

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April
Tuesday

22

12

Week 16 Day (102-263)

09.00

CLASSIFICATION

10.00

- SIMPLE CLASSIFICATION - BINARY CLASSIFICATION
- CLASSIFY IN MORE THAN 2 CLASSES

11.00

EXAMPLE - SPAM OR NOT SPAM

12.00

FRAUD TRANSACTION OR NOT

EYE COLOR (BLUE, GREEN, BROWN)

13.00

REGRESSION VS CLASSIFICATION

14.00

- REGRESSION - RESPONSE VARIABLE IS QUANTITATIVE
- CLASSIFICATION - RESPONSE VARIABLE IS CATEGORICAL

15.00

OR QUALITATIVE.

16.00

PERFORMING CLASSIFICATION WITH LOGISTIC REGRESSION

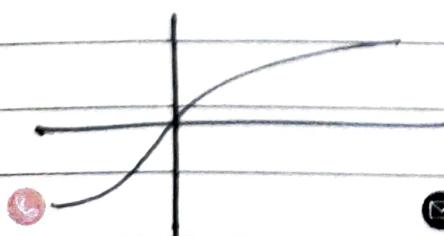
17.00

- DURING LOGISTIC REGRESSION WE WANT TO DETERMINE THE PROBABILITY OF AN OBSERVATION TO BE PART OF A CLASS OR NOT.

18.00

- OUTPUT BETWEEN 0 AND 1 (1 MEANS VERY LIKELY)

THE FUNCTION TO DO THAT - SIGMOID FUNCTION



$$P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$



* HERE WE ARE ASSUMING ONLY ONE PREDICTOR 'X'.

13

'22

April

Wednesday

Week 16 Day (103-262)

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
				1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

09.00

$$\beta_0 + \beta_1 x$$

$$\underline{p(x)} = e$$

10.00

$$1 - p(x)$$

take log both side

11.00

$$\log \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x_1$$

12.00

$$\log \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

13.00

In Python, Logistic Regression can be performed Like this -

14.00

15.00

`logistic-reg = LogisticRegression()`

16.00

`logistic-reg.fit(X-train, y-train.ravel())`

17.00

`y-prob = logistic-reg.predict_proba
(X-test)[:, 1]`

18.00

`y-pred = np.where(y-prob > 0.5, 1, 0)`

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31											

April

22

Thursday

14

Week 16 Day (104-261)

09.00 LINEAR DISCRIMINANT ANALYSIS (LDA)

10.00 CAVEATS (WARNING) OF LOGISTIC REGRESSION

- WHEN CLASSES ARE WELL SEPARATED, PARAMETERS OF LOGISTIC REGRESSION ARE UNSTABLE.
- UNSTABLE FOR SMALL DATASETS
- NOT THE BEST TO PREDICT MORE THAN 2 CLASSES.
- LOGISTIC REGRESSION CAN ONLY BE USED FOR BINARY CLASSIFICATION.

14.00 ↳ BUT WITH LDA WE CAN OVERCOME THESE ISSUES BECAUSE

- MODELS THE DISTRIBUTION PREDICTORS SEPARATELY FOR EACH CLASS.
- USE BAYES'S THEOREM TO ESTIMATE THE PROBABILITY

16.00 BAYES THEOREM FOR CLASSIFICATION

17.00 SUPPOSE WE WANT

TO CLASSIFY AN

18.00 OBSERVATION INTO

ONE OF 'K' CLASS

Let $K \geq 2$, then

- $\prod_k f_k(x)$ is overall probability

- f_k is density function

if $X=x$, f_k is large

$$P(Y=k | X=x)$$

$$= \frac{\prod_k f_k(x)}{\sum_{l=1}^K \prod_l f_l(x)}$$

15

'22

April

Week 16 Day (105-260)

Friday

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10				
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

09.00

CHALLENGE IS APPROXIMATING THE DENSITY FUNCTION

10.00

→ LDA FOR ONE PREDICTOR

11.00

$$f_K(x) = \frac{1}{\sqrt{2\pi}\sigma_K} \exp\left(-\frac{(x-\mu_K)^2}{2\sigma^2}\right)$$

12.00

if we plot the upper function in the previous formula,
and then take log, we get —

14.00

$$\delta_K(x) = \underbrace{x\mu_K}_{\text{discriminant}} - \frac{\mu_K^2}{\sigma^2} + \log(\pi_K) - \frac{1}{2\sigma^2}$$

15.00

WHEN WE APPLY LDA, WE NEED TO BE AWARE OF
ASSUMPTIONS —

- EACH CLASS IS DRAWN FROM A GAUSSIAN DISTRIBUTION
- EACH CLASS HAS ITS OWN MEAN
- ASSUME A COMMON VARIANCE
- ASSUME A COMMON COVARIANCE MATRIX

LDA IN PYTHON →

(ON NEXT PAGE) *



MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN
			1	2	3	4
9	10	11	12	13	14	15
23	24	25	26	27	28	29

April

'22

16

Saturday

Week 16 Day (106-259)

09.00

LOA IN PYTHON

10.00

lدا = Linear Discriminant Analysis()

11.00

lدا =
lدا.fit(x-train, y-train.ravel())

12.00

y-prob-lدا = lدا.predict_proba
(x-test)[:, 1]

13.00

y-pred-lدا = np.where(y-prob-lدا
> 0.5, 1, 0)

14.00

15.00

16.00

QUADRATIC DISCRIMINANT ANALYSIS (QDA)

17.00

ASSUMPTIONS OF QDA

Sunday 17

- EACH CLASS IS DRAWN FROM A MULTIVARIATE GAUSSIAN

18.00

DISTRIBUTION

- EACH CLASS HAS ITS OWN VECTOR (MEAN VECTOR)
- EACH CLASS HAS ITS OWN COVARIANCE MATRIX

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT
11	12	13	14	15	16	17	18	19	20	21	22	23
25	26	27	28	29	30							

18

'22

April

Monday

Week 17 Day (108-257)

09.00

DISCRIMINANT WITH QDA

10.00

$$\delta(x) = -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log(\pi_k)$$

11.00

ADVANTAGE OF QDA

12.00

- BETTER FOR LARGE DATASET

- LOWER BIAS, HIGHER VARIANCE

13.00

QDA IN PYTHON

14.00

qda = QuadraticDiscriminantAnalysis()

15.00

qda.fit(X_train, y_train.ravel())

16.00

y_prob_qda = qda.predict_proba(X_test)
[:, 1]

17.00

y_pred_qda = np.where(y_prob_qda
0.5, 1, 0)

18.00



MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April
19

Tuesday

Week 17 Day (109-256)

09.00

CLASSIFICATION

HOW TO ASSESS THE PERFORMANCE OF A MODEL

10.00

- SENSITIVITY - TRUE POSITIVE RATE. PROPORTION OF ACTUAL POSITIVE IDENTIFIED.

11.00

EX- PROPORTION OF FRAUDENT TRANSACTION THAT ARE ACTUALLY FRAUDULENT.

12.00

- SPECIFICITY - TRUE NEGATIVE RATE. PROPORTION OF ACTUAL NEGATIVES IDENTIFIED.

14.00

EX- PROPORTION OF NON- FRAUDULENT TRANSACTIONS THAT ARE ACTUALLY NON- FRAUDULENT.

15.00

WE CAN ALSO USE ROC CURVE (RECIEVER OPERATING CHARACTERISTIC)

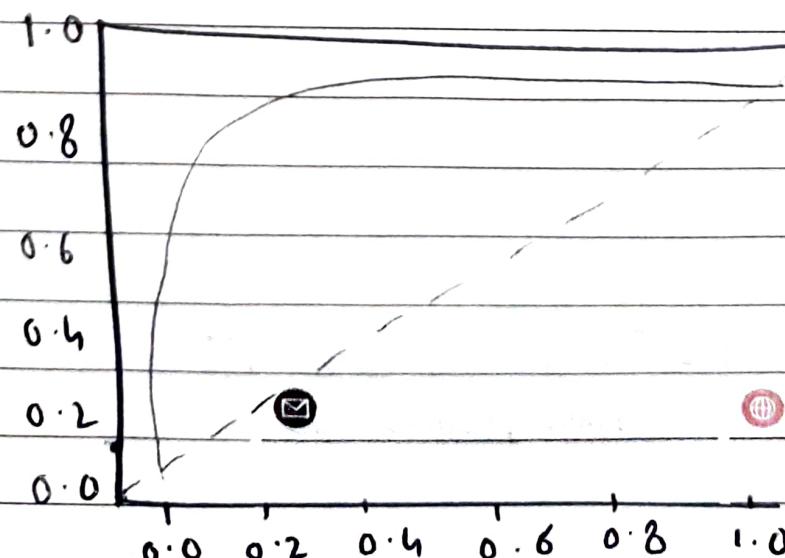
16.00

- TAKE AREA UNDER THE CURVE (AUC)
- THE CLOSER TO 1, THE BETTER IT IS

17.00

ROC Curve

18.00

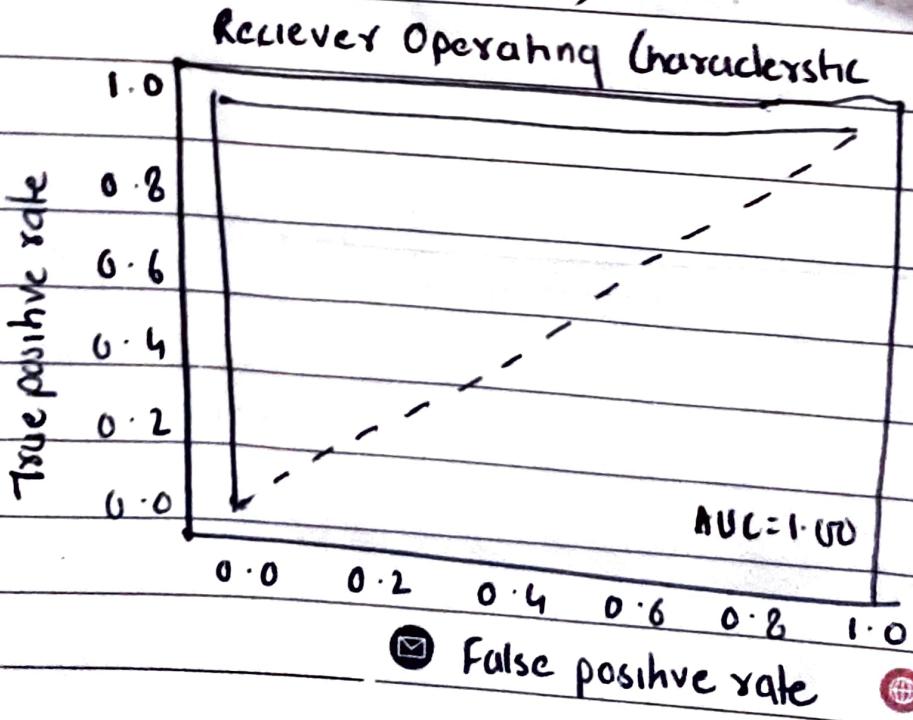


09.00

ROC IN PYTHON

```
10.00 def plot_roc(rocauc):
11.00     plt.figure(figsize=(7,7))
12.00     plt.title('Receiver Operating Characteristic')
13.00     plt.plot(false_positive_rate,
14.00             true_positive_rate,
15.00             color='red', label='AUC = %.2f' % rocauc)
16.00     plt.legend(loc='lower right')
17.00     plt.plot([0,1], [0,1], linestyle='--')
18.00     plt.axis('tight')
19.00     plt.ylabel('True positive rate')
20.00     plt.xlabel('False positive rate')
```

ROC IN PYTHON (OUTPUT)



09.00

- HANDS ON IN PYTHON -

10.00

CLASSIFICATION MODEL # CODE DANIEL

11.00

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
from sklearn.preprocessing import  
LabelEncoder  
from sklearn.model_selection import  
train_test_split, cross_val_score  
from sklearn.metrics import roc_curve,  
auc, confusion_matrix  
%matplotlib inline
```

16.00

17.00

pandas - Library for data manipulation and analysis.

function for importing & exporting data from various file formats.

numpy - functions for working with arrays and matrices.

matplotlib - plotting library for creating interactive visualization

seaborn - Library built on top of matplotlib for visualization

sklearn - Library for Machine Learning



22

22

April

Week 17 Day (112-253)

Friday

SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3
11	12	13	14	15	16	17	18
18	19	20	21	22	23	24	
25	26	27	28	29	30		

- 09.00 LabelEncoder - class from sklearn preprocessing that transforms categorical data into numerical data by assigning unique integers to each category.
- 10.00 train-test-split - function from sklearn.model_selection that splits class into training and testing sets.
- 11.00 cross_val_score - perform cross validation, a technique of model by splitting dataset into multiple training & testing sets.
- 12.00 roc_curve and auc - function that calculates the operating characteristic (ROC) curve and the area under curve (AUC), which is used to evaluate performance of binary classifiers.
- 13.00 confusion_matrix - computes confusion matrix.

16.00
17.00
18.00

```
datapath = 'url/path'  
data = pd.read_csv(datapath)
```

↳ additionally for Seeing

data.head()

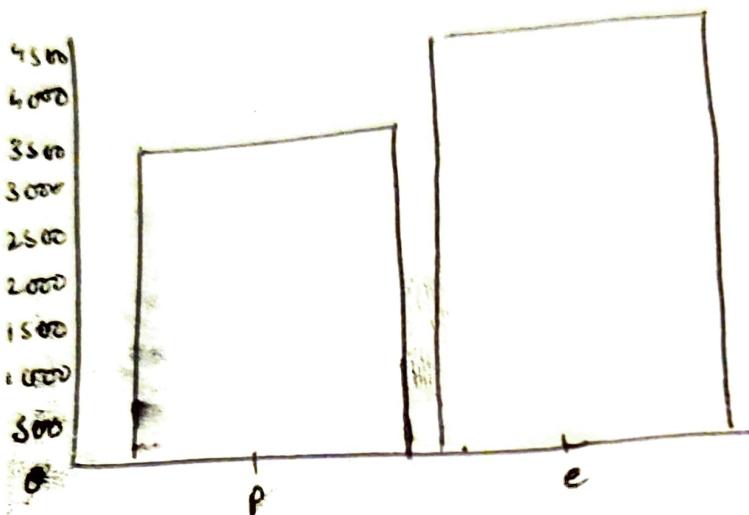
09.00

`X = data['class']``ax = sns.countplot(x=X, data=data)`

10.00

— 1

11.00



12.00

13.00

14.00

15.00

— 2

1) extracts the column of data with the label 'class' from pandas dataframe data and assign it to variable X.

2) it creates a countplot using the Seaborn Library Sunday 24 which is barplot that shows the count of each unique value

X parameter specifies the categorical variable to be plotted which in this case X variable was previously extracted from the data dataframe

data parameter specifies that dataframe that the categorical variable is coming from , which is also data in this case.

resulting plot is stored in variable ax, which can be used



to further



customize the plot.



09.00

```
def plot_data(hue, data):  
    for i, col in enumerate(data.columns):  
        plt.figure(i)  
        ax = sns.countplot(x=data[col],  
                            hue=hue,  
                            data=data)
```

10.0

11.0

12.0

13.0

14.0

15.0

16.00

plot_data takes 2 parameters -

hue - string specifying the name of column to be used for

coloring the countplot. This is the column that contains the categorical variable that will be used to group data and

color the bars in the plot.

data - pandas dataframe contains data to be plotted

data dataframe - create separate countplot with specified hue parameter.

enumerate function - get index of each column into data dataframe which is used to create separate figure for each

plot. (i.e. plt.figure(i))

sns.countplot is used to create countplot, with x parameter set to the name of column (col) and the hue parameter set to the specified hue parameter.

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15
23	24	25	26	27	28	29
20	21	22	23	24	25	31

April

22

Tuesday

Week 18 Day (116-249)

26

09.00

Resulting plot for each column is stored in variable ax, but it is not returned by the function. If function is called, it will simply create series of countplot in the dataframe with specified hue parameter

11.00

12.00

hue = data['class']

—1

data-to-plot = data.drop('class', axis=1)

—2

plot-data(hue, data-to-plot)

—3

14.00

L01
L02
L03
L04

15.00

16.00

17.00

18.00 1) The code creates new variable named hue that extracts the column of data with Label 'class' from pandas dataframe data and assigns it to hue. This column is used for coloring the countplots.

2) New dataframe named data-to-plot drops column with Label 'class' from the data dataframe and assigns it to data-to-plot. This new dataframe contains all columns from original data dataframe except 'class' column.

3) This will create series of countplots, one for each column in data-to-plot with 'class' column used for coloring the bars in plot.

Wednesday

09.00

PREPROCESSING -

10.00

```
for col in data.columns:  
    print(f'{col}: {data[col].  
        isnull().sum()}' )
```

11.00

12.00

13.00 Code loops through each column in the pandas dataframe data and for each column it prints out the number of missing value using the isnull () method to find missing values and sum () method to count them.

14.00

```
le = LabelEncoder()  
data['class'] = le.fit_transform  
(data['class'])  
data.head()
```

15.00

— 1

— 2

- 1) Create an instance of LabelEncoder class & assign it to variable le
- 2) Then use fit-transform method of LabelEncoder object to encode the value in the 'class' column of pandas dataframe data. fit_transform method fits the encoder to data and transform the data in single step.
It replaces values in 'class' with encoded integer value.

09.00

MODELLING —

10.00

`y = data['class'].values.reshape(-1, 1)` — 1

11.00

`X = encoded_data.drop(['class'], axis=1)` — 2

12.00

`X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)` — 3

13.00

14.00

15.00

16.00

1) first creates a variable by selecting the "class" column from pandas dataframe using indexing, and then reshaping it using the reshape method of numpy array. The reshape method is used to convert 1D array to 2D array with single column. Because scikit learn expects the target variable to be in 2D array format.

2) creates a variable X by dropping class column from pandas dataframe encoded-data using drop method with axis parameter set to 1. new dataframe without 'class' column.

17.00

18.00

19.00

20.00

21.00

22.00

23.00

24.00

25.00

26.00

27.00

28.00

29.00

30.00

31.00

32.00

33.00

34.00

35.00

36.00

37.00

38.00

39.00

40.00

41.00

42.00

43.00

44.00

45.00

46.00

47.00

48.00

49.00

50.00

51.00

52.00

53.00

54.00

55.00

56.00

57.00

58.00

59.00

60.00

61.00

62.00

63.00

64.00

65.00

66.00

67.00

68.00

69.00

70.00

71.00

72.00

73.00

74.00

75.00

76.00

77.00

78.00

79.00

80.00

81.00

82.00

83.00

84.00

85.00

86.00

87.00

88.00

89.00

90.00

91.00

92.00

93.00

94.00

95.00

96.00

97.00

98.00

99.00

100.00

101.00

102.00

103.00

104.00

105.00

106.00

107.00

108.00

109.00

110.00

111.00

112.00

113.00

114.00

115.00

116.00

117.00

118.00

119.00

120.00

121.00

122.00

123.00

124.00

125.00

126.00

127.00

128.00

129.00

130.00

131.00

132.00

133.00

134.00

135.00

136.00

137.00

138.00

139.00

140.00

141.00

142.00

143.00

144.00

145.00

146.00

147.00

148.00

149.00

150.00

151.00

152.00

153.00

154.00

155.00

156.00

157.00

158.00

159.00

160.00

161.00

162.00

163.00

164.00

165.00

166.00

167.00

168.00

169.00

170.00

171.00

172.00

173.00

174.00

175.00

176.00

177.00

178.00

179.00

180.00

181.00

182.00

183.00

184.00

185.00

186.00

187.00

188.00

189.00

190.00

191.00

192.00

193.00

194.00

195.00

196.00

197.00

198.00

199.00

200.00

201.00

202.00

203.00

204.00

205.00

206.00

207.00

208.00

209.00

210.00

211.00

212.00

213.00

214.00

215.00

216.00

217.00

218.00

219.00

220.00

221.00

222.00

223.00

224.00

225.00

226.00

227.00

228.00

229.00

230.00

231.00

232.00

233.00

234.00

235.00

236.00

237.00

238.00

239.00

240.00

241.00

242.00

243.00

244.00

245.00

246.00

247.00

248.00

249.00

250.00

251.00

252.00

253.00

254.00

255.00

256.00

257.00

258.00

259.00

260.00

261.00

262.00

263.00

264.00

265.00

266.00

267.00

268.00

269.00

270.00

271.00

272.00

273.00

274.00

275.00

276.00

277.00

278.00

279.00

280.00

281.00

282.00

283.00

284.00

285.00

286.00

287.00

288.00

289.00

290.00

291.00

292.00

293.00

294.00

295.00

296.00

297.00

298.00

299.00

300.00

301.00

302.00

303.00

304.00

305.00

306.00

307.00

308.00

309.00

310.00

311.00

312.00

313.00

314.00

315.00

316.00

317.00

318.00

319.00

320.00

321.00

322.00

323.00

324.00

325.00

326.00

327.00

328.00

329.00

330.00

331.00

332.00

09.00

LOGISTIC REGRESSION -

10.00

From sklearn.linear_model import
LogisticRegression

→ 1

11.00

logistic-reg = LogisticRegression()
logistic-reg.fit(X-train, y-train.ravel())

→ 2

12.00

y-prob = logistic-reg.predict_proba
(X-test)[:,1]

→ 3

13.00

y-pred = np.where(y-prob > 0.5, 1, 0)

→ 4

14.00

log-confusion-matrix = confusion_matrix
(y-test, y-pred)

]

15.00

log-confusion-matrix

5

16.00

1) It provides variety of linear models for use in regression
and classification

17.00

2) Create instance of LogisticRegression class and assign it to
variable logistic-reg , ravel method flatten 1D array

18.00

3) code uses predict_proba to predict probability of positive
class , predict_proba returns 2D numpy array where each
row corresponds to class. class1 → 0 , class2 → 1

4) np.where function convert the predicted probabilities in
y-prob to binary predictions (y-pred) by applying threshold



5) it calculates confusion matrix for LR (Logistic Regression)

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
6	7	8	9	10	11	12	13	14	15	16	17	18	19
23	24	25	26	27	28	29	30	31					

April 22
Saturday
Week 18 Day (120-245)

09:00

false-positive-rate, true-positive-rate, threshold = ROC-curve (y-test, y-prob)
 $\text{ROC-AUC} = \text{auc}(\text{false-positive-rate}, \text{true-positive-rate})$
 ROC-AUC

12:00

CODE CALCULATE RECEIVING OPERATING CHARACTERISTIC (ROC) curve for Logistic Regression Model.

13:00

ROC curve is a graphical representation of tradeoff b/w TPR & FPR. for different threshold of binary classifier.

14:00

```
def plot_roc(roc_auc):
    plt.figure(figsize=(7,7))
    plt.title('Receiver Operating Characteristic')
    plt.plot(false_positive_rate, true_positive_rate, color='red', label='AUC = %0.2f' % roc_auc)
    plt.legend(loc='lower-right')
    plt.plot([0,1], [0,1], linestyle='--')
    plt.axis('tight')
    plt.ylabel('True Positive Rate')
    plt.xlabel('False Positive Rate')
    plot_roc(roc_auc)
```

function used
to define
properties of
GRAPH

Same will be applied to LDA & QDA
with slight changes.