

## Week 1

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from scipy import stats
5 from scipy.stats import chi2_contingency
6 import statsmodels.api as sm
7 from statsmodels.formula.api import ols
8 from sklearn import tree
9 from sklearn.cluster import KMeans
10 from sklearn.impute import SimpleImputer
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.model_selection import train_test_split
13 from sklearn.linear_model import LinearRegression
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.metrics import confusion_matrix, roc_curve, auc, classification_report, accuracy_score

```

## Load the dataset into a pandas DataFrame

```

1 df = pd.read_csv('/content/odi.csv')
2 df = df.dropna()
3 # Display the first few rows of the dataframe
4 print(df.head())

```

	mid	date	venue	bat_team	bowl_team	\
0	1	2006-06-13	Civil Service Cricket Club, Stormont	England	Ireland	
1	1	2006-06-13	Civil Service Cricket Club, Stormont	England	Ireland	
2	1	2006-06-13	Civil Service Cricket Club, Stormont	England	Ireland	
3	1	2006-06-13	Civil Service Cricket Club, Stormont	England	Ireland	
4	1	2006-06-13	Civil Service Cricket Club, Stormont	England	Ireland	

	batsman	bowler	runs	wickets	overs	runs_last_5	\
0	ME Trescothick	DT Johnston	0	0	0.1	0	
1	ME Trescothick	DT Johnston	0	0	0.2	0	
2	ME Trescothick	DT Johnston	4	0	0.3	4	
3	ME Trescothick	DT Johnston	6	0	0.4	6	
4	ME Trescothick	DT Johnston	6	0	0.5	6	

	wickets_last_5	striker	non-striker	total
0	0	0	0	301
1	0	0	0	301
2	0	0	0	301
3	0	0	0	301
4	0	0	0	301

## Week 2

## Calculate the probability of a team scoring more than 300 runs

```

1 # Assuming 'total' is the target variable we want to analyze
2 total_runs = df['total']
3 probability = (total_runs > 300).mean()
4 print(f"The probability of scoring more than 300 runs is {probability:.2f}")
5

```

The probability of scoring more than 300 runs is 0.23

## Average runs per over

```

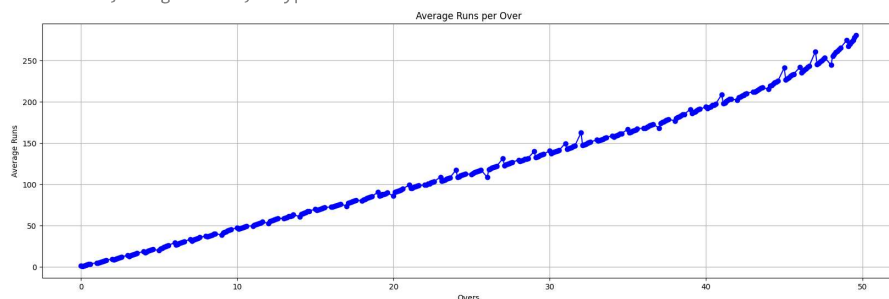
1 average_runs_per_over = df.groupby('overs')['runs'].mean()
2 print(average_runs_per_over)
3
4 plt.figure(figsize=(20, 6))
5 plt.plot(average_runs_per_over.index, average_runs_per_over.values, marker='o', linestyle='-', color='b')
6 plt.title('Average Runs per Over')
7 plt.xlabel('Overs')
8 plt.ylabel('Average Runs')
9 plt.grid(True)
10 plt.show()

```

```

overs
0.0      1.211864
0.1      0.635788
0.2      1.167340
0.3      1.793239
0.4      2.353323
...
49.2    270.046358
49.3    272.824661
49.4    275.153216
49.5    278.147129
49.6    281.134271
Name: runs, Length: 350, dtype: float64

```



Double-click (or enter) to edit

✓ Week 3

✓ Sampling

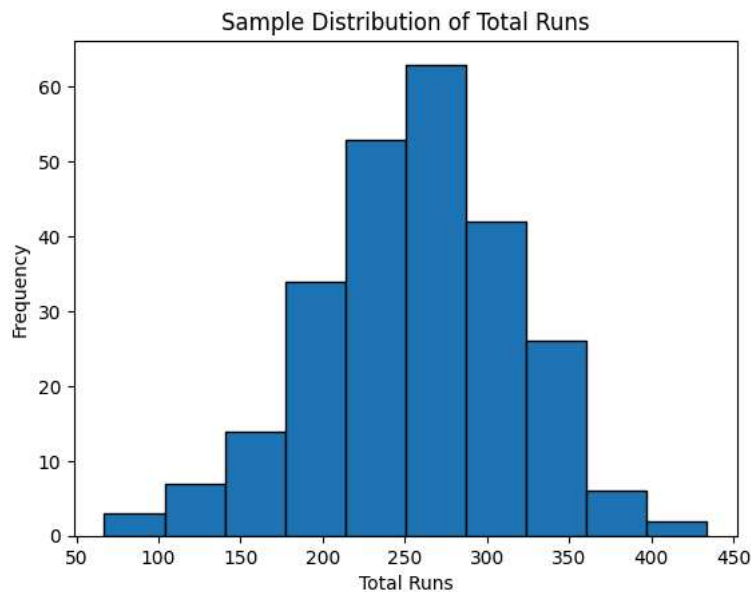
```

1 # Define the sample size you want to draw
2 sample_size = 250
3
4 # Perform simple random sampling
5 sample_df = df.sample(n=sample_size, random_state=1)
6
7 # Calculate sample mean and standard deviation for 'total' runs
8 sample_mean = sample_df['total'].mean()
9 sample_std = sample_df['total'].std()
10
11 print(f"Sample Mean: {sample_mean}")
12 print(f"Sample Standard Deviation: {sample_std}")
13
14 plt.hist(sample_df['total'], bins=10, edgecolor='black')
15 plt.title('Sample Distribution of Total Runs')
16 plt.xlabel('Total Runs')
17 plt.ylabel('Frequency')
18 plt.show()
19

```

Sample Mean: 256.292

Sample Standard Deviation: 60.77160098223327



Double-click (or enter) to edit

## Week 4

### Hypothesis Testing

```

1 # Let's say we want to test the hypothesis that the average total score for England is 260
2 # Null hypothesis (H0): mean total score for England is 260
3 # Alternative hypothesis (H1): mean total score for England is not 260
4
5 # Filter the dataset for England's batting team
6 england_df = df[df['bat_team'] == 'England']
7
8 # Group by match ID and take the maximum 'total' to get the total score per match for England
9 england_total_scores = england_df.groupby('mid')['total'].max().reset_index()
10
11 # Calculate the mean and standard deviation of total scores for England
12 mean_score = england_total_scores['total'].mean()
13 std_deviation = england_total_scores['total'].std()
14 print(f"Mean: {mean_score}")
15 print(f"Standard Deviation: {std_deviation}")
16
17 # Plot histogram of total scores for England
18 plt.hist(england_total_scores['total'], bins=20, alpha=0.7, color='blue', edgecolor='black')
19

```

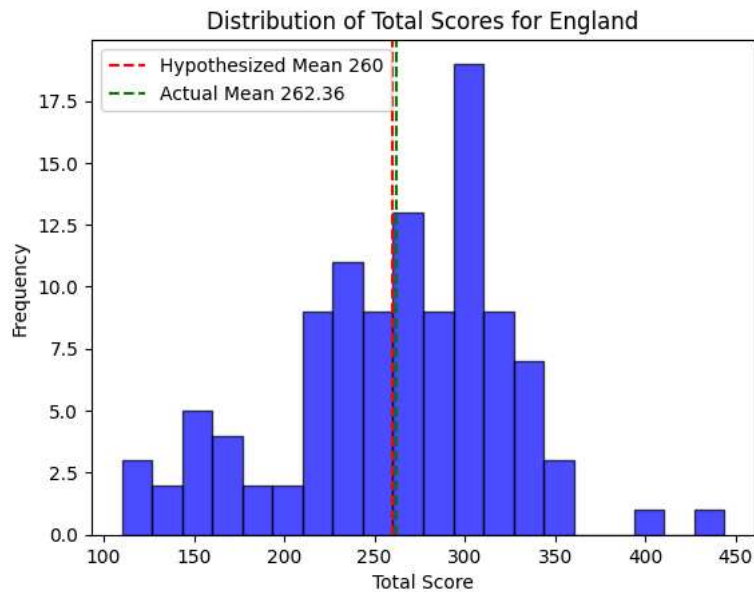
```

20 hypothesized_mean = 260
21 # Add a vertical line for the hypothesized mean score
22 plt.axvline(x=hypothesized_mean, color='red', linestyle='--', label=f'Hypothesized Mean {hypothesized_mean}')
23 plt.axvline(x=mean_score, color='green', linestyle='--', label=f"Actual Mean {mean_score:.2f}")
24 # Add labels and title
25 plt.xlabel('Total Score')
26 plt.ylabel('Frequency')
27 plt.title('Distribution of Total Scores for England')
28
29 # Add legend
30 plt.legend()
31
32 # Show plot
33 plt.show()
34
35 # Perform a one-sample t-test
36 t_statistic, p_value = stats.ttest_1samp(england_total_scores['total'], hypothesized_mean)
37
38 # Set your significance level
39 alpha = 0.05
40
41 # Check if we reject or fail to reject the null hypothesis
42 if p_value < alpha:
43     print(f"Reject the null hypothesis. p-value: {p_value}")
44 else:
45     print(f"Fail to reject the null hypothesis. p-value: {p_value}")

```

Mean: 262.35779816513764

Standard Deviation: 61.79798674291424



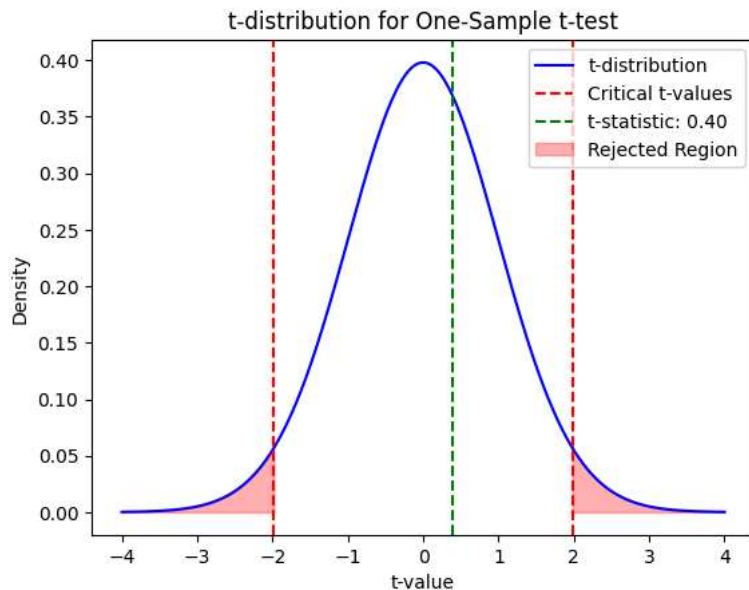
Fail to reject the null hypothesis. p-value: 0.6911716118103672

✓ T test

```

1 x = np.linspace(-4, 4, 1000)
2
3 # Calculate the t-distribution with degrees of freedom (df) equal to the sample size minus 1
4 df = len(england_total_scores) - 1
5 t_distribution = stats.t.pdf(x, df)
6
7 # Plot the t-distribution
8 plt.plot(x, t_distribution, 'b-', label='t-distribution')
9
10 # Plot the critical t-values for two-tailed test at alpha = 0.05
11 critical_t_values = [stats.t.ppf(0.025, df), stats.t.ppf(0.975, df)]
12 plt.axvline(critical_t_values[0], color='red', linestyle='--', label='Critical t-values')
13 plt.axvline(critical_t_values[1], color='red', linestyle='--')
14
15 # Plot the t-statistic
16 plt.axvline(t_statistic, color='green', linestyle='--', label=f't-statistic: {t_statistic:.2f}')
17
18 # Calculate the p-value
19 p_value = stats.t.sf(np.abs(t_statistic), df) * 2 # two-tailed test
20
21 # Shade the rejected region if p-value is less than alpha
22 alpha = 0.05
23 x_accepted_left = np.linspace(-4, critical_t_values[0], 100)
24 x_accepted_right = np.linspace(critical_t_values[1], 4, 100)
25 plt.fill_between(x_accepted_left, stats.t.pdf(x_accepted_left, df), color='red', alpha=0.3, label='Rejected Region')
26 plt.fill_between(x_accepted_right, stats.t.pdf(x_accepted_right, df), color='red', alpha=0.3)
27
28 # Add labels and title
29 plt.xlabel('t-value')
30 plt.ylabel('Density')
31 plt.title('t-distribution for One-Sample t-test')
32 plt.legend()
33
34 # Show plot
35 plt.show()
36
37 # Print p-value
38 print(f"P-value: {p_value:.4f}")

```



P-value: 0.6912

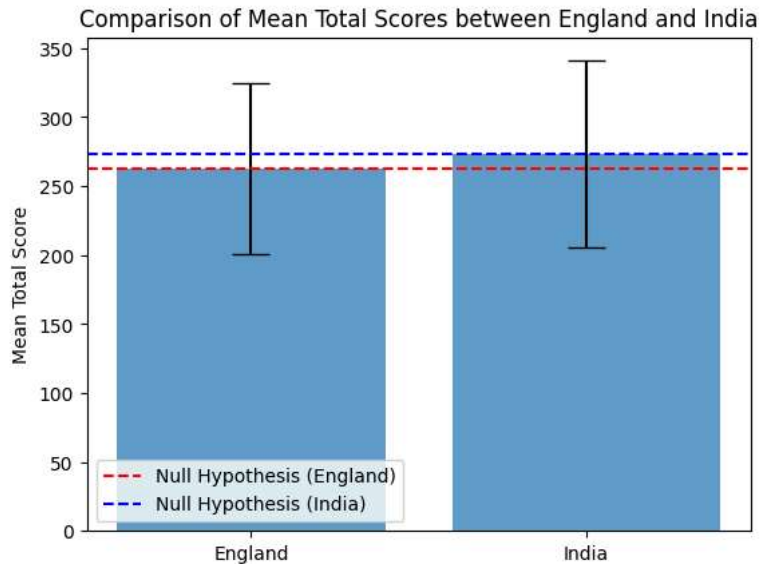
✓ Week 5

✓ Two Sample Testing and ANOVA

```

1 # say we want to compare the average total scores of two teams, England and India
2 # Null hypothesis (H0): There is no difference in the average total scores of England and India
3 # Alternative hypothesis (H1): There is a difference in the average total scores of England and India
4
5 # Filter the dataset for each team
6 england_df = df[df['bat_team'] == 'England']
7 england_total_scores = england_df.groupby('mid')['total'].max().reset_index()
8 india_df = df[df['bat_team'] == 'India']
9 india_total_scores = india_df.groupby('mid')['total'].max().reset_index()
10
11 # Calculate means and standard deviations for each team
12 england_mean = england_total_scores['total'].mean()
13 india_mean = india_total_scores['total'].mean()
14 england_std = england_total_scores['total'].std()
15 india_std = india_total_scores['total'].std()
16
17 # Set up the plot
18 teams = ['England', 'India']
19 x_pos = np.arange(len(teams))
20 means = [england_mean, india_mean]
21 std_devs = [england_std, india_std]
22
23 # Plot the bar chart
24 plt.bar(x_pos, means, yerr=std_devs, align='center', alpha=0.7, ecolor='black', capsize=10)
25 plt.xticks(x_pos, teams)
26 plt.ylabel('Mean Total Score')
27 plt.title('Comparison of Mean Total Scores between England and India')
28
29 # Add a horizontal line representing the null hypothesis
30 plt.axhline(y=england_mean, color='red', linestyle='--', label='Null Hypothesis (England)')
31 plt.axhline(y=india_mean, color='blue', linestyle='--', label='Null Hypothesis (India)')
32
33 # Add legend
34 plt.legend()
35
36 # Show plot
37 plt.show()
38
39 # Perform a two-sample t-test
40 t_statistic, p_value = stats.ttest_ind(england_total_scores['total'], india_total_scores['total'])
41
42 # Print p-value
43 print(f"P-value: {p_value:.4f}")
44
45 # Set your significance level
46 alpha = 0.05
47
48 # Check if we reject or fail to reject the null hypothesis
49 if p_value < alpha:
50     print(f'Reject the null hypothesis. p-value: {p_value:.4f}')
51 else:
52     print(f'Fail to reject the null hypothesis. p-value: {p_value:.4f}')
53

```



P-value: 0.2059

Fail to reject the null hypothesis. p-value: 0.2059

## Week 6

## Linear regression

```

1 def custom_accuracy(y_test,y_pred,threshold):
2     right = 0
3
4     l = len(y_pred)
5     for i in range(0,l):
6         if(abs(y_pred[i]-y_test[i]) <= threshold):
7             right += 1
8     return ((right/l)*100)
9
10
11 X = df.iloc[:,[7,8,9,12,13]].values
12 y = df.iloc[:, 14].values
13
14
15 # Splitting the dataset into the Training set and Test set
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
17
18 # Feature Scaling
19 sc = StandardScaler()
20 X_train = sc.fit_transform(X_train)
21 X_test = sc.transform(X_test)
22
23 # Training the dataset
24 lin = LinearRegression()
25 lin.fit(X_train,y_train)
26
27 # Testing the dataset on trained model
28 y_pred = lin.predict(X_test)
29 score = lin.score(X_test,y_test)*100
30 print("R square value:" , score)
31 print("Custom accuracy:" , custom_accuracy(y_test,y_pred,20))
32
33 # Testing with a custom input
34 new_prediction = lin.predict(sc.transform(np.array([[100,0,13,50,50]])))
35 print("Prediction score:" , new_prediction)
36
37

```

R square value: 52.737657811129445  
 Custom accuracy: 43.354801937874036  
 Prediction score: [322.42983935]

## Week 7

### Linear Regression and Multiple Regression

```

1 # predict the 'total' score using multiple predictors.
2 # For this example, we'll use 'runs', 'wickets', 'overs', 'runs_last_5', and 'wickets_last_5'.
3
4 # Define the predictors and the response variable
5 predictors = ['runs', 'wickets', 'overs', 'runs_last_5', 'wickets_last_5']
6 X = df[predictors]
7 y = df['total']
8
9 # Add a constant to the model (intercept)
10 X = sm.add_constant(X)
11
12 # Fit the multiple linear regression model
13 model = sm.OLS(y, X).fit()
14
15 # Get the summary of the regression
16 summary = model.summary()
17 print(summary)
18
19 # You can also use the model to make predictions
20 # Here's an example of predicting the 'total' for a new observation
21 new_observation = {'runs': 50, 'wickets': 2, 'overs': 10, 'runs_last_5': 25, 'wickets_last_5': 1}
22 new_observation['total'] = 0
23 new_X = pd.DataFrame([new_observation])
24 new_X = sm.add_constant(new_X)
25
26 predicted_total = model.predict(new_X)
27 print(f"Predicted Total: {predicted_total[0]}")

```

#### OLS Regression Results

```

=====
Dep. Variable:          total    R-squared:                0.527
Model:                  OLS      Adj. R-squared:           0.527
Method:                 Least Squares    F-statistic:         7.834e+04
Date:                  Mon, 08 Apr 2024    Prob (F-statistic):       0.00
Time:                  09:29:32    Log-Likelihood:        -1.8166e+06
No. Observations:      350899    AIC:                   3.633e+06
Df Residuals:          350893    BIC:                   3.633e+06
Df Model:               5
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	252.8874	0.214	1182.966	0.000	252.468	253.306
runs	1.0811	0.003	313.979	0.000	1.074	1.088
wickets	-12.9809	0.063	-205.223	0.000	-13.105	-12.857
overs	-3.5606	0.021	-172.843	0.000	-3.601	-3.520
runs_last_5	0.2146	0.009	22.615	0.000	0.196	0.233
wickets_last_5	-3.6459	0.106	-34.291	0.000	-3.854	-3.438

```

=====
Omnibus:                12201.805    Durbin-Watson:           0.016
Prob(Omnibus):           0.000    Jarque-Bera (JB):        30831.871
Skew:                   -0.158    Prob(JB):                0.00
Kurtosis:                4.417    Cond. No.                 429.
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
 Predicted Total: 12427.924298030672

## Week 8

### Concepts of MLE (Maximum Likelihood Estimation) and Logistic Regression

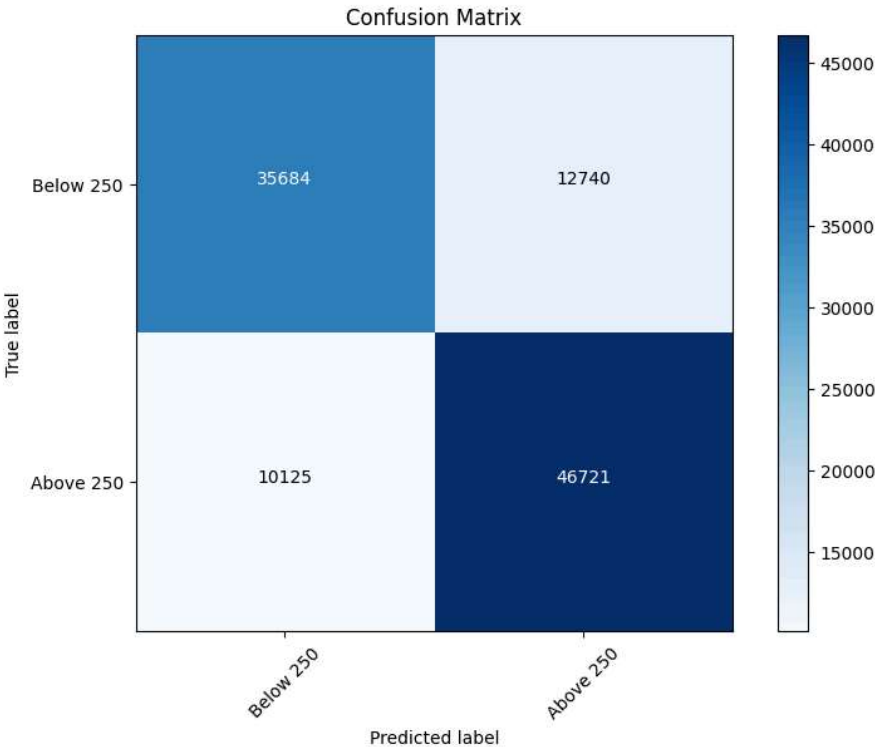


```

1 df['high_score'] = df['total'].apply(lambda x: 1 if x > 250 else 0)
2
3 # Define the predictors and the binary target variable
4 predictors = ['runs', 'wickets', 'overs', 'runs_last_5', 'wickets_last_5']
5 X = df[predictors]
6 y = df['high_score']
7
8 # Split the data into training and testing sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
10
11 # Initialize the Logistic Regression model
12 log_reg = LogisticRegression()
13
14 # Fit the model to the training data
15 log_reg.fit(X_train, y_train)
16
17 # Predict on the testing data
18 y_pred = log_reg.predict(X_test)
19
20 # Evaluate the model
21 print(confusion_matrix(y_test, y_pred))
22 print(classification_report(y_test, y_pred))
23
24 # Visualize Confusion Matrix
25 cm = confusion_matrix(y_test, y_pred)
26 plt.figure(figsize=(8, 6))
27 plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
28 plt.title('Confusion Matrix')
29 plt.colorbar()
30 classes = ['Below 250', 'Above 250']
31 tick_marks = np.arange(len(classes))
32 plt.xticks(tick_marks, classes, rotation=45)
33 plt.yticks(tick_marks, classes)
34
35 thresh = cm.max() / 2.
36 for i in range(cm.shape[0]):
37     for j in range(cm.shape[1]):
38         plt.text(j, i, format(cm[i, j], 'd'),
39                 horizontalalignment="center",
40                 color="white" if cm[i, j] > thresh else "black")
41
42 plt.tight_layout()
43 plt.ylabel('True label')
44 plt.xlabel('Predicted label')
45 plt.show()
46
47

```

[[35684 12740]					
[10125 46721]]					
	precision	recall	f1-score	support	
0	0.78	0.74	0.76	48424	
1	0.79	0.82	0.80	56846	
accuracy			0.78	105270	
macro avg	0.78	0.78	0.78	105270	
weighted avg	0.78	0.78	0.78	105270	

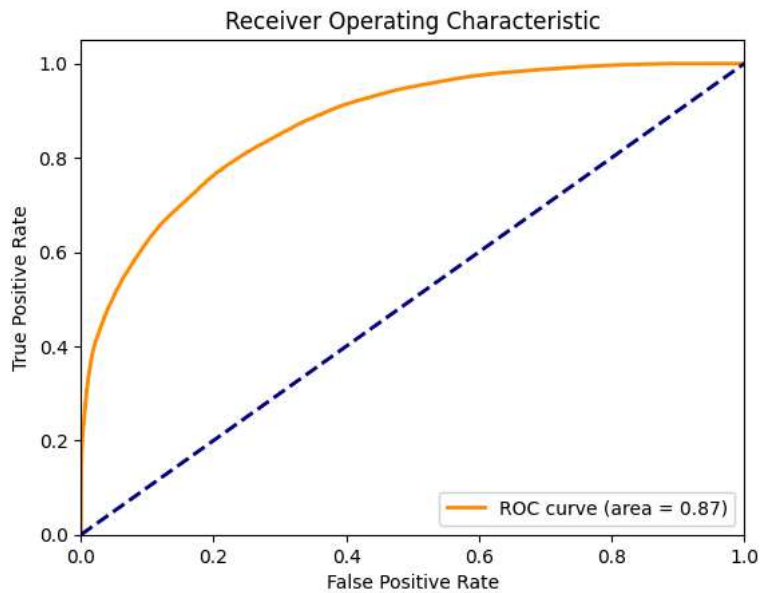


- Week 9
- ROC (Receiver Operating Characteristic) and Regression Analysis Model Building

```

1 # For logistic regression, we need a binary target variable. Let's create one for illustration.
2 # For example, we'll predict if the total score will be above or below 250.
3 df['high_score'] = df['total'].apply(lambda x: 1 if x > 250 else 0)
4
5 # Define the predictors and the binary target variable
6 predictors = ['runs', 'wickets', 'overs', 'runs_last_5', 'wickets_last_5']
7 X = df[predictors]
8 y = df['high_score']
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
12
13 # Initialize the Logistic Regression model
14 log_reg = LogisticRegression()
15
16 # Fit the model to the training data
17 log_reg.fit(X_train, y_train)
18
19 # Predict probabilities for the test data
20 y_probs = log_reg.predict_proba(X_test)[:, 1]
21
22 # Compute ROC curve and ROC area
23 fpr, tpr, thresholds = roc_curve(y_test, y_probs)
24 roc_auc = auc(fpr, tpr)
25
26 # Plot ROC curve
27 plt.figure()
28 plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
29 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
30 plt.xlim([0.0, 1.0])
31 plt.ylim([0.0, 1.05])
32 plt.xlabel('False Positive Rate')
33 plt.ylabel('True Positive Rate')
34 plt.title('Receiver Operating Characteristic')
35 plt.legend(loc="lower right")
36 plt.show()

```



✓ Week 10

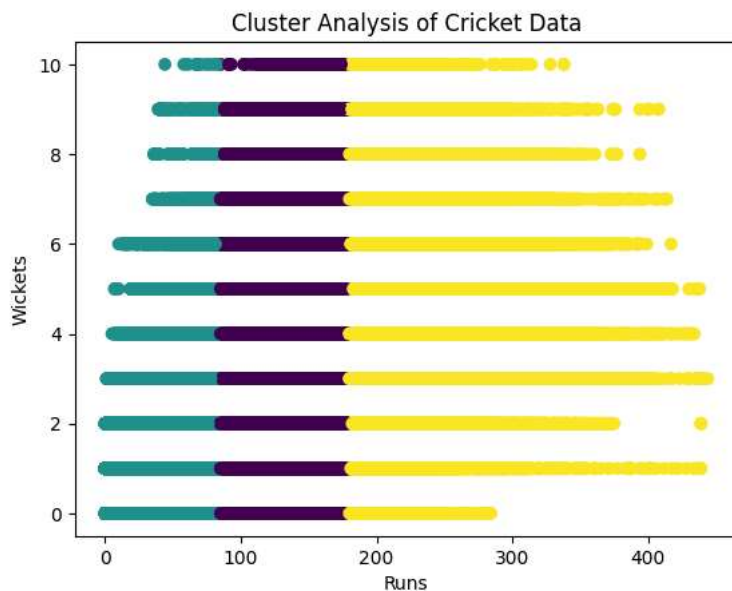
✓ Chi-square Test and Introduction to Cluster Analysis

```

1 # Chi-square Test
2 # to test if there's a significant association between 'bat_team' and 'bowl_team' winning.
3
4 # Create a contingency table
5 contingency_table = pd.crosstab(df['bat_team'], df['bowl_team'])
6
7 # Perform the Chi-square test
8 chi2, p, dof, expected = chi2_contingency(contingency_table)
9
10 # Set your significance level
11 alpha = 0.05
12
13 # Check if we reject or fail to reject the null hypothesis
14 print(f'Chi-square Statistic: {chi2}, p-value: {p}')
15 if p < alpha:
16     print('Reject the null hypothesis, indicating a significant association between teams.')
17 else:
18     print('Fail to reject the null hypothesis, no significant association was found.')
19
20 # Introduction to Cluster Analysis
21 # For cluster analysis, let's use 'runs' and 'wickets' for clustering the data points.
22
23 # Define the data for clustering
24 cluster_data = df[['runs', 'wickets']]
25
26 # Initialize KMeans
27 kmeans = KMeans(n_clusters=3, random_state=0)
28
29 # Fit the model
30 kmeans.fit(cluster_data)
31
32 # Predict the clusters
33 labels = kmeans.predict(cluster_data)
34
35 # Plot the clusters
36 plt.scatter(cluster_data['runs'], cluster_data['wickets'], c=labels)
37 plt.xlabel('Runs')
38 plt.ylabel('Wickets')
39 plt.title('Cluster Analysis of Cricket Data')
40 plt.show()

```

Chi-square Statistic: 920349.3894003384, p-value: 0.0  
 Reject the null hypothesis, indicating a significant association between teams.  
 /usr/local/lib/python3.10/dist-packages/sklearn/cluster/\_kmeans.py:870: FutureWarning:  
 warnings.warn(



```
1 # Separate bat_team and bowl_team by country
2 # Extract the country from the bat_team column
3 df['bat_country'] = df['bat_team'].apply(lambda x: x.split()[0])
4 # Extract the country from the bowl_team column
5 df['bowl_country'] = df['bowl_team'].apply(lambda x: x.split()[0])
6
7 # Chi-square Test
8
9 # Create a contingency table
10 contingency_table = pd.crosstab(df['bat_country'], df['bowl_country'])
11 # Perform the Chi-square test
12 chi2, p, dof, expected = chi2_contingency(contingency_table)
13
14 # Set the significance level
15 alpha = 0.05
16
17 print(f'Chi-square Statistic: {chi2}, p-value: {p}')
18 if p < alpha:
19     print('Reject the null hypothesis, indicating a significant association between teams.')
20 else:
21     print('Fail to reject the null hypothesis, no significant association was found.')
22
23 # Cluster Analysis
24
25 # Select relevant columns for clustering
26 cluster_data = df[['runs', 'wickets', 'bat_country', 'bowl_country']]
27 # Select matches where India is the batting team and England is the bowling team
28 india_data = cluster_data[(cluster_data['bat_country'] == 'India') & (cluster_data['bowl_country'] == 'England')]
29
30 # Fit KMeans model
31
32 # Initialize KMeans model with 3 clusters
33 kmeans = KMeans(n_clusters=3, random_state=0)
34 # Fit the model to the selected data
35 kmeans.fit(india_data[['runs', 'wickets']])
36
37 # Predict clusters
38
39 # Predict cluster labels for the data points
40 labels = kmeans.predict(india_data[['runs', 'wickets']])
41
42 # Plot clusters
43
44 # Scatter plot of runs vs. wickets with clusters colored according to labels
45 plt.scatter(india_data['runs'], india_data['wickets'], c=labels)
46 plt.xlabel('Runs')
47 plt.ylabel('Wickets')
48 plt.title('Cluster Analysis of Cricket Data (India)')
49 plt.show()
50
```

## Week 11

### Cluster Analysis of Cricket Data (India)

#### Clustering Analysis

```

1 # Selecting features for clustering (for example, 'runs' and 'wickets')
2 features = df[['runs', 'wickets']]
3
4 # Standardizing the features for better clustering performance
5 scaler = StandardScaler()
6 features_scaled = scaler.fit_transform(features)
7
8 # Using the elbow method to find the optimal number of clusters
9 wcss = []
10 for i in range(1, 11):
11     kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
12     kmeans.fit(features_scaled)
13     wcss.append(kmeans.inertia_)
14
15 # Plotting the results onto a line graph to observe 'The elbow'
16 plt.plot(range(1, 11), wcss)
17 plt.title('The Elbow Method')
18 plt.xlabel('Number of clusters')
19 plt.ylabel('WCSS') # Within cluster sum of squares
20 plt.show()
21
22 # Applying K-Means to the dataset with the optimal number of clusters
23 optimal_clusters = 3 # This is an example, you should choose the number based on the elbow method
24 kmeans = KMeans(n_clusters=optimal_clusters, init='k-means++', max_iter=300, n_init=10, random_state=0)
25 y_kmeans = kmeans.fit_predict(features_scaled)
26
27 # Visualizing the clusters
28 plt.scatter(features_scaled[y_kmeans == 0, 0], features_scaled[y_kmeans == 0, 1], s=100, c='red', label='Cluster 1')
29 plt.scatter(features_scaled[y_kmeans == 1, 0], features_scaled[y_kmeans == 1, 1], s=100, c='blue', label='Cluster 2')
30 plt.scatter(features_scaled[y_kmeans == 2, 0], features_scaled[y_kmeans == 2, 1], s=100, c='green', label='Cluster 3')
31 plt.title('Clusters of players')
32 plt.xlabel('Runs (standardized)')
33 plt.ylabel('Wickets (standardized)')
34 plt.legend()
35 plt.show()

```

