

Do (wo)men talk too much in films? Project in Machine Learning

Anonymous Author(s)
Number of group members: 2

Abstract

In this document, we looked at the typical 1037 films data set of Hollywood movies and their leading roles. We analyzed the data set, brought insights out of it, and tried to come up with the assumption that we could predict the lead role in a movie based on what features (given in the data set). For example, Revenue generated, year, and the number of dialogues could be one of the main features to predict the leading or main role in the films. Later we implemented Logistic Regression and Random Forest models and checked how authentic our assumptions were based on the data analysis we did earlier. We also added a naive classifier and compared it with our models. Further, using Cross Validation, we checked the accuracy of our models and came up with one model that gave us more accurate predictions. This model will be used on the test data set to check how well it performs.

1 Introduction

Film Industry has been here for a very long now, and it has had a huge impact on any country's image and contributed to its economy as a separate industry and provides a huge amount of viewership. While this industry has been here, there have been some flaws in the system, which has been lacking in keeping fair standards and equal opportunities, causing the serious issue of gender biasness. For reference, let's have a look at the graph of the Star War movie series and how gender biasness exists there, and how fewer female roles are given in all the series of Star Wars. ¹.

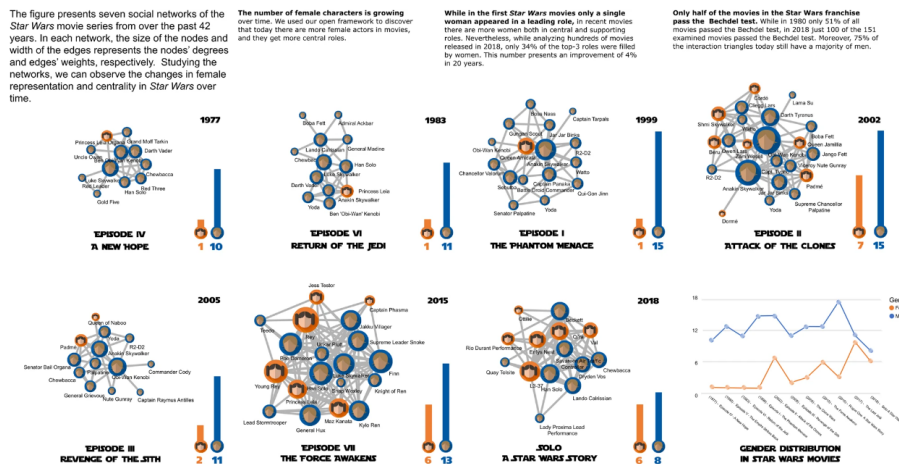


Figure 1: Star Wars Films

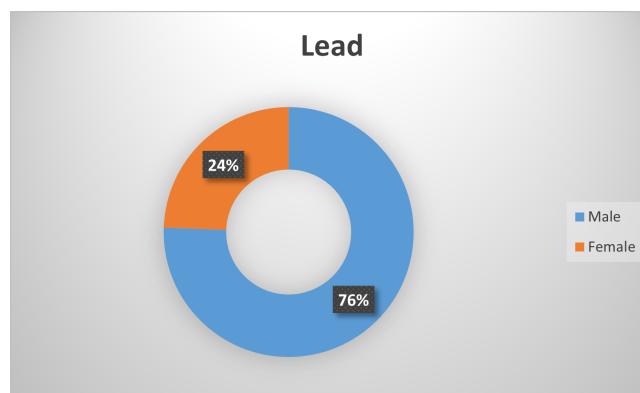
¹<https://www.nature.com/articles/s41599-020-0436-1/figures/1>

21 We can clearly see in the graph above that over the past few years of the episodes of star wars, the
22 number of female roles has been very low as compared to men, and they have been kept aside. We
23 will be closely looking at a similar data set where we have a list of films and their characteristics, and
24 we will be evaluating if the gender of the lead role is based on certain characteristics or features of
25 the films. If that is true, then what are those main features, and how is it predicted? We have a data
26 set of 1037 films and the gender of lead roles in some Hollywood movies(Male/Female). Along with
27 the main leading roles, we have multiple features in the data set that may help us predict the gender
28 of the leading role, e.g., the number of dialogues spoken by male actors. Female actors, the year in
29 which the movie was released, the revenue movie created, the age of lead roles and co-lead roles,
30 the number of male and female actors in the movie, etc. We will be analyzing this data and coming
31 up with some key features and insights from the data set, checking if we can predict the gender of a
32 leading role in any film.

33 2 Data Analysis Task

34 Q:1: Do men or women dominate speaking roles in Hollywood movies?

35 From sample data, we can see that over the years from 1939-2015, in most films, male actors have
36 dominated compared to females. Talking about the ratio, as seen in the graph, 76% of lead roles
37 have been covered by males, and only 24% is there for women. The training data also shows that the
38 ratio of male and female actors working in a film is very different. In the majority of the films, the
39 number of male actors working is very high compared to the females, thus giving male actors more
40 probability to be the lead role. So, we can see the clear dominance of the male actors.



41 Figure 2: Male and Female Roles (Designed and Created by group using Matplotlib)

43 Q:2: Has gender balance in speaking roles changed over time (i.e., years)?

44 If we look at the data from 2000-2015, nothing much has changed in these 15 years. Numbers are
45 still the same. Male speaking roles have been dominating (as shown in the graph below). In the year
46 2015, we can see some exceptions that the females were slightly dominant; to see how this goes
47 further, we'll need some more data to analyze further.

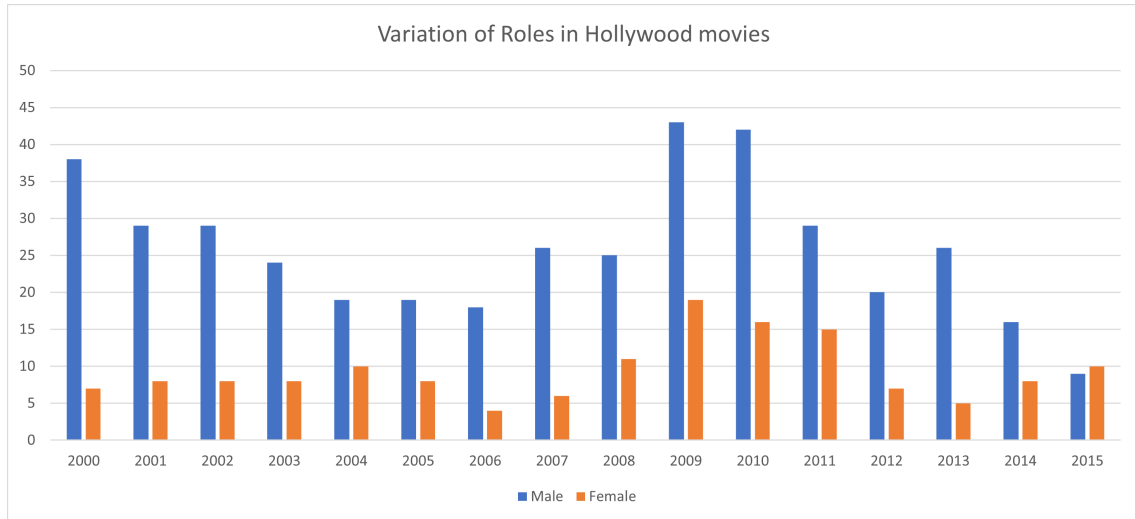


Figure 3: Male and Female gender balance (Designed and Created by group using Matplotlib)

Q:3: Do films in which men do more speaking make a lot more money than films in which women speak more?

Yes, the Below graph clearly shows that leading male films has done some really good business as compared to the leading female roles. Comparing the data from 1975 to 2015, we saw that there were clear and significant gaps between the revenue generated by leading male films and leading female films.

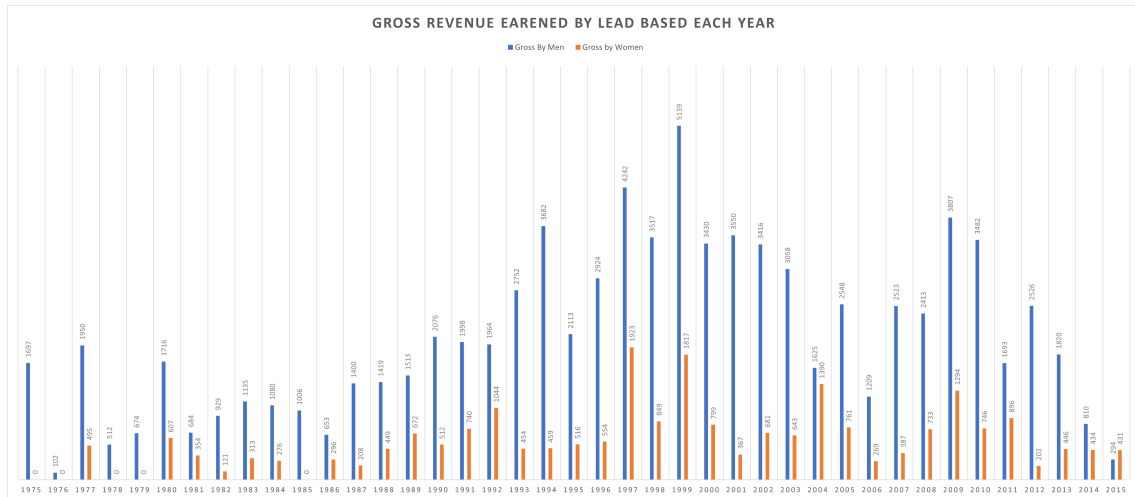


Figure 4: Role and Revenue (Designed and Created by group using Matplotlib)

3 Methodology

3.1 Logistic Regression

In Logistic Regression, we use the statistical analysis method on prior observation of the data set to get a binary outcome, like 1 or 0. The resultant prediction of logistic regression is a dependent data variable that we get after analyzing the relationship between one or more independent variables. ²

Application and Evaluation:

As said in the definition, Logistic regression is used to predict the outcomes as true/false or 0/1, so in our case as well, we are here to predict the lead role in the films as Male(1)/Female(0). First, we

²<https://www.techtarget.com/searchbusinessanalytics/definition/logistic-regression>

imported the whole Train.csv file into the jupyter notebook, analyzed the dataset, saw how many male and female leads there were, and analyzed the dataset based on the age of leads. Further, since the Logistic Regression model predicts 0 and 1, we remade the Lead column in the form of 0(female) and 1(male). Based on the evaluation, we dropped some columns which might not be useful in making predictions. We took dependent and independent variables for our model from the dataset. After that, From the sklearn library, we applied the train test split module to our data set to stimulate how our model is going to work, having a test size of 33% and a random state of 42. After that we used, we used the LogisticRegression method from sklearn linear model on our dataset, through which we got the following results:
 Train Accuracy = 85.9%
 Test Accuracy = 84.0%
 We also used the confusion metric to see how our model works. In order to tune our methods, we used hyperparameter regression with the following parameters penalty, C, solver, max iter, and using GridSearchCV module; we got the Accuracy = 87.9%.

3.2 Random Forest

Random Forest, which we use for classification or regression, is an supervised learning algorithm that implements an ensemble learning method that constitutes a large number of decision trees, the outputted result is the consensus of best answer to that problem. It is mostly used for classification and regression.
 Ensemble learning, which is based on collective opinion, consists of many machine learning algorithms that combine to produce a better result- the wisdom of the crowd. Just like if we say many people with lesser knowledge on a topic can produce a better result by helping each other instead of one person with more knowledge on something performing all in all.³

Application and Evaluation:

As done for the logistic Regression, we used the process here as well till using the sklearn library for applying the train test split module to our data set with the test size of 33% and a random state of 5. After that we used, we used the ensemble.RandomForestClassifier method with number of estimators = 100 from sklearn, through which we got the following results:
 Train Accuracy = 100%
 Test Accuracy = 82.2%
 As the model is overfitting, we tried to fix the model using the tuning method of hyper parameterization using the following parameters: estimators, max features, max depth, min samples split, min samples leaf, and bootstrap. After tuning the method with these parameters and then checking the accuracy of our model, we got the following results : Train Accuracy = 84.8%
 Test Accuracy = 79.0%

3.3 Naive Classifier

Naive Bayes is actually a group of classification algorithms that uses Bayes' Theorem. The Principal mainly followed by that every pair of classified features is independent of each other.

Bayes' Theorem computes the probability for an event occurring given that the probability of another which already occurred. Naive Bayes' theorem is stated below in mathematical form as the following equation: Where A and B are events and $P(B) \neq 0$ $P(A | B) = \frac{P(B|A)P(A)}{P(B)}$

P(A/B) posterior probability P(B/A) Likelihood P(A) Prior Probability P(B) Prior Probability that the Evidence is True

Multinomial Naive Bayes have a representation of frequencies using feature vectors which are used to produce certain events. This is the event model that is typically used for the purpose of document classification.

Bernoulli Naive Bayes have input features representation in the form of booleans variables that are independent of each other. Similar to the multinomial event model, this is also used for document

³<https://www.nvidia.com/en-us/glossary/data-science/random-forest/>

115 classification tasks, but it's more popular, where binary means yes or no in words example, we say
 116 either words exist or not. We use features in there instead of term frequencies.
 117 **In Gaussian, Naive Bayes** uses continuous values that are associated with features, and they are
 118 assumed to be distributed according to the Gaussian distribution. There is another name for this
 119 distribution Normal distribution. When we plot it, it shows a curve that appears bell-shaped, and it is
 120 symmetric about the mean of the feature values.

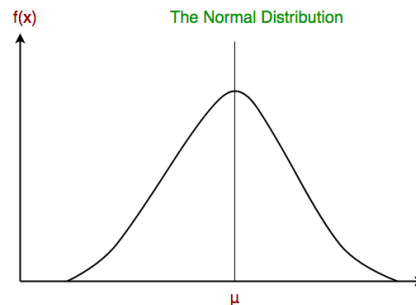


Figure 5: Gaussian Naive Bayes ⁴

123 Application and Evaluation:

124 As done in the other algorithms, we used the process here as well, using the sklearn library for
 125 applying the train test split module to our data set with the test size of 33% and a random state of
 126 5. After that we used, we used three classification algorithms from the Bayes theorem family of
 127 classification. We got the following results with algorithms.

129 Gaussian Naive Bayes:

130 Accuracy Score = 72%

131 Posterior Probability of POS Label 'Male' = 80%

133 Multinomial Naive Bayes:

134 Accuracy Score = 54%

135 Posterior Probability of POS Label 'Male' = 74%

137 Bernoulli Naive Bayes:

138 Accuracy Score = 77%

139 Posterior Probability of POS Label 'Male' = 87%

141 3.4 Method for Production

142 We will use it for production to decide on one model out of these two selected models; we carried out
 143 the cross-validation technique. The model which gives better accuracy with cross-validation will be
 144 the one to be taken in for production. So we carried out the K Fold Cross Validation technique for
 145 both models. By using the K Fold Cross Validation(KFold(10)) technique for Logistic Regression,
 146 we got 86% accuracy. But for Random forests, using the same technique, we got an accuracy of
 147 85.36%. Based on these evaluations we got from the cross-validation technique, we can see that
 148 Logistic Regression is ahead, and we can go with this method for production.

149 References

- 150 <https://www.nature.com/articles/s41599-020-0436-1/figures/1>
- 151 <https://www.techtarget.com/searchbusinessanalytics/definition/logistic-regression>
- 152 <https://www.nvidia.com/en-us/glossary/data-science/random-forest/>
- 153 <https://www.geeksforgeeks.org/naive-bayes-classifiers/>

⁴<https://www.geeksforgeeks.org/naive-bayes-classifiers/>

154 4 Appendix

155 Explanation of data

156 Year: That the film was released.
157 Number of female actors : With major speaking roles.
158 Number of male actors : With major speaking roles.
159 Gross : Profits made by film.
160 Total words : Total number of words spoken in the film.
161 Number of words male : Number of words spoken by all other male actors in the film (excluding
162 lead if lead is male)
163 Number of words female : Number of words spoken by all other female actors in the film (exclud-
164 ing lead if lead is female)
165 Number of words lead : Number of words spoken by lead.
166 Difference in words lead and co-lead : Difference in number of words by lead and the actor of
167 opposite gender who speaks most.
168 Lead Age : Age of lead actor.
169 Co-lead Age : Age of co-lead actor.
170 Mean Age Male : Mean age of all male characters.
171 Mean Age Female : Mean age of all female characters.

172 4.1 Code for Logistic Regression

```
173
174 #step 1 : Import Modules
175 from sklearn.datasets import make_classification
176 from matplotlib import pyplot as plt
177 from sklearn.linear_model import LogisticRegression
178 from sklearn.model_selection import train_test_split
179 from sklearn.metrics import confusion_matrix
180 import pandas as pd
181 import numpy as np
182 import seaborn as sns
183
184
185 #step 2 : Import Dataset
186 df = pd.read_csv("train.csv")
187
188
189 #step 3 : Viewing the data
190 df.head()
191
192 #Data Analysis
193 #Visually Analysing the data using Seaborn
194
195 #Check how many MALE and Female are there using countplot method
196
197 sns.countplot(x='Lead', data = df)
198
199 #null values in the dataset
200
201 df.isna().sum()
202
203
204 #Visualizing the null values
205 sns.heatmap(df.isna())
206
207 #data of age of leads where Lead is Male
208 df_Male = df.where(df['Lead'] == 'Male')
209 df_M=df_Male.dropna(how='all')
```

```

210
211 #data of age of leads where Lead is Female
212 df_Female = df.where(df['Lead'] == 'Female')
213 df_F=df_Female.dropna(how='all')
214
215 #find the distribution of age where Lead is Male
216 sns.displot(x='Age Lead', data= df_M)
217
218
219 #find the distribution of age where Lead is Female
220 sns.displot(x='Age Lead', data= df_F)
221
222
223 #Preparing data for model
224 #Convert Lead gender coloumn to binary values
225 pd.get_dummies(df['Lead'])
226
227 Gender_B = pd.get_dummies(df['Lead'], drop_first=True)
228 df['Gender_B'] = Gender_B
229
230 df.head()
231
232 #drop column which are not required
233 df.drop(['Year', 'Gross', 'Age Co-Lead'],axis=1,inplace=True)
234
235 df.head()
236
237 #Seperate dependent and dependent variable
238 #independent variable
239 x = df[['Total words', 'Number words male', 'Number words female', 'Number of words lead', 'Differenc
240 #dependent variable
241 y = df['Gender_B']
242 print(y)
243
244 #Data Modeling
245
246 from sklearn.model_selection import train_test_split
247
248 #train test split
249 x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.33, random_state=42)
250
251 import logistic regression
252 from sklearn.linear_model import LogisticRegression
253
254 #Fit Logistic Regression
255 lr = LogisticRegression()
256
257 lr.fit(x_train, y_train)
258
259 #check accuracy
260 print(f'Train Accuracy = {lr.score(x_train, y_train):.3f}')
261 print(f'Test Accuracy = {lr.score(x_test, y_test):.3f}')
262
263 #prediction
264 predict = lr.predict(x_test)
265
266 #Probabiltiy for +tive outcome is kept
267 lr_prob = lr.predict_proba(x_test)[: ,1]
268

```

```

269 #Complete AREA UNDER THE ROC CURVE values
270 from sklearn.metrics import roc_curve, roc_auc_score
271 lr_auc = roc_auc_score(y_test, lr_prob)
272
273 #Display the AREA UNDER THE ROC CURVE score
274 print("Logistic Regression : AUROC = %.3f" %(lr_auc))
275
276 #Calculate AREA UNDER THE ROC CURVE score
277 lr_fpr, lr_tpr, _ = roc_curve(y_test, lr_prob)
278
279 #plot the curve
280 plt.figure(figsize = (15,10))
281 plt.plot(lr_fpr, lr_tpr, linestyle = '--', label = 'Random forest (AUROC = %0.3f)' %lr_auc)
282
283 plt.title('ROC Plot')
284 plt.xlabel('False Positive Rate')
285 plt.ylabel('True Positive Rate')
286
287 plt.legend()
288 plt.show()
289
290 #print Confusion Matrix to see how well your model works
291 from sklearn.metrics import confusion_matrix
292
293
294 pd.DataFrame(confusion_matrix(y_test,predict), columns = ['Predicted No', 'Predicted Yes'], index = ['Actual No', 'Actual Yes'])
295
296
297 from sklearn.metrics import classification_report
298 print(classification_report(y_test, predict))
299
300
301 #Tunning the method using hyperparameter regression
302
303 logModel = LogisticRegression()
304
305 param_grid = [
306     {'penalty' : ['l1','l2','elasticnet','none'],
307      'C' : np.logspace(-4, 4, 20),
308      'solver' : ['lbfgs','newton-cg','liblinear','sag','saga'],
309      'max_iter' : [100, 1000, 2500, 5000]
310     }
311 ]
312
313
314
315 from sklearn.model_selection import GridSearchCV
316 clf = GridSearchCV(logModel, param_grid = param_grid , cv = 3, verbose = True, n_jobs = -1)
317
318 best_clf = clf.fit(x,y)
319 ]:
320
321 best_clf.best_estimator_
322
323
324 #check accuracy
325 print(f'Accuracy = {best_clf.score(x,y):.3f}')
326
327

```



```

328 from sklearn.model_selection import KFold
329 model = LogisticRegression()
330 kfold_validation = KFold(10)
331
332
333 from sklearn.model_selection import cross_val_score
334 result = cross_val_score(model, x, y, cv = kfold_validation)
335 print(result)
336 print(np.mean(result))
337
338 from sklearn.model_selection import ShuffleSplit
339 model = LogisticRegression()
340 ssplit = ShuffleSplit(n_splits = 10, test_size = 0.30)
341 results = cross_val_score(model, x, y, cv=ssplit)
342
343 results
344 print(f'Results - {results}')
345 print(f'Mean Result - {np.mean(results):.5f}')

```

346 **4.2 Code for Random Forest Regression**

```

347 #importing the modules
348 from sklearn.datasets import make_classification
349 from matplotlib import pyplot as plt
350 from sklearn.linear_model import LogisticRegression
351 from sklearn.model_selection import train_test_split
352 from sklearn.metrics import confusion_matrix
353 import pandas as pd
354 import numpy as np
355 import seaborn as sns
356
357 #reading the dataset
358 df = pd.read_csv("train.csv")
359 step 3 : Viewing the data
360 df.head()
361
362
363 ##Data Analysis
364 ##Visually Analysisng the data using Seaborn
365 ##Check how many MALE and Female are there using countplot method
366
367 sns.countplot(x='Lead', data = df)
368 #Checking for null values in the dataset
369
370 df.isna().sum()
371
372 #Visualizing the null values in the dataset
373
374
375 sns.heatmap(df.isnull())
376
377 #data of age of leads where Lead is Male
378 df_Male = df.where(df['Lead'] == 'Male')
379 df_M=df_Male.dropna(how='all')
380
381 #data of age of leads where Lead is Female
382 df_Female = df.where(df['Lead'] == 'Female')
383 df_F=df_Female.dropna(how='all')
384

```

```

385 #find the distribution of age where Lead is Male
386 sns.displot(x='Age Lead', data= df_M)
387
388 #find the distribution of age where Lead is Female
389 sns.displot(x='Age Lead', data= df_F)
390
391
392 #Preparing data for model
393 #Convert Lead gender coloumn to Numerical values
394 pd.get_dummies(df['Lead'])
395
396
397 Gender_B = pd.get_dummies(df['Lead'], drop_first=True)
398
399
400 df['Gender_B'] = Gender_B
401
402
403 df.head()
404
405
406 #drop column which are not required for the model
407 df.drop(['Year', 'Gross', 'Age Co-Lead'], axis=1, inplace=True)
408
409 df.head()
410
411 #Seperate dependent and independent variable
412 #independent variable
413 x = df[['Total words', 'Number words male', 'Number words female', 'Number of words lead', 'Difference']]
414 #dependent variable
415 y = df['Gender_B']
416
417 print(y)
418
419
420 #Data Modeling
421
422 from sklearn.model_selection import train_test_split
423
424 #train test split
425 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.33, random_state=5)
426
427 import Random Forest
428 from sklearn import ensemble
429 rf_clf = ensemble.RandomForestClassifier(n_estimators = 100)
430 rf_clf.fit(x_train, y_train)
431
432
433 #check accuracy
434 print(f'Train Accuracy = {rf_clf.score(x_train, y_train):.3f}')
435 print(f'Test Accuracy = {rf_clf.score(x_test, y_test):.3f}')
436
437 #Predict Probabiltiy
438 rf_prob = rf_clf.predict_proba(x_test)
439
440 #Probabiltiy for +tive outcome is kept
441 rf_prob = rf_clf.predict_proba(x_test)[: ,1]
442
443 #Complete AREA UNDER THE ROC CURVE values

```

```

444 from sklearn.metrics import roc_curve, roc_auc_score
445 rf_auc = roc_auc_score(y_test, rf_prob)
446
447 #Display the AREA UNDER THE ROC CURVE score
448 print("Random forest : AUROC = %.3f" %(rf_auc))
449
450 #Calculate AREA UNDER THE ROC CURVE score
451 rf_fpr, rf_tpr, _ = roc_curve(y_test, rf_prob)
452
453 #plot the curve
454 plt.figure(figsize = (15,10))
455 plt.plot(rf_fpr, rf_tpr, linestyle = '--', label = 'Random forest (AUROC = %0.3f)' %rf_auc)
456
457 plt.title('ROC Plot')
458 plt.xlabel('False Positive Rate')
459 plt.ylabel('True Positive Rate')
460
461 plt.legend()
462 plt.show()
463
464 #No. of trees
465 n_estimators, = [int(x) for x in np.linspace(start = 10, stop = 80, num = 10)]
466 #No. of features
467 max_features = ['auto', 'sqrt']
468 #Max number of levels in tree
469 max_depth = [2,4]
470 #Mini no. of samples to split node
471 min_samples_split = [2, 5]
472 #Mini no. of samples required at each leaf node
473 min_samples_leaf = [1, 2]
474 #Method of selecting samples for training each tree
475 bootstrap = [True, False]
476
477 #Create the param grid
478 param_grid = {'n_estimators': n_estimators,
479               'max_features': max_features,
480               'max_depth': max_depth,
481               'min_samples_split': min_samples_split,
482               'min_samples_leaf': min_samples_leaf,
483               'bootstrap': bootstrap}
484
485
486 rf_Model = ensemble.RandomForestClassifier()
487
488
489 from sklearn.model_selection import GridSearchCV
490 rf_Grid = GridSearchCV(estimator = rf_Model, param_grid = param_grid , cv = 3, verbose = 2, n_jobs=
491
492
493 rf_Grid.fit(x_train,y_train)
494
495
496 rf_Grid.best_params_
497
498 #check accuracy
499 print(f'Train Accuracy = {rf_Grid.score(x_train, y_train):.3f}')
500 print(f'Test Accuracy = {rf_Grid.score(x_test, y_test):.3f}')
501
502 from sklearn.model_selection import KFold

```

```

503 model = ensemble.RandomForestClassifier()
504 kfold_validation = KFold(10)
505
506
507 from sklearn.model_selection import cross_val_score
508 result = cross_val_score(model, x, y, cv = kfold_validation)
509 print(result)
510 print(np.mean(result))
511
512
513
514 from sklearn.model_selection import ShuffleSplit
515 model = ensemble.RandomForestClassifier()
516 ssplit = ShuffleSplit(n_splits = 10, test_size = 0.30)
517 results = cross_val_score(model, x, y, cv=ssplit)
518
519 print(f'Results = {results}')
520 print(f'Mean Result = {np.mean(results):.5f}')
521

```

522 4.3 Code for Naive Classifier

```

523 #step 1 : Import Modules
524 import numpy as np
525 import pandas as pd
526 from sklearn.model_selection import train_test_split
527 from sklearn.naive_bayes import GaussianNB
528 from sklearn.naive_bayes import BernoulliNB
529 from sklearn.naive_bayes import MultinomialNB
530 from sklearn.metrics import accuracy_score
531 from sklearn.metrics import confusion_matrix, f1_score
532
533 #step 2 : Import Dataset
534 df = pd.read_csv("train.csv")
535 # checking within the columns if we have missing values
536 df.info()
537
538 #step 3 : Viewing the data
539 df.head(5)
540
541 #step 4 : Separating the Labels for our continues other attributes
542 x = df.drop('Lead', axis=1)
543 y = df['Lead']
544
545 #step 5 : Splitting the training & test data
546 x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.25,random_state=42)
547
548
549 # step 6 : Using Bernoulli Naive Bayes for score and also showing accuracy
550 BernNB = BernoulliNB(binarize= 0.1)
551 BernNB.fit(x_train, y_train)
552 print(BernNB)
553
554 y_expect = y_test
555 y_pred = BernNB.predict(x_test)
556 print("Accuracy Score: ",accuracy_score(y_expect, y_pred))
557 print("POS Label Male: ", f1_score(y_expect, y_pred, average="binary", pos_label="Male"))
558
559 # step 7 : Using Multinomial Naive Bayes for score and also showing accuracy

```

```

560 MultiNB = MultinomialNB()
561 MultiNB.fit(x_train, y_train)
562 print(MultiNB)
563 y_pred = MultiNB.predict(x_test)
564 print("Accuracy Score: ",accuracy_score(y_expect, y_pred))
565 print("POS Label Male: ", f1_score(y_expect, y_pred,average="binary", pos_label="Male"))
566
567 # step 8 : Using Gaussian Naive Bayes for score and also showing accuracy
568 GausNB = GaussianNB()
569 GausNB.fit(x_train, y_train)
570 print(GausNB)
571 y_pred = GausNB.predict(x_test)
572 print("Accuracy Score: ",accuracy_score(y_expect, y_pred))
573 print("POS Label Male: ", f1_score(y_expect, y_pred, average="binary", pos_label="Male"))
574
575

```