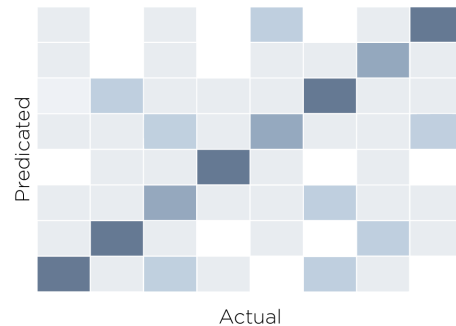**Confusion Matrix**



**What is our GOAL for this MODULE?**

The goal of this module is to check for the accuracy of the prediction model in linear and multi linear regression using the confusion matrix.

**What did we ACHIEVE in the class TODAY?**

We checked for the accuracy of the prediction model in linear regression and multi linear regression using the confusion matrix and heatmap

**Which CONCEPTS/CODING BLOCKS did we cover today?**

- Usage of sklearn library
- Training and testing of the prediction model
- Using the heatmap and creating a confusion matrix

**How did we DO the activities?**
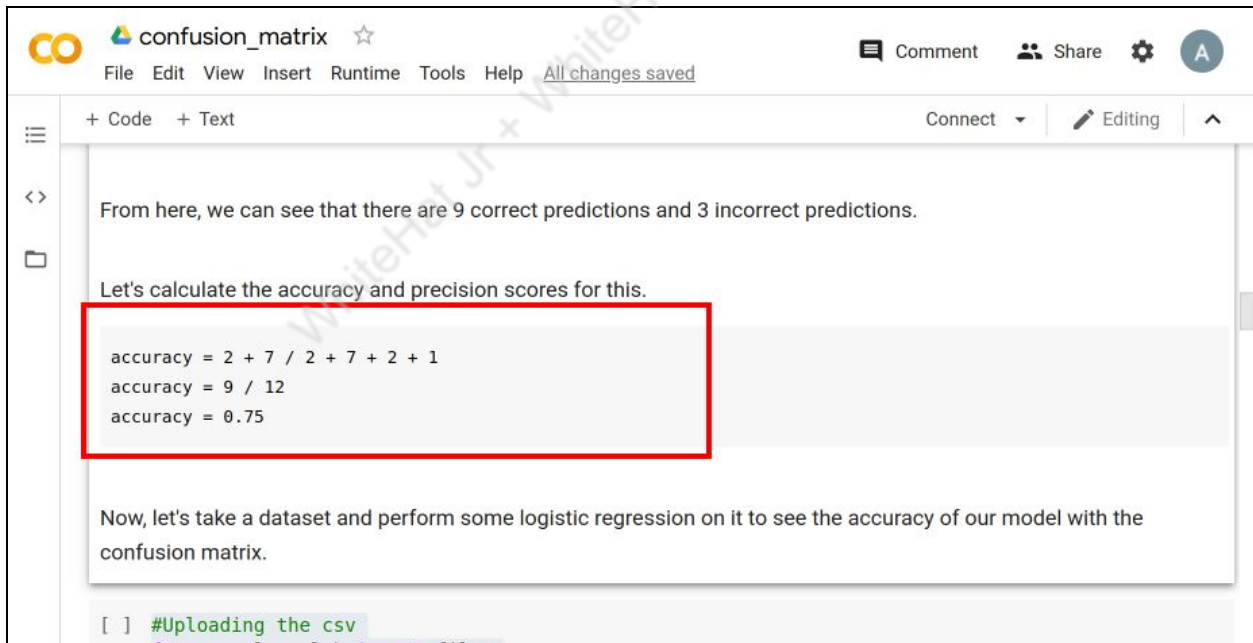
1. We learned about the confusion matrix.

## Actual Values

|  | Yes | No |
|---|---|---|
| **Yes** | True Positive | False Positive |
| **No** | False Negative | True Negative |

**Predicted Values**

2. We saw how our prediction matrix looks like.

| Predicted Class | Actual Class |
|---|---|
| No Cancer | No Cancer |
| Has Cancer | Has Cancer |
| No Cancer | No Cancer |
| No Cancer | No Cancer |
| No Cancer | Has Cancer |
| Has Cancer | Has Cancer |
| No Cancer | No Cancer |
| Has Cancer | No Cancer |
| No Cancer | No Cancer |
| No Cancer | No Cancer |
| Has Cancer | Has Cancer |
| No Cancer | No Cancer |

**Actual Class**

Predicted Class

|  | Has Cancer | No Cancer |
|---|---|---|
| Has Cancer | 2 | 2 |
| No Cancer | 1 | 7 |

3. We calculated the accuracy of the prediction model.

From here, we can see that there are 9 correct predictions and 3 incorrect predictions.

Let's calculate the accuracy and precision scores for this.

```
accuracy = 2 + 7 / 2 + 7 + 2 + 1
accuracy = 9 / 12
accuracy = 0.75
```

Now, let's take a dataset and perform some logistic regression on it to see the accuracy of our model with the confusion matrix.

[ ] #Uploading the csv

4. We uploaded the heart attack rate csv file in the colab notebook and using the pandas library we read it and print it's content.

5.  Then we split the data into 75% and 25% to train and test the prediction model.



Training the data:

6. Then we test the model using the remaining data and substitute the value of 0 as No and 1 as Yes and create 2 arrays for the values in heart attack predictions and actual_values and create a labels array with values Yes and No for the confusion matrix.

7. We plotted the heatmap for the data and calculated its accuracy.





8. We also checked how multiple factors( variables) affect the  accuracy of a prediction model.
9. We took multiple factors and created two different data frames and then split that data into 75% and 25%.

10. Then we made all the factors scalar as all had different measurements.



11. We trained our model on the remaining data.

```
confusion_matrix
File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code  + Text                                          Connect ▾   🖊 Editing  ⌃

[ ]  heart_attack = df["target"]

     factors_train, factors_test, heart_attack_train, heart_attack_test = train_test_split(factors, heart_attack, t
```

Since all of age, sex, cp and chol have different measurement units, let's make them scaler to analyse them well.

```
[ ]  from sklearn.preprocessing import StandardScaler
     sc_x = StandardScaler()

     factors_train = sc_x.fit_transform(factors_train)
     factors_test = sc_x.transform(factors_test)
```

```
[ ]  classifier2 = LogisticRegression(random_state = 0)
     classifier2.fit(factors_train, heart_attack_train)

 👤  LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                        intercept_scaling=1, l1_ratio=None, max_iter=100,
                        multi_class='auto', n_jobs=None, penalty='l2',
                        random_state=0, solver='lbfgs', tol=0.0001, verbose=0,
                        warm_start=False)
```

12. We then substituted the values of 0 by No and 1 by Yes and created it's list for the heat map.

```
confusion_matrix
File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code  + Text                                          Connect ▾   🖊 Editing  ⌃

[   heart_attack_prediction_1 = classifier2.predict(factors_test)

    predicted_values_1 = []
    for i in heart_attack_prediction_1:
      if i == 0:
        predicted_values_1.append("No")
      else:
        predicted_values_1.append("Yes")

    actual_values_1 = []
    for i in heart_attack_test.ravel():
      if i == 0:
        actual_values_1.append("No")
      else:
        actual_values_1.append("Yes")

[ ]  cm = confusion_matrix(actual_values_1, predicted_values_1, labels)
```
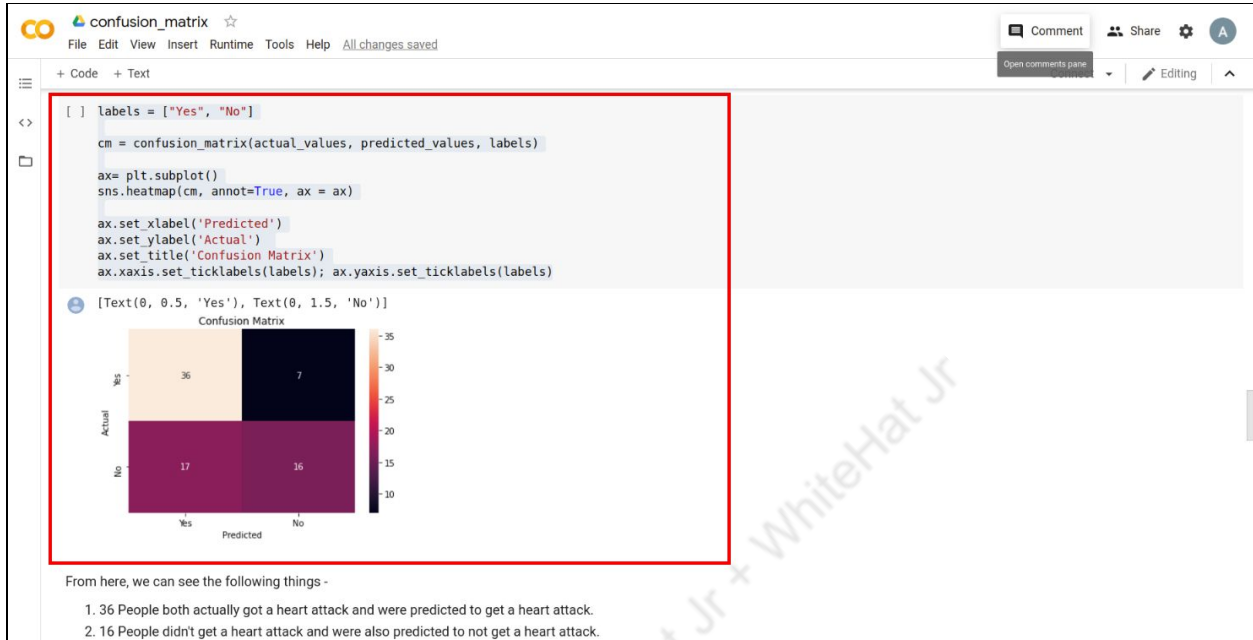
13. We plotted the heatmap.

```
[ ] cm = confusion_matrix(actual_values_1, predicted_values_1, labels)

    ax= plt.subplot()
    sns.heatmap(cm, annot=True, ax = ax)

    ax.set_xlabel('Predicted')
    ax.set_ylabel('Actual')
    ax.set_title('Confusion Matrix')
    ax.xaxis.set_ticklabels(labels); ax.yaxis.set_ticklabels(labels)
```
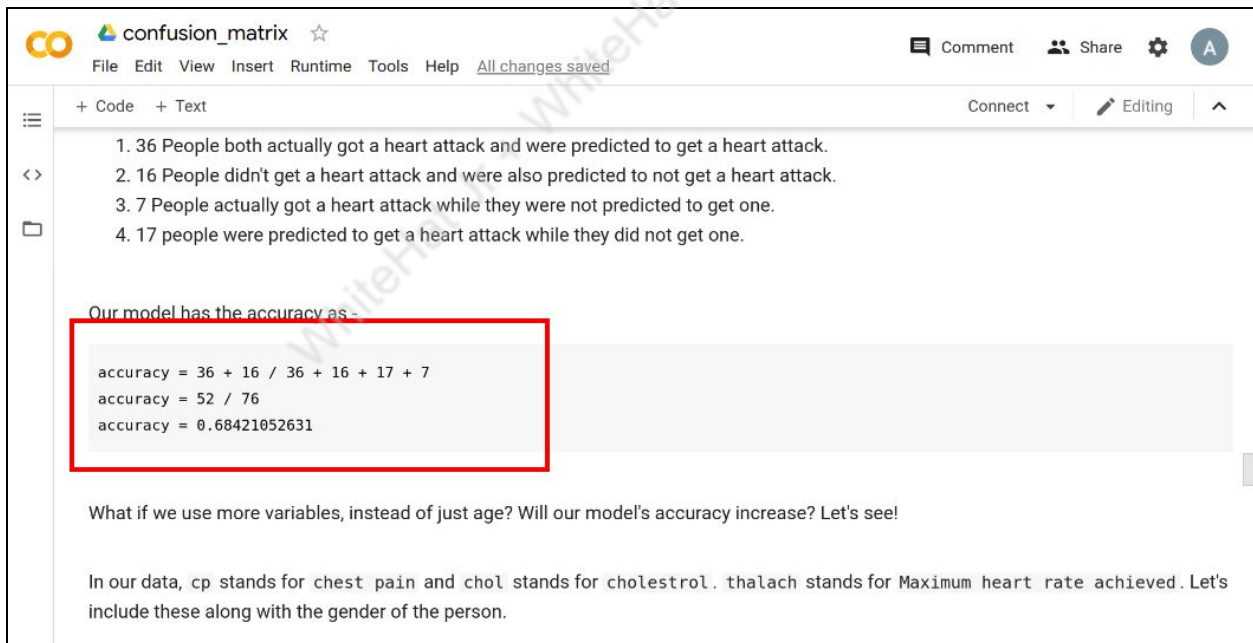
```
[Text(0, 0.5, 'Yes'), Text(0, 1.5, 'No')]
```

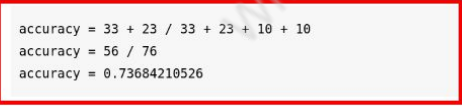14. Calculated the accuracy of the model.

3. 10 People actually got a heart attack while they were not predicted to get one.
4. 10 people were predicted to get a heart attack while they did not get one.

Our model has the accuracy as -

```
accuracy = 33 + 23 / 33 + 23 + 10 + 10
accuracy = 56 / 76
accuracy = 0.73684210526
```

With the new model that we just built, we have a higher accuracy to detect if a person will get a heart attack or not.

You can try this out with a combination of different set of variables, or add more variables to this to see if that improves the accuracy of the model?

**We concluded that more number of variables affects the accuracy of the model.**

**What's NEXT?**

In the next class, we will learn about clustering.

**EXTEND YOUR KNOWLEDGE:**

Learn more about the confusion matrix from the following doc:
https://www.geeksforgeeks.org/confusion-matrix-machine-learning/