

DATA WAREHOUSE

Final Project

Analysis of Joe Root's all International Cricket Centuries

KUHDSE22.1F – 018

KUHDSE22.1F – 023

KUHDSE22.1F – 004

KUHDSE22.1F – 019

Objective

Our dataset is Joe Root's All International Cricket Centuries. Joe Root was the captain of the England cricket team and a good performer of the team. He is a all-round cricketer in the England team. We are going to analysis specifically targeting Joe Root's international career and his all centuries of all formats.

We can use data mining algorithms to this dataset to extract meaningful insights and patterns to get many potential benefits.

- Can analyze the performance of international cricket of him. Getting to know the trends, patterns, and factors which help to score successful centuries.
- Analysis of this dataset may give valuable insights to other international teams and coaches to prepare the teams and plan the games against England team.
- Can Identify the strengths and weaknesses of Joe Root according to the century patterns.
- Coaches and selectors are able to understand the things to improve and development of skills of the Joe Root.
- Also, they can understand the countries that the Joe Root has performed well. It will help to take decisions when select the team against those countries.
- This will also be helpful to predict the matches results like the game changer of the match.

Data set description

Link to dataset :-

<https://www.kaggle.com/drahulsingh/joe-root-all-international-cricket-centuries>

This dataset includes Joe Root's all century details of all formats (Test, T20, and ODI). The details this dataset is included are the country that the century had been taken, position that Joe Root had played, the inning that the century has been taken, venue of the game played, the ground, date, score of Joe Root and the result of the match. Like that there are 9 columns in this dataset.

Original Data Set Features:

No	Number of the row.
Score	Score that has been taken by the Joe Root of each match that he got centuries.
Against	The country against which Joe Root scored a century.
Position	The place in the team that Joe Root has played at that match.
Innings	The Inning that the Joe Root got the century.
Venue	The place that the match had been played

Ground	The ground place where the match had been played (Home/Away/Neutral)
Date	The date that the match had been hold.
Result	Th result of the match (Won/Lost/Drawn)

Data Preprocessing

In this dataset there are no any missing values of each attribute, the missing percentage of all attributes is 0%.

We have removed an unnecessary column that we do not need. We removed No column that column denoted the number of each row.

In this dataset we wanted to check the outliers. We used the Interquartile Range algorithm to check the outliers.

As the result of that algorithm in this dataset there is not any outliers or any extreme values found. We can keep the dataset 46 instances as it is.

Filter

Choose **InterquartileRange** -R first-last -O 3.0 -E 6.0

Current relation
Relation: Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-R1-weka.filters.uns... Attributes: 10
Instances: 46 Sum of weights: 46

Attributes

All None Invert Pattern

No.	Name
1	<input type="checkbox"/> Score
2	<input type="checkbox"/> Against
3	<input type="checkbox"/> Position
4	<input type="checkbox"/> Innings
5	<input type="checkbox"/> Venue
6	<input type="checkbox"/> Ground
7	<input type="checkbox"/> Date
8	<input type="checkbox"/> Result
9	<input type="checkbox"/> Outlier
10	<input checked="" type="checkbox"/> ExtremeValue

Applying the Interquartile Range

Selected attribute

Name: ExtremeValue
Missing: 0 (0%)
Distinct: 1
Type: Nominal
Unique: 0 (0%)

No.	Label	Count	Weight
1	no	46	46
2	yes	0	0

No any Extreme Value Found

Selected attribute

Name: Outlier
Missing: 0 (0%)
Distinct: 1
Type: Nominal
Unique: 0 (0%)

No.	Label	Count	Weight
1	no	46	46
2	yes	0	0

No any Outliers Found

Data Transformation

In our dataset some of attributes are in numeric type and one attribute is a string type

- Score – String
- Position – Numeric
- Innings – Numeric

```
@relation Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-R1
@attribute Score string
@attribute Against {' New Zealand',' Australia',' West Indies',' Sri Lanka',' India',' South Africa',' Pakistan',' Bangladesh'}
@attribute Position numeric
@attribute Innings numeric
@attribute Venue {'Headingley, Leeds','Lord\'s, London','Sir Vivian Richards Stadium, North Sound','Trent Bridge, Nottingham','The Oval, London','Pallekele International Cricket Stadium, Kandy','Wellington Regional Stadium, Wellington','National Cricket Stadium, St. George\'s','Edgbaston, Birmingham','Sophia Gardens, Cardiff','Wanderers Stadium, Johannesburg','SuperSport Park, Centurion','Old Trafford, Manchester','Saurashtra Cricket Association Stadium, Rajkot','Kensington Oval, Bridgetown','Edgbaston Cricket Ground, Birmingham','University Oval, Dunedin','Daren Sammy Cricket Ground, Gros Islet','Rose Bowl, Southampton','Seddon Park, Hamilton','Galle International Stadium, Galle','M. A. Chidambaram Stadium, Chennai','Basin Reserve, Wellington'}
@attribute Ground {'Home,Away,Neutral'}
@attribute Date {'24-May-13,18-Jul-13,05-Mar-14,12-Jun-14,09-Jul-14,15-Aug-14,05-Sep-14,10-Dec-14,01-Mar-15,21-Apr-15,09-Jun-15,17-Jun-15,08-Jul-15,06-Aug-15,14-Jan-16,09-Feb-16,12-Feb-16,22-Jul-16,09-Nov-16,09-Mar-17,01-Jun-17,06-Jul-17,17-Aug-17,07-Mar-18,14-Jul-18,17-Jul-18,07-Sep-18,14-Nov-18,09-Feb-19,20-Feb-19,03-Jun-19,14-Jun-19,29-Nov-19,14-Jan-21,22-Jan-21,05-Feb-21,04-Aug-21,12-Aug-21,25-Aug-21,08-Mar-22,16-Mar-22,02-Jun-22,10-Jun-22,01-Jul-22,24-Feb-23,16-Jun-23'}
@attribute Result {'Won,Drawn,Lost'}
```

Here, we want to convert the numeric types to nominal and the string type to nominal.

We used **StringToNominal** algorithm to convert the string value of score to nominal.

Filter		Choose StringToNominal -R first		Apply Stop	
Current relation		Relation: Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-R1-weka.filters.uns...		Attributes: 10 Instances: 46 Sum of weights: 46	
Attributes		All None Invert Pattern		Selected attribute	
No. Name		1 Score		Name: Score Missing: 0 (0%) Distinct: 38 Type: Nominal Unique: 31 (67%)	
No.	Label	Count	Weight		
1	104	2	2		
2	180	1	1		
3	107	2	2		
4	200*	1	1		

After convert to nominal

```
@attribute Score {'104,180,107,'200*','154*','149*','113','104*','121','182*','106*','134,130,110,125,109,254,124,101,'133*','190,136,102,'113*','100*','122,'100*','226,228,186,218,'180*','153,'115*','176,142*','153*','118*}'}
```

We use the method call **discretize** to convert the numerical values to categorical values. We can use equal width binning or equal frequency binning to do this.

We took 2 bins to do the discretize the Position because he has only played as Top order and Middle order batsman and took 4 bins to discretize the Innings because there are only maximum 4 inning in a match. After that the Position and Innings looked like this.

```
@attribute Position {'\(-inf-3.5]\','\'(3.5-inf)\'}
```

```
@attribute Innings {'\(-inf-1.75]\','\'(1.75-2.5]\','\'(2.5-3.25]\','\'(3.25-inf)\'}
```

Then we replaced the categories names manually using notepad

```
@attribute Position {Top-Order,Middle-Order}
```

```
@attribute Innings {1st,2nd,3rd,4th}
```

After the transformation the dataset

```
@relation Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-R1-weka.filters.unsupervised.attribute.InterquartileRange-Rfirst-last-03.0-E6.0-weka.filters.unsupervised.attribute.StringToNominal-Rfirst-weka.filters.unsupervised.attribute.Discretize-B2-M-1.0-R3-precision6-weka.filters.unsupervised.attribute.Discretize-B4-M-1.0-R4-precision6

@attribute Score {104,180,107,'200*','154*','149*','113','104*','121','182*','106*','134','130','110','125','109','254','124','101','133*','190','136','102','113*','100*','122','100*','226','228','186','218','180*','153','115*','176','142','153*','118*'}
@attribute Against {'New Zealand','Australia','West Indies','Sri Lanka','India','South Africa','Pakistan','Bangladesh'}
@attribute Position {Top-Order,Middle-Order}
@attribute Innings {1st,2nd,3rd,4th}
@attribute Venue {'Headingley, Leeds','Lord's, London','Sir Vivian Richards Stadium, North Sound','Trent Bridge, Nottingham','The Oval, London','Pallekele International Cricket Stadium, Kandy','Wellington Regional Stadium, Wellington','National Cricket Stadium, St. George's','Edgbaston, Birmingham','Sophia Gardens, Cardiff','Wanderers Stadium, Johannesburg','SuperSport Park, Centurion','Old Trafford, Manchester','Saurashtra Cricket Association Stadium, Rajkot','Kensington Oval, Bridgetown','Edgbaston Cricket Ground, Birmingham','University Oval, Dunsedin','Daren Sammy Cricket Ground, Gros Islet','Rose Bowl, Southampton','Seddon Park, Hamilton','Galle International Stadium, Galle','M. A. Chidambaram Stadium, Chennai','Basin Reserve, Wellington'}
@attribute Ground {Home,Away,Neutral}
@attribute Date {24-May-13,18-Jul-13,05-Mar-14,12-Jun-14,09-Jul-14,15-Aug-14,05-Sep-14,10-Dec-14,01-Mar-15,21-Apr-15,09-Jun-15,17-Jun-15,08-Jul-15,06-Aug-15,14-Jan-16,09-Feb-16,12-Feb-16,22-Jul-16,09-Nov-16,00-Mar-17,01-Jun-17,06-Jul-17,17-Aug-17,07-Mar-18,14-Jul-18,17-Jul-18,07-Sep-18,14-Nov-18,09-Feb-19,20-Feb-19,03-Jun-19,14-Jun-19,29-Nov-19,14-Jan-21,22-Jan-21,05-Feb-21,04-Aug-21,12-Aug-21,25-Aug-21,08-Mar-22,16-Mar-22,02-Jun-22,10-Jun-22,01-Jul-22,24-Feb-23,16-Jun-23}
@attribute Result {Won,Drawn,Lost}
@attribute Outlier {no,yes}
@attribute ExtremeValue {no,yes}

@data
104,'New Zealand',Middle-Order,1st,'Headingley, Leeds',Home,24-May-13,Won,no,no
180,'Australia',Top-Order,3rd,'Lord's, London',Home,18-Jul-13,Won,no,no
107,'West Indies',Middle-Order,1st,'Sir Vivian Richards Stadium, North Sound',Away,05-Mar-14,Won,no,no
'200*','Sri Lanka',Middle-Order,1st,'Lord's, London',Home,12-Jun-14,Drawn,no,no
'154*','India',Middle-Order,2nd,'Trent Bridge, Nottingham',Home,09-Jul-14,Drawn,no,no
'149*','India',Middle-Order,2nd,'The Oval, London',Home,15-Aug-14,Won,no,no
113,'India',Middle-Order,1st,'Headingley, Leeds',Home,05-Sep-14,Won,no,no
'104*','Sri Lanka',Middle-Order,2nd,'Pallekele International Cricket Stadium, Kandy',Away,10-Dec-14,Won,no,no
121,'Sri Lanka',Middle-Order,1st,'Wellington Regional Stadium, Wellington',Neutral,01-Mar-15,Lost,no,no
```

Once all above processes were done, the preprocessing and transformation parts were completed. Now we are able to use the processed data set to do any type of data mining techniques.

Data Mining Process

Now it can be applied some data mining algorithms to these processed dataset. Here, we are going to use Apriori, Decision tree (J48), and Naive Bayes algorithms.

Apriori algorithm

This is an algorithm that is used to association rule mining. This can be used to this dataset to discover interesting associations or patterns among the attributes.

The result of this algorithm after adding it to the dataset will help to get the patterns or associations. From this it can be understood that Joe Root tends to score centuries more often against certain teams or in specific match formats.

First applied this algorithm to dataset with the min_confidence = 0.9 and min_sup = 0.1. But considering the patterns of that those are not reliable. Therefore we tried many times with changing the min_confidence and the min_sup.

```

Choose Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

St... Stop
Result list (right-click...)
21:52:36 - Apriori

Associator output

Position
Innings
Venue
Ground
Date
Result

=== Associator model (full training set) ===

Apriori
=====

Minimum support: 0.1 (5 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

Size of set of large itemsets L(1): 16

Size of set of large itemsets L(2): 33

Size of set of large itemsets L(3): 21

Size of set of large itemsets L(4): 3

Best rules found:

1. Against= Sri Lanka 6 ==> Position=Middle-Order 6 <conf:(1)> lift:(1.59) lev:(0.05) [2] conv:(2.22)
2. Venue=Lord's, London 6 ==> Ground=Home 6 <conf:(1)> lift:(1.77) lev:(0.06) [2] conv:(2.61)
3. Venue=Trent Bridge, Nottingham 6 ==> Ground=Home 6 <conf:(1)> lift:(1.77) lev:(0.06) [2] conv:(2.61)
4. Against= New Zealand Result=Won 5 ==> Ground=Home 5 <conf:(1)> lift:(1.77) lev:(0.05) [2] conv:(2.17)
5. Against= New Zealand Ground=Home 5 ==> Result=Won 5 <conf:(1)> lift:(1.48) lev:(0.04) [1] conv:(1.63)
6. Against= India Innings=2nd 5 ==> Ground=Home 5 <conf:(1)> lift:(1.77) lev:(0.05) [2] conv:(2.17)
7. Innings=2nd Venue=Trent Bridge, Nottingham 5 ==> Ground=Home 5 <conf:(1)> lift:(1.77) lev:(0.05) [2] conv:(2.17)

```

Minimum support = 0.1 | Minimum confidence = 0.9 | Number of cycles performed = 18

Like this we applied the apriori algorithm so many times to the dataset.

Finally we can got patterns that can be reliable after checking the pattern.

```

St... Stop
Result list (right-click...)
21:52:36 - Apriori
23:03:06 - Apriori
23:03:44 - Apriori
23:04:19 - Apriori
23:05:10 - Apriori
23:05:19 - Apriori
23:06:21 - Apriori
23:06:43 - Apriori
23:07:29 - Apriori

Associator output

Scheme: weka.associations.Apriori -N 10 -T 0 -C 0.6 -D 0.05 -U 1.0 -M 0.4 -S -1.0 -c -1
Relation: Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-RJ
Instances: 46
Attributes: 8
Score
Against
Position
Innings
Venue
Ground
Date
Result

=== Associator model (full training set) ===

Apriori
=====

Minimum support: 0.4 (18 instances)
Minimum metric <confidence>: 0.6
Number of cycles performed: 12

Generated sets of large itemsets:

Size of set of large itemsets L(1): 6

Size of set of large itemsets L(2): 2

Best rules found:

1. Ground=Home 26 ==> Result=Won 20 <conf:(0.77)> lift:(1.14) lev:(0.05) [2] conv:(1.21)
2. Position=Middle-Order 29 ==> Result=Won 21 <conf:(0.72)> lift:(1.07) lev:(0.03) [1] conv:(1.05)
3. Result=Won 31 ==> Position=Middle-Order 21 <conf:(0.68)> lift:(1.07) lev:(0.03) [1] conv:(1.04)
4. Result=Won 31 ==> Ground=Home 20 <conf:(0.65)> lift:(1.14) lev:(0.05) [2] conv:(1.12)

```

Here, we put the min_confidence = 0.6 and min_sup = 0.4. Then we can got a reliable patterns as Joe Root's century taking patterns.

Here, we can say Joe Root took many centuries in home and most of the matches were won. And he has best performance as a middle-order batsman. He is performing well at home grounds rather than away from the England.

Classifying

Before adding the J48 we will applied zeroR algorithm to the dataset. It is a very simple algorithm and we can compare the accuracy of other algorithms with **zeroR** algorithm.

The screenshot shows the Orange3 software interface with the ZeroR classifier selected. The 'Test options' panel on the left shows 'Cross-validation' selected with 10 folds. The 'Classifier output' panel on the right displays the results of a 10-fold cross-validation. The results include a summary of performance metrics and a detailed accuracy by class table.

Test options:

- Use training set
- Supplied test set
- ☒ Cross-validation Folds: 10
- Percentage split %: 66
- More options...

Classifier output:

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

ZeroR predicts class value: Won

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	31	67.3913 %
Incorrectly Classified Instances	15	32.6087 %
Kappa statistic	0	
Mean absolute error	0.3374	
Root mean squared error	0.4074	
Relative absolute error	100	%
Root relative squared error	100	%
Total Number of Instances	46	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	1.000	1.000	0.674	1.000	0.805	?	0.404	0.632	Won
	0.000	0.000	?	0.000	?	?	0.335	0.144	Drawn
	0.000	0.000	?	0.000	?	?	0.357	0.157	Lost
Weighted Avg.	0.674	0.674	?	0.674	?	?	0.386	0.475	

=== Confusion Matrix ===

```
a b c <-- classified as
31 0 0 | a = Won
7 0 0 | b = Drawn
8 0 0 | c = Lost
```

After applying zeroR algorithm we can see there are 31 correctly classified instances and 15 of incorrectly classified instances. Correctly classified instances as a percentage 67.39% and incorrectly classified instances 32.6%. Accuracy by class the target variable we have used is Result. According to this accuracy we can compare this with other algorithms.

Decision tree (J48)

Now we are going to use J48 algorithm to this dataset. It may classify whether a given match which Joe Root scored century had won or not.

Classifier

Choose J48 -U -M 10

Test options

☐ Use training ...
☐ Supplied test... Set...
☒ Cross-validation... Folds 10
☐ Percentage ... % 66
More options...

(Nom) Result

Start Stop

Result list (right-click for options)

00:41:52 - trees.J48
00:43:09 - trees.J48
00:43:32 - trees.J48
00:45:10 - trees.J48
00:45:25 - trees.J48
00:48:04 - trees.J48
00:50:18 - trees.J48
00:52:56 - trees.J48
00:53:16 - trees.J48
00:53:19 - trees.J48
00:53:33 - trees.J48
00:53:46 - trees.J48
00:55:27 - trees.J48
00:55:56 - trees.J48
00:56:35 - trees.J48
00:57:10 - trees.J48
00:58:07 - trees.J48
00:58:28 - trees.J48
00:58:47 - trees.J48
00:59:29 - trees.J48
00:59:49 - trees.J48
01:00:02 - trees.J48
01:00:24 - trees.J48
01:00:54 - trees.J48

Classifier output

Innings = 4th: Won (2.0)
Number of Leaves : 6
Size of the tree : 8
Time taken to build model: 0 seconds
=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances	29	63.0435 %
Incorrectly Classified Instances	17	36.9565 %
Kappa statistic	-0.0303	
Mean absolute error	0.3422	
Root mean squared error	0.4417	
Relative absolute error	101.4204 %	
Root relative squared error	108.4292 %	
Total Number of Instances	46	

=== Detailed Accuracy By Class ===

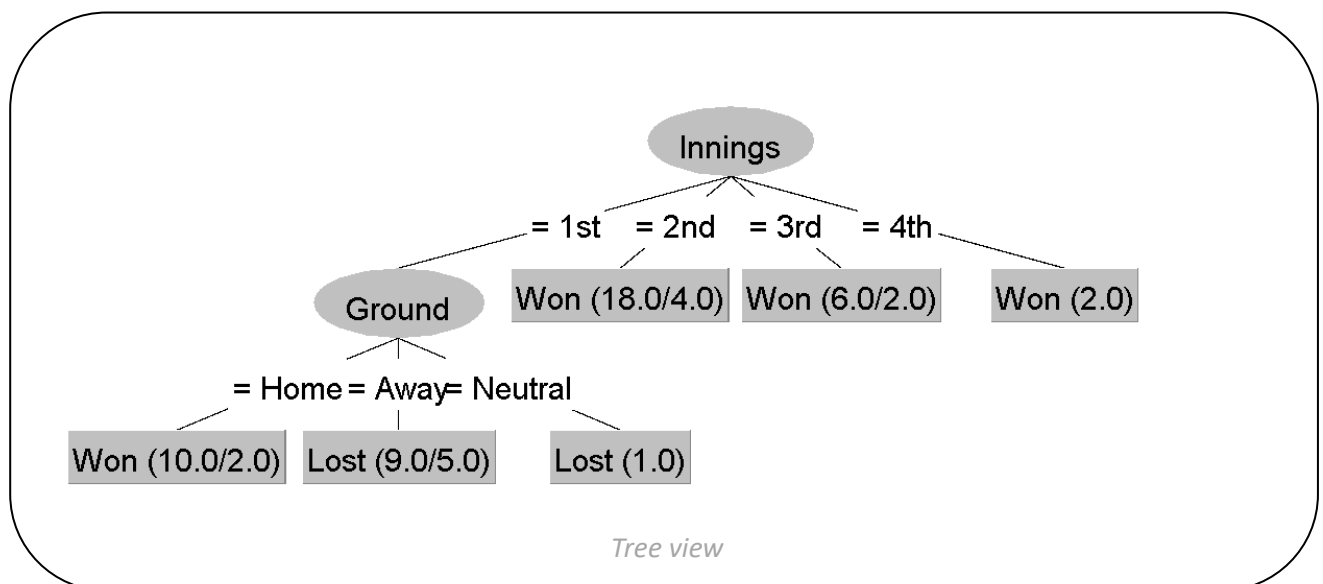
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.935	0.933	0.674	0.935	0.784	0.004	0.415	0.622	Won
	0.000	0.000	?	0.000	?	?	0.311	0.127	Drawn
	0.000	0.079	0.000	0.000	0.000	-0.121	0.536	0.200	Lost
Weighted Avg.	0.630	0.643	?	0.630	?	?	0.420	0.473	

=== Confusion Matrix ===

	a	b	c	<-- classified as
29	0	2		a = Won
6	0	1		b = Drawn
8	0	0		c = Lost

After applied the J48 algorithm among 46 of instances the correctly classified instances are 29 and the percentage is 63.04%. The incorrectly classified instances are 17 and the percentage is 36.95%.

The decision tree view is helps to understand how the algorithm partitions the dataset based on different attributes and creates decision rules for classification or predication tasks.



J48 Percentage split

Now, to the dataset again going to be applied the J48. Here, it will be applied with percentage split. It means this algorithm will be applied to 66% percentage of data from this dataset. 66% is the training and 34% is testing split. The dataset is partitioned 70% of the instances will be using to training the J48 model remaining 34% is used to evaluate the model's performance.

The screenshot shows the WEKA Classifier window. The 'Classifier' dropdown is set to 'J48 -U -M 8'. Under 'Test options', 'Percentage split' is selected with a value of 66%. The 'Classifier output' pane displays the following information:

```
Size of the tree :      8

Time taken to build model: 0 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances      14           87.5 %
Incorrectly Classified Instances     2           12.5 %
Kappa statistic                     0
Mean absolute error                  0.334
Root mean squared error              0.3882
Relative absolute error              100.9542 %
Root relative squared error          107.1686 %
Total Number of Instances           16

=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	1.000	1.000	0.875	1.000	0.933	?	0.500	0.875	Won
	?	0.000	?	?	?	?	?	?	Drawn
	0.000	0.000	?	0.000	?	?	0.500	0.125	Lost
Weighted Avg.	0.875	0.875	?	0.875	?	?	0.500	0.781	

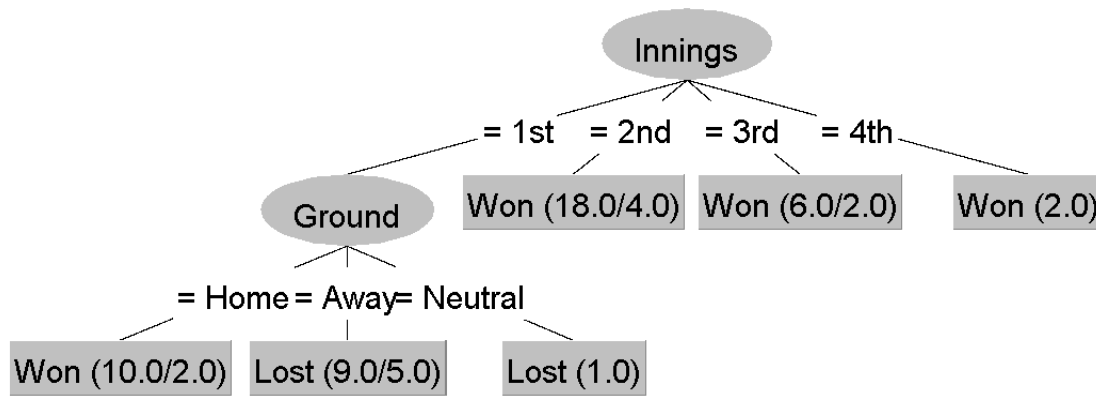
```
=== Confusion Matrix ===

 a  b  c  <-- classified as
14  0  0 | a = Won
 0  0  0 | b = Drawn
 2  0  0 | c = Lost
```

Here, we can see the correctly classified instances are 14 and the percentage is 87.5%. Incorrectly classified instances are 2 the percentage is 12.5%. Total of instances that has been taken to applied the J48 is 16 from 46.

When the accuracy of this result is compared to the previous result the percentage split accuracy has increased than previous one.

The decision tree view will give a clear idea about how the algorithm has worked. Below is the decision tree view of the percentage split J48



Tree view

Naive bayes algorithm

After applying J48 algorithm we applied the naive bayes algorithm to the dataset using training set.

Classifier: Choose **NaiveBayes**

Test options:

- ☒ Use training set
- ☐ Supplied test set Set...
- ☐ Cross-validation Folds 3
- ☐ Percentage split % 66
- More options...

(Nom) Result: Start Stop

Result list (right-click for options):

- 14:04:38 - trees.J48
- 14:04:59 - trees.J48
- 16:00:43 - bayes.NaiveBayes

Classifier output:

[total] 77.0 53.0 54.0

Time taken to build model: 0 seconds

=== Evaluation on training set ===

Time taken to test model on training data: 0.34 seconds

=== Summary ===

Correctly Classified Instances	44	95.6522 %
Incorrectly Classified Instances	2	4.3478 %
Kappa statistic	0.9075	
Mean absolute error	0.1232	
Root mean squared error	0.1828	
Relative absolute error	36.7181 %	
Root relative squared error	45.0959 %	
Total Number of Instances	46	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	1.000	0.133	0.939	1.000	0.969	0.902	0.989	0.995	Won
	0.857	0.000	1.000	0.857	0.923	0.914	1.000	1.000	Drawn
	0.875	0.000	1.000	0.875	0.933	0.923	0.993	0.975	Lost
Weighted Avg.	0.957	0.090	0.959	0.957	0.956	0.908	0.992	0.992	

=== Confusion Matrix ===

```

a b c <-- classified as
31 0 0 | a = Won
1 6 0 | b = Drawn
1 0 7 | c = Lost
  
```

Here the correctly classified instances are 44 and the percentage of it is 95.6522%. The incorrectly classified instances is 2 and the percentage of it is 4.3478%. We can say this is the highest accurate algorithm among the other algorithms that we have applied so far.

Naive bayes using percentage split

Next we applied the naïve bayes algorithms using 66% percentage of training data and the remaining 34% testing data.

Classifier
Choose **NaiveBayes**

Test options
☐ Use training set
☐ Supplied test set Set...
☐ Cross-validation Folds 3
☒ Percentage split % 66
 More options...

(Nom) Result
 Start Stop

Result list (right-click for options)
 14:04:38 - trees.J48
 14:04:59 - trees.J48
 16:00:43 - bayes.NaiveBayes
 16:12:32 - bayes.NaiveBayes

Classifier output

```
[total] 77.0 53.0 54.0
```

Time taken to build model: 0 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances	8	50	%
Incorrectly Classified Instances	8	50	%
Kappa statistic	0.0857		
Mean absolute error	0.3885		
Root mean squared error	0.4669		
Relative absolute error	117.4344	%	
Root relative squared error	128.8901	%	
Total Number of Instances	16		

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.500	0.500	0.875	0.500	0.636	0.000	0.321	0.811	Won
	?	0.375	0.000	?	?	?	?	?	Drawn
	0.500	0.071	0.500	0.500	0.500	0.429	0.464	0.313	Lost
Weighted Avg.	0.500	0.446	0.828	0.500	0.619	0.054	0.339	0.749	

=== Confusion Matrix ===

```
a b c  <-- classified as
7 6 1 | a = Won
0 0 0 | b = Drawn
1 0 1 | c = Lost
```

Here, the both correctly classified instances and incorrectly classified instances are the same 50% and 50%. It has been used only 16 instances from 46. We can see the accuracy of naïve bayes using percentage split is very lower than naïve bayes with training set.

Clustering

To this dataset we are going to apply the clustering algorithm. Here, we will use **SimpleKmeans** algorithm to do the clustering. This will give data exploration by organizing the data into meaningful groups. It will be easy to understand the distribution and relationships between this dataset's attributes.

Firstly we applied the Kmeans to the dataset having 3 clusters. But the squared error was high. And the number of iteration is 4.

Within cluster sum of squared errors: 167.0

The screenshot shows the Orange3 SimpleKmeans widget interface. On the left, the 'Cluster mode' section has 'Use training set' selected. The 'Clusterer output' section on the right shows 'Number of iterations: 4' and 'Within cluster sum of squared errors: 167.0'. Below this, the 'Initial starting points (random):' section lists five clusters with their respective data points. The 'Missing values globally replaced with mean/mode' section shows the final cluster centroids for each attribute. The 'Time taken to build model (full training data) : 0.01 seconds' is displayed. The '==== Model and evaluation on training set ====' section shows the clustered instances.

Attribute	Full Data (46.0)	Cluster# 0 (6.0)
Score	109	109
Against	India	West Indies
Position	Middle-Order	Top-Order
Innings	1st	1st
Venue	Lord's, London	Kensington Oval, Bridgetown
Ground	Home	Away
Date	24-May-13	12-Feb-16
Result	Won	Drawn

Time taken to build model (full training data) : 0.01 seconds

==== Model and evaluation on training set ====

Clustered Instances

Cluster	Count	Percentage
0	6	(13%)

Then we generated many models. After that we could found a model with low squared error. Here we needed to put 45 clusters to reduce the squared error. Number of iterations is 2.

Within cluster sum of squared errors: 3.00

The screenshot shows the WEKA GUI with the SimpleKMeans classifier selected. The 'Cluster mode' section on the left has 'Use training set' selected. The 'Cluster output' pane on the right displays the following information:

```
==== Run information ====
Scheme: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 45 -A "weka.core.EuclideanDistance" -R first-last -I 500 -num-slots 1 -S 10
Relation: Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-RI-weka.filters.unsupervised.attribute.InterquartileRange-Rfilter
Instances: 46
Attributes: 8
Score
Against
Position
Innings
Venue
Ground
Date
Ignored: Result
Test mode: Classes to clusters evaluation on training data

--- Clustering model (full training set) ---

kMeans
=====
Number of iterations: 2
Within cluster sum of squared errors: 3.0000000000000004

Initial starting points (random):
Cluster 0: 109, ' West Indies',Top-Order,3rd,'Sir Vivian Richards Stadium, North Sound',Away,08-Mar-22
Cluster 1: 125, ' South Africa',Top-Order,1st,'SuperSport Park, Centurion',Away,09-Feb-16
Cluster 2: 190, ' South Africa',Middle-Order,1st,'Lord's, London',Home,06-Jul-17
Cluster 3: '180* ', ' India',Middle-Order,2nd,'Lord's, London',Home,12-Aug-21
Cluster 4: 122, ' West Indies',Middle-Order,3rd,'Daren Sammy Cricket Ground, Gros Islet',Away,09-Feb-19
Cluster 5: '200* ', ' Sri Lanka',Middle-Order,1st,'Lord's, London',Home,12-Jun-14
Cluster 6: 176, ' New Zealand',Middle-Order,2nd,'Trent Bridge, Nottingham',Home,10-Jun-22
Cluster 7: 100*, ' India',Top-Order,2nd,'Headingley, Leeds',Home,17-Jul-18
Cluster 8: 153, ' West Indies',Top-Order,1st,'Kensington Oval, Bridgetown',Away,16-Mar-22
```

Clustering with percentage split

Then we applied the clustering to the dataset. We took 60% of data to train and the 40% of remaining data to test. Here, we could reduce the number of clusters to 27 and the squared error was reduced to 0.00. Number of iterations is 2.

The screenshot shows the WEKA GUI with the SimpleKMeans classifier selected. The 'Cluster mode' section on the left has 'Percentage split' selected with a value of 60. The 'Cluster output' pane on the right displays the following information:

```
==== Run information ====
Scheme: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 27 -A "weka.core.EuclideanDistance" -R first-last -I 500 -num-slots 1 -S 10
Relation: Joe-Root-All-International-Cricket-Centuries-weka.filters.unsupervised.attribute.Remove-RI-weka.filters.unsupervised.attribute.InterquartileRange-Rfilter
Instances: 46
Attributes: 8
Score
Against
Position
Innings
Venue
Ground
Date
Ignored: Result
Test mode: Classes to clusters evaluation on training data

--- Clustering model (full training set) ---

kMeans
=====
Number of iterations: 2
Within cluster sum of squared errors: 0.0

Initial starting points (random):
Cluster 0: 125, ' India',Middle-Order,3rd,'The Oval, London',Home,07-Sep-18,Won
Cluster 1: 124, ' Sri Lanka',Middle-Order,3rd,'Pallekele International Cricket Stadium, Kandy',Away,14-Nov-18,Won
Cluster 2: 109, ' South Africa',Top-Order,1st,'Wanderers Stadium, Johannesburg',Away,12-Feb-16,Lost
Cluster 3: '104* ', ' Sri Lanka',Middle-Order,2nd,'Pallekele International Cricket Stadium, Kandy',Away,10-Dec-14,Won
Cluster 4: '182* ', ' West Indies',Middle-Order,2nd,'National Cricket Stadium, St. George's',Away,21-Apr-15,Won
Cluster 5: 109, ' India',Middle-Order,3rd,'Trent Bridge, Nottingham',Home,04-Aug-21,Drawn
Cluster 6: 190, ' South Africa',Middle-Order,1st,'Lord's, London',Home,06-Jul-17,Won
Cluster 7: 113, ' India',Middle-Order,1st,'Headingley, Leeds',Home,05-Sep-14,Won
Cluster 8: 254, ' Pakistan',Top-Order,1st,'Old Trafford, Manchester',Home,22-Jul-16,Won
Cluster 9: 142*, ' India',Middle-Order,4th,'Edgbaston Cricket Ground, Birmingham',Home,01-Jul-22,Won
Cluster 10: 180, ' Australia',Top-Order,3rd,'Lord's, London',Home,19-Jul-19,Won
Cluster 11: 102, ' New Zealand',Top-Order,1st,'University Oval, Dunedin',Away,07-Mar-18,Lost
Cluster 12: 154*, ' India',Middle-Order,2nd,'Trent Bridge, Nottingham',Home,09-Jun-14,Drawn
Cluster 13: 124, ' India',Top-Order,1st,'Saurashtra Cricket Association Stadium, Rajkot',Away,09-Nov-16,Drawn
Cluster 14: 106*, ' New Zealand',Top-Order,2nd,'Trent Bridge, Nottingham',Home,17-Jun-15,Won
Cluster 15: 118*, ' Australia',Middle-Order,1st,'Edgbaston Cricket Ground, Birmingham',Home,16-Jun-23,Lost
Cluster 16: 153, ' West Indies',Top-Order,1st,'Kensington Oval, Bridgetown',Away,16-Mar-22,Drawn
Cluster 17: 134, ' Australia',Middle-Order,1st,'Sophia Gardens, Cardiff',Home,08-Jul-15,Won
Cluster 18: 121, ' Sri Lanka',Middle-Order,1st,'Wellington Regional Stadium, Wellington',Neutral,01-Mar-15,Lost
Cluster 19: 130, ' Australia',Middle-Order,2nd,'Trent Bridge, Nottingham',Home,06-Aug-15,Won
Cluster 20: 226, ' New Zealand',Middle-Order,2nd,'Seddon Park, Hamilton',Away,29-Nov-19,Drawn
Cluster 21: 125, ' South Africa',Top-Order,1st,'SuperSport Park, Centurion',Away,09-Feb-16,Lost
Cluster 22: '100* ', ' West Indies',Top-Order,2nd,'Rose Bowl, Southampton',Home,14-Jun-19,Won
Cluster 23: '180* ', ' India',Middle-Order,2nd,'Lord's, London',Home,12-Aug-21,Lost
Cluster 24: 104, ' New Zealand',Top-Order,1st,'Edgbaston, Birmingham',Home,09-Jun-15,Won
Cluster 25: 176, ' New Zealand',Middle-Order,2nd,'Trent Bridge, Nottingham',Home,10-Jun-22,Won
Cluster 26: '200* ', ' Sri Lanka',Middle-Order,1st,'Lord's, London',Home,12-Jun-14,Drawn

Missing values globally replaced with mean/mode

Final cluster centroids:
```

Then we used **classes to cluster evaluation** to result attribute. We used 3 clusters initially but the incorrectly clustered instances were little bit high (25.00). As a percentage it was 54.3478 %. Then we decided to reduce clusters to 2. Then the incorrectly clustered instances were 18.0.

As a percentage it was 39.1304 %. Here the squared error was high when the number of clusters reduce but the incorrectly clustered instances became low.

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 60

☒ Classes to clusters evaluation

(Nom) Result

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

13:54:04 - SimpleKMeans

13:54:26 - SimpleKMeans

13:54:37 - SimpleKMeans

13:55:14 - SimpleKMeans

13:55:24 - SimpleKMeans

13:55:40 - SimpleKMeans

13:56:05 - SimpleKMeans

13:56:16 - SimpleKMeans

13:56:26 - SimpleKMeans

13:58:18 - SimpleKMeans

13:59:29 - SimpleKMeans

13:59:57 - SimpleKMeans

14:00:40 - SimpleKMeans

14:47:32 - SimpleKMeans

14:47:48 - SimpleKMeans

14:55:23 - SimpleKMeans

14:56:52 - SimpleKMeans

15:20:11 - SimpleKMeans

15:20:32 - SimpleKMeans

15:20:51 - SimpleKMeans

15:21:03 - SimpleKMeans

15:21:21 - SimpleKMeans

15:21:33 - SimpleKMeans

15:21:41 - SimpleKMeans

15:26:33 - SimpleKMeans

15:31:43 - SimpleKMeans

15:32:24 - SimpleKMeans

15:32:35 - SimpleKMeans

15:32:57 - SimpleKMeans

Clusterer output

Score	109	107	104	121
Against	India	West Indies	New Zealand	India
Position	Middle-Order	Middle-Order	Top-Order	Middle-Order
Innings	1st	2nd	1st	1st
Venue	Lord's, London	Trent Bridge, Nottingham	Kensington Oval, Bridgetown	Lord's, London
Ground	Home	Away	Away	Home
Date	24-May-13	05-Mar-14	09-Jun-15	24-May-13

Time taken to build model (full training data) : 0 seconds

=== Model and evaluation on training set ===

Clustered Instances

	0	18 (39%)
1	9 (20%)	
2	19 (41%)	

With 3 clusters

Class attribute: Result

Classes to Clusters:

```

0 1 2 <-- assigned to cluster
14 3 14 | Won
3 2 2 | Drawn
1 4 3 | Lost

Cluster 0 <-- Drawn
Cluster 1 <-- Lost
Cluster 2 <-- Won

```

Out put of the algorithm

Incorrectly clustered instances : 25.0 54.3478 %

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 60

☒ Classes to clusters evaluation

(Nom) Result

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

13:54:04 - SimpleKMeans

13:54:26 - SimpleKMeans

13:54:37 - SimpleKMeans

13:55:14 - SimpleKMeans

13:55:24 - SimpleKMeans

13:55:40 - SimpleKMeans

13:56:05 - SimpleKMeans

13:56:16 - SimpleKMeans

13:56:26 - SimpleKMeans

13:58:18 - SimpleKMeans

13:59:29 - SimpleKMeans

13:59:57 - SimpleKMeans

14:00:40 - SimpleKMeans

14:47:32 - SimpleKMeans

14:47:48 - SimpleKMeans

14:55:23 - SimpleKMeans

14:56:52 - SimpleKMeans

15:20:11 - SimpleKMeans

15:20:32 - SimpleKMeans

15:20:51 - SimpleKMeans

15:21:03 - SimpleKMeans

15:21:21 - SimpleKMeans

15:21:33 - SimpleKMeans

15:21:41 - SimpleKMeans

15:26:33 - SimpleKMeans

15:31:43 - SimpleKMeans

15:32:24 - SimpleKMeans

15:32:35 - SimpleKMeans

15:32:57 - SimpleKMeans

Clusterer output

Score	(46.0)	(29.0)	(17.0)
Against	India	West Indies	New Zealand
Position	Middle-Order	Middle-Order	Middle-Order
Innings	1st	2nd	1st
Venue	Lord's, London	Trent Bridge, Nottingham	Lord's, London
Ground	Home	Home	Home
Date	24-May-13	18-Jul-13	24-May-13

Time taken to build model (full training data) : 0 seconds

=== Model and evaluation on training set ===

Clustered Instances

	0	29 (63%)
1	17 (37%)	

With 2 clusters

Class attribute: Result

Classes to Clusters:

```

0 1 <-- assigned to cluster
22 9 | Won
5 2 | Drawn
2 6 | Lost

Cluster 0 <-- Won
Cluster 1 <-- Lost

```

Out put of the algorithm

Incorrectly clustered instances : 18.0 39.1304 %

Studying above algorithms we can take some most accurate algorithms that worked well for this particular dataset:

- Apriori algorithm
- J48 algorithm with percentage split a
- [Naive bayes algorithm using training set](#)
- clustering with classes to cluster evaluation (result attribute)

Results and patterns

Considering above algorithms applied to the dataset. We can identify Joe Root is performing well at home grounds. And he has taken most of the centuries at the home. He is mostly preferred and performed to play as a middle order batsman. Most of the matches that Joe Root performed well had won by the England. We can identify Joe Root has little bit low performance of matches that is playing away from home. Also, away from home some of matches were lost even Joe Root performed well. From that we can say England team is most comfortable to win the matches at home.