



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

Performance of Logistic Model-III

Dr A. RAMESH

DEPARTMENT OF MANAGEMENT STUDIES



Agenda

Python demo for accuracy prediction in logistic regression model using Receiver operating characteristics curve

Sensitivity and Specificity

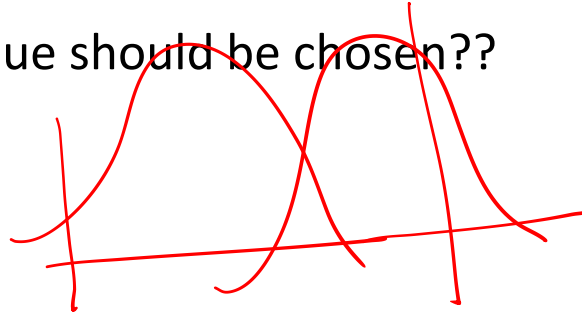
- For checking, what type of error we are making; we use two parameters-
 1. Sensitivity = $tp/(tp+fn)$ \longrightarrow True Positive Rate(tpr)
 2. Specificity = $tn/(tn+fp)$ \longrightarrow True Negative Rate (tnr)

Specificity and Sensitivity Relationship with Threshold

Threshold (Lower)	Sensitivity (\uparrow)	Specificity (\downarrow)
Threshold (Higher)	Sensitivity (\downarrow)	Specificity (\uparrow)



Which threshold value should be chosen??



Measuring Accuracy, Specificity and Sensitivity

```
In [20]: 1 Accuracy = (tp + tn) / (tp + tn + fp + fn)
          2 print("Accuracy {:.2f}".format(Accuracy))
```

Accuracy 0.76

```
In [21]: 1 Specificity = tn/(tn+fp)
          2 print("Specificity {:.2f}".format(Specificity))
```

Specificity 0.94

```
In [22]: 1 Sensitivity = tp/(tp+fn)
          2 print("Sensitivity {:.2f}".format(Sensitivity))
```

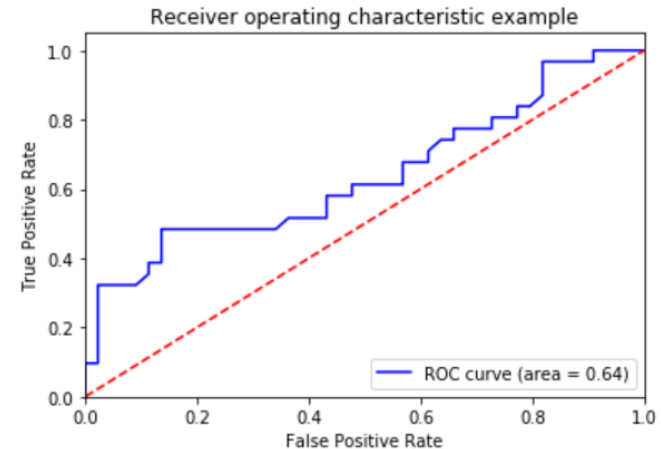
Sensitivity 0.44

$t = 0.5$

ROC Curve for Training dataset

```
In [23]: 1 from sklearn.metrics import roc_auc_score
2 from sklearn.metrics import roc_curve, auc
3 log_ROC_AUC1 = roc_auc_score(y_train, y_predict_train)
4 fpr1, tpr1, thresholds1 = roc_curve(y_train, y_prob_train)
5 roc_auc1 = auc(fpr1, tpr1)
```

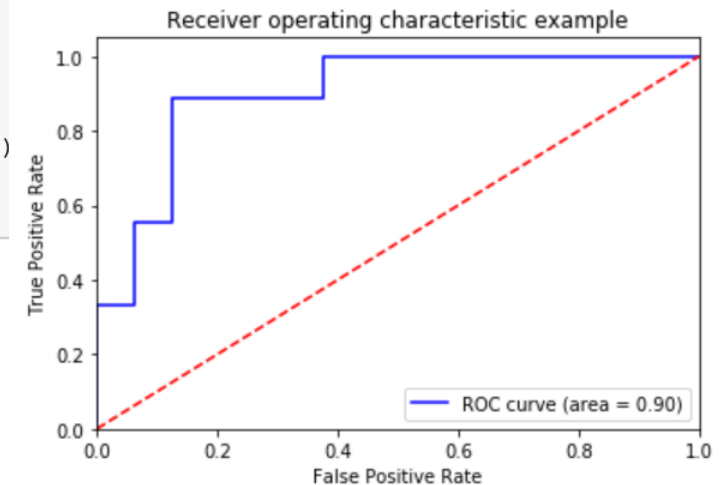
```
In [24]: 1 plt.figure()
2 plt.plot(fpr1, tpr1, color='blue', label='ROC curve (area = %0.2f)' % roc_auc1)
3 plt.plot([0, 1], [0, 1], 'r--')
4 plt.xlim([0.0, 1.0])
5 plt.ylim([0.0, 1.05])
6 plt.xlabel('False Positive Rate')
7 plt.ylabel('True Positive Rate')
8 plt.title('Receiver operating characteristic example')
9 plt.legend(loc="lower right")
10 plt.show()
```



ROC Curve for Test data set

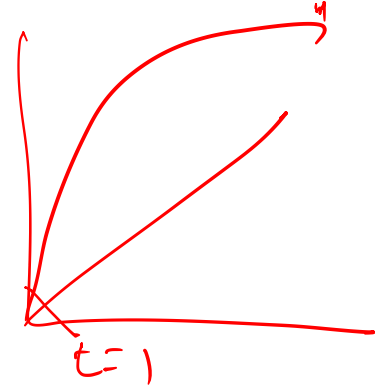
```
In [25]: 1 log_ROC_AUC = roc_auc_score(y_test, y_predict)
          2 fpr, tpr, thresholds = roc_curve(y_test, y_prob)
          3 roc_auc = auc(fpr, tpr)
```

```
In [26]: 1 plt.figure()
          2 plt.plot(fpr, tpr, color='blue', label='ROC curve (area = %0.2f)' % roc_auc)
          3 plt.plot([0, 1], [0, 1], 'r--')
          4 plt.xlim([0.0, 1.0])
          5 plt.ylim([0.0, 1.05])
          6 plt.xlabel('False Positive Rate')
          7 plt.ylabel('True Positive Rate')
          8 plt.title('Receiver operating characteristic example')
          9 plt.legend(loc="lower right")
         10 plt.show()
```



Threshold value selection

- The outcome of logistic regression model is a probability.
- Selecting a good threshold value is often challenging.
- Threshold values on ROC curve –



Threshold = 1	TPR = 0	FPR = 0
Threshold = 0	TPR = 1	FPR = 1

- Threshold values are often selected based on which errors are better.

Accuracy checking for different threshold values

```
In [27]: 1 from sklearn.preprocessing import binarize
          2 y_predict_class1 = binarize(y_prob.reshape(1, -1), 0.35)[0]
          3 y_predict_class1
```

```
Out[27]: array([1., 1., 1., 0., 0., 0., 0., 1., 1., 1., 0., 0., 1., 1., 1., 1., 1.,
                1., 0., 1., 1., 1., 1., 1., 0.])
```

```
In [28]: 1 #converting the array from float data type to integer data type
          2 y_predict_class1 = y_predict_class1.astype(int)
          3 y_predict_class1
```

```
Out[28]: array([1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1,
                1, 1, 0])
```

```
In [29]: 1 confusion_matrix_2 = confusion_matrix(y_test, y_predict_class1)
          2 print(confusion_matrix_2)
```

```
[[8 8]
 [0 9]]
```

Accuracy checking for different threshold values

```
In [30]: 1 tn, fp, fn, tp = confusion_matrix(y_test, y_predict_class1).ravel()  
2 print("True Negatives: ",tn)  
3 print("False Positives: ",fp)  
4 print("False Negatives: ",fn)  
5 print("True Positives: ",tp)
```

```
True Negatives: 8  
False Positives: 8  
False Negatives: 0  
True Positives: 9
```

```
In [31]: 1 from sklearn.metrics import classification_report  
2 print(classification_report(y_test, y_predict_class1))
```

	<u>precision</u>	<u>recall</u>	f1-score	support
0	1.00	0.50	0.67	16
1	0.53	1.00	0.69	9
micro avg	0.68	0.68	0.68	25
macro avg	0.76	0.75	0.68	25
weighted avg	0.83	0.68	0.68	25

Accuracy checking for different threshold values

```
In [32]: 1 from sklearn.preprocessing import binarize
          2 y_predict_class2 = binarize(y_prob.reshape(1,-1) (0.50))[0]
          3 y_predict_class2
```

```
Out[32]: array([1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 1., 0., 0.])
```

```
In [33]: 1 confusion_matrix_3 = confusion_matrix(y_test, y_predict_class2)
          2 print(confusion_matrix_3)
```

```
[[15  1]
 [ 5  4]]
```

```
In [34]: 1 from sklearn.metrics import classification_report
          2 print(classification_report(y_test, y_predict_class2))
```

	precision	recall	f1-score	support
0	0.75	0.94	0.83	16
1	0.80	0.44	0.57	9
micro avg	0.76	0.76	0.76	25
macro avg	0.78	0.69	0.70	25
weighted avg	0.77	0.76	0.74	25

Accuracy checking for different threshold values

```
In [35]: 1 from sklearn.preprocessing import binarize
2 y_predict_class3 = binarize(y_prob.reshape(1,-1), 0.70)[0]
3 y_predict_class3
```

```
Out[35]: array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,  
                0., 0., 0., 0., 0., 0., 0., 0.])
```

```
In [36]: 1 confusion_matrix_4 = confusion_matrix(y_test, y_predict_class3)
          2 print(confusion_matrix_4)
```

$$\begin{bmatrix} 16 & 0 \\ 9 & 0 \end{bmatrix}$$

```
In [37]: 1 from sklearn.metrics import classification_report
          2 print(classification_report(y_test, y_predict_class3))
```

	precision	recall	f1-score	support
0	0.64	1.00	0.78	16
1	0.00	0.00	0.00	9
micro avg	0.64	0.64	0.64	25
macro avg	0.32	0.50	0.39	25
weighted avg	0.41	0.64	0.50	25

Calculating Optimal Threshold Value

```
In [38]: 1 from sklearn.metrics import roc_curve, auc
```

```
In [39]: 1 fpr, tpr, thresholds= roc_curve(y_test, y_prob)
        2 roc_auc = auc(fpr, tpr)
```

```
In [40]: 1 print("Area under the ROC curve : %f" % roc_auc)
```

```
In [41]: 1 import numpy as np
        2 i = np.arange(len(tpr)) # index for df
        3 roc = pd.DataFrame({'fpr' : pd.Series(fpr, index=i), 'tpr' : pd.Series(tpr, index = i),
        4                  '1-fpr' : pd.Series(1-fpr, index = i), 'tf' : pd.Series(tpr - (1-fpr), index = i),
        5                  'thresholds' : pd.Series(thresholds, index = i)})
        6 roc.iloc[(roc.tf-0).abs().argsort()[:1]]
```

Out[41]:

	fpr	tpr	1-fpr	tf	thresholds
7	0.125	0.888889	0.875	0.013889	0.457033

Optimal Threshold Value in ROC Curve

In [42]:

```
1 fig, ax = plt.subplots()
2 plt.plot(roc['tpr'])
3 plt.plot(roc['1-fpr'], color = 'red')
4 plt.xlabel('1-False Positive Rate')
5 plt.ylabel('True Positive Rate')
6 plt.title('Receiver operating characteristic')
7 ax.set_xticklabels([])
```

Out [42]: []



Classification Report using Optimal Threshold Value

```
In [43]: 1 from sklearn.preprocessing import binarize
          2 y_predict_class4 = binarize(y_prob.reshape(1,-1), 0.45)[0]
          3 y_predict_class4
```

```
Out[43]: array([1., 1., 1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 1., 0., 0., 0., 1.,
                1., 0., 1., 0., 1., 1., 0., 0.])
```

```
In [44]: 1 confusion_matrix_5 = confusion_matrix(y_test, y_predict_class4)
          2 print(confusion_matrix_5)
```

```
[[14  2]
 [ 1  8]]
```

```
In [45]: 1 from sklearn.metrics import classification_report
          2 print(classification_report(y_test, y_predict_class4))
```

	precision	recall	f1-score	support
0	0.93	0.88	0.90	16
1	0.80	0.89	0.84	9
micro avg	0.88	0.88	0.88	25
macro avg	0.87	0.88	0.87	25
weighted avg	0.89	0.88	0.88	25

Thank You

