model1.best_estimator_.get_params()['C'])
648

```

```

649 #improving model 1:
650 #apply this penalty:
651 model_best1=skl_lm.LogisticRegression(penalty = 'l1', C =
652     2.7825594022071245,random_state = 0)
653 model_best1.fit(X_train, Y_train)
654 predict_probbest1=model_best1.predict_proba(X_test)
655
656
657 prediction1=np.empty(len(X_test), dtype=object)
658 prediction1=np.where(predict_probbest1[:,0]>=0.5,'0','1')
659 prediction1[0:5]
660
661 #Cross Validation
662 model_best1c=skl_lm.LogisticRegression(penalty = 'l1', C =
663     2.7825594022071245,random_state = 0)
664
665 from sklearn.model_selection import cross_val_score
666 accuracies = cross_val_score(estimator =model_best1c , X = X_train, y = Y_train, cv
667     = 10)
668 accuracies.mean()
669
670 randomize_indices=np.random.choice(songs.shape[0], songs.shape[0], replace=False)
671 misclassification=np.zeros((10,200))
672
673 for i in range(10):
674     n=np.ceil(songs.shape[0]/10)# number of samples in each fold
675     validationIndex=np.arange(i*n,min(i*n+n, songs.shape[0]),1).astype('int')
676     randomize_validationIndex=randomize_indices[validationIndex]
677     train=songs.iloc[~songs.index.isin(randomize_validationIndex)]
678     validation=songs.iloc[randomize_validationIndex]
679
680     X_train=train[['acousticness', 'danceability', 'duration', 'energy',
681         'instrumentalness', 'key', 'liveness', 'loudness', 'mode',
682         'speechiness', 'tempo', 'time_signature', 'valence']]
683     Y_train=train['label']
684     X_validation=validation[['acousticness', 'danceability', 'duration', 'energy',
685         'instrumentalness', 'key', 'liveness', 'loudness', 'mode',
686         'speechiness', 'tempo', 'time_signature', 'valence']]
687     Y_validation=validation['label']
688
689     model_f1=skl_lm.LogisticRegression(penalty = 'l1', C =
690         2.7825594022071245,random_state = 0,solver='liblinear')
691     model_f1.fit(X_train, Y_train)
692     prediction=model_f1.predict(X_validation)
693     err=np.mean(prediction!=Y_validation)
694     print('Error rate for logistic regression:'+str(err))
695
696 # Fit grid search for second reduced model:
697 best_model2 = clf.fit(X_test2,Y_test2)
698
699 # View best hyperparameters
700 print('Best Penalty:', best_model2.best_estimator_.get_params()['penalty'])
701 print('Best C:', best_model2.best_estimator_.get_params()['C'])
702
703 #improving model 2:
704 #apply this penalty:
705 model_best2=skl_lm.LogisticRegression(penalty = 'l2', C =59.94842503189409
706     ,random_state = 0)
707 model_best2.fit(X_train2, Y_train2)
708 predict_probbest2=model_best2.predict_proba(X_test2)
709
710
711 prediction2=np.empty(len(X_test2), dtype=object)
712 prediction2=np.where(predict_probbest2[:,0]>=0.5,'0','1')
713 prediction2[0:5]

```

```

714
715 #Cross Validation
716 model_best2c=skl_lm.LogisticRegression(penalty = 'l2', C =
717     59.94842503189409,random_state = 0)
718
719 from sklearn.model_selection import cross_val_score
720 accuracies = cross_val_score(estimator =model_best2c , X = X_train2, y = Y_train2,
721     cv = 10)
722 accuracies.mean()
723
724 #for second penalized model:
725 randomize_indices2=np.random.choice(songs.shape[0], songs.shape[0], replace=False)
726 misclassification2=np.zeros((10,200))
727
728 for i in range(10):
729     n=np.ceil(songs.shape[0]/10)# number of samples in each fold
730     validationIndex2=np.arange(i*n,min(i*n+n, songs.shape[0]),1).astype('int')
731     randomize_validationIndex2=randomize_indices2[validationIndex2]
732     train2=songs.iloc[~songs.index.isin(randomize_validationIndex2)]
733     validation2=songs.iloc[randomize_validationIndex2]
734
735     X_train2=train2[['acousticness',
736         'danceability','energy','loudness','speechiness']]
737     Y_train2=train2['label']
738     X_validation2=validation2[['acousticness',
739         'danceability','energy','loudness','speechiness']]
740     Y_validation2=validation2['label']
741
742 model_f2=skl_lm.LogisticRegression(penalty = 'l2', C =59.94842503189409
743     ,random_state = 0,solver='liblinear')
744 model_f2.fit(X_train2, Y_train2)
745 prediction2=model_f2.predict(X_validation2)
746 err2=np.mean(prediction2!=Y_validation2)
747 print('Error rate for logistic regression:'+str(err2))
748
749 # Fit grid search for third reduced model:
750 best_model3 = clf.fit(X_test3,Y_test3)
751
752 # View best hyperparameters
753 print('Best Penalty:', best_model3.best_estimator_.get_params()['penalty'])
754 print('Best C:', best_model3.best_estimator_.get_params()['C'])
755
756 #improving model 3:
757 #apply this penalty:
758 model_best3=skl_lm.LogisticRegression(penalty = 'l1', C
759     =2.7825594022071245,random_state = 0)
760 model_best3.fit(X_train3, Y_train3)
761 predict_probbest3=model_best3.predict_proba(X_test3)
762
763 prediction3=np.empty(len(X_test3), dtype=object)
764 prediction3=np.where(predict_probbest3[:,0]>=0.5,'0','1')
765
766 #for second penalized model:
767 randomize_indices3=np.random.choice(songs.shape[0], songs.shape[0], replace=False)
768 misclassification3=np.zeros((10,200))
769
770 for i in range(10):
771     n=np.ceil(songs.shape[0]/10)# number of samples in each fold
772     validationIndex3=np.arange(i*n,min(i*n+n, songs.shape[0]),1).astype('int')
773     randomize_validationIndex3=randomize_indices3[validationIndex3]
774     train3=songs.iloc[~songs.index.isin(randomize_validationIndex3)]
775     validation3=songs.iloc[randomize_validationIndex3]
776
777     X_train3=train3[['acousticness','energy','loudness','speechiness']]
778     Y_train3=train3['label']

```

```

779     X_validation3=validation3[['acousticness','energy','loudness','speechiness']]
780     Y_validation3=validation3['label']
781
782     model_f3=skl_lm.LogisticRegression(penalty = 'l1', C = 2.7825594022071245
783     ,random_state = 0,solver='liblinear')
784     model_f3.fit(X_train3, Y_train3)
785     prediction3=model_f3.predict(X_validation3)
786     err3=np.mean(prediction3!=Y_validation3)
787     print('Error rate for logistic regression:'+str(err3))
788     prediction3[0:5]
789
790     #Cross Validation
791     model_best3c=skl_lm.LogisticRegression(penalty = 'l1', C =
792     2.7825594022071245,random_state = 0)
793
794     from sklearn.model_selection import cross_val_score
795     accuracies = cross_val_score(estimator =model_best3c , X = X_train3, y = Y_train3,
796     cv = 10)
797     accuracies.mean()

```

799 5.4 Random Forest Classifier with AdaBoost

```

800     #Read Data
801     np.random.seed(1)
802     data = pd.read_csv('Data/training_data.csv');
803     print(data.head())
804
805     #MainData
806     X = data.drop(columns=['label'])
807     y = data['label']
808
809     #Split Using Train test split
810     X_train, X_test, y_train, y_test = train_test_split(X, y,
811     test_size=0.2,random_state=0)
812
813     #Feature Scaling
814     from sklearn.preprocessing import StandardScaler
815     sc = StandardScaler()
816     X_train = sc.fit_transform(X_train)
817     X_test = sc.fit_transform(X_test)
818
819     #Fitting the model
820     rfc = RandomForestClassifier()
821     ada = AdaBoostClassifier(n_estimators=100, base_estimator = rfc,learning_rate=0.001)
822     ada.fit(X_train,y_train)
823
824     #Cross Validation and AccuracyScore
825     from sklearn.model_selection import cross_val_score
826     accuracies = cross_val_score(estimator = ada, X = X_train, y = y_train, cv = 10)
827     accuracies.mean()
828
829     #Predictions and string output for inputting in the leaderboard
830     ada_predict = ada.predict(X_test)
831     string = ''
832     for i in ada_predict:
833         string += str(i)
834     print(string)
835
836     #Confusion Matrix and Classification Report
837     print("=== Confusion Matrix ===")
838     print(confusion_matrix(y_test, ada_predict))
839     print('\n')
840     print("=== Classification Report ===")

```



```

842 print(classification_report(y_test, ada_predict))
843 print('\n')
844
845 #Plotting Validation Curve
846 param_range = np.arange(1,600,50)
847 train_scores, test_scores = validation_curve(rfc,
848                                             X_train,
849                                             y_train,
850                                             param_name="n_estimators",
851                                             param_range=param_range,
852                                             cv=5,
853                                             scoring="accuracy",
854                                             n_jobs=-1)
855 # Calculate mean and standard deviation for training set scores
856 train_mean = np.mean(train_scores, axis=1)
857 train_std = np.std(train_scores, axis=1)
858
859 # Calculate mean and standard deviation for test set scores
860 test_mean = np.mean(test_scores, axis=1)
861 test_std = np.std(test_scores, axis=1)
862
863 # Plot mean accuracy scores for training and test sets
864 plt.plot(param_range, train_mean, label="Training score", color="black")
865 plt.plot(param_range, test_mean, label="Cross-validation score", color="dimgrey")
866
867 # Plot accuracy bands for training and test sets
868 plt.fill_between(param_range, train_mean - train_std, train_mean + train_std,
869                 color="gray")
870 plt.fill_between(param_range, test_mean - test_std, test_mean + test_std,
871                 color="gainsboro")
872
873 # Create plot
874 plt.title("Validation Curve With Random Forest")
875 plt.xlabel("Number Of Trees")
876 plt.ylabel("Accuracy Score")
877 plt.tight_layout()
878 plt.legend(loc="best")
879 plt.show()
880

```

881 References

- 882 1 [1] Article title: Machine Learning Basics with the K-Nearest Neighbors Algorithm Website title: Towards
883 Data Science URL: [https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-](https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761)
884 [algorithm-6a6e71d01761](https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761)
- 885 [2]Article title: Hyperparameter Optimization in Machine Learning Website title: DataCamp Community URL:
886 <https://www.datacamp.com/community/tutorials/parameter-optimization-machine-learning-models>
- 887 [3]Article title: 3.5. Validation curves: plotting scores to evaluate models — scikit-learn 0.20.2 documentation
888 Website title: Scikit-learn.org
URL: https://scikit-learn.org/stable/modules/learning_curve.html