Taimoor khan

Id: `13579

Data Mining

Sir: Zain shukat

| Classification | Regression |
|---|---|
| Classification is the discovery of model or functions where the mapping of objects is done into predefined classes. | Regression is a devised model in which the mapping of objects is done into values. |
| It involves prediction of discrete values | It involves prediction of continuous values |
| Nature of the predicted data is unordered | Nature of the predicted data is ordered |
| Algorithms used are Decision tree, logistic regression, etc. | Algorithms used are Regression tree (Random forest), Linear regression, etc. |
| Method of calculation is measuring accuracy | Method of calculation is Measurement of root mean square error |

Classification Scenario Example:

Suppose from your past data (train data) you come to know that your best friend likes the above movies. Now one new movie (test data) released. Hopefully, you want to know your best friend like it or not. If you strongly confirmed about the chances of your friend like the move.  You can take your friend to a movie this weekend.

If you clearly observe the problem it is just whether your friend like or not. Finding a solution to this type of problem is called as classification. This is because we are classifying the things to their belongings (yes or no, like or dislike). Keep in mind here we are forecasting target class (classification) and the other thing this classification belongs to supervised learning. This is because you are learning this from your train data.

In this case, the problem is a binary classification in which we have to predict whether output belongs to class 1 or class 2 (class 1 : yes, class 2: no ). As we have discussed earlier we can use classification for predicting more classes too. Like ( colour prediction: red,green,blue,yellow,orange).
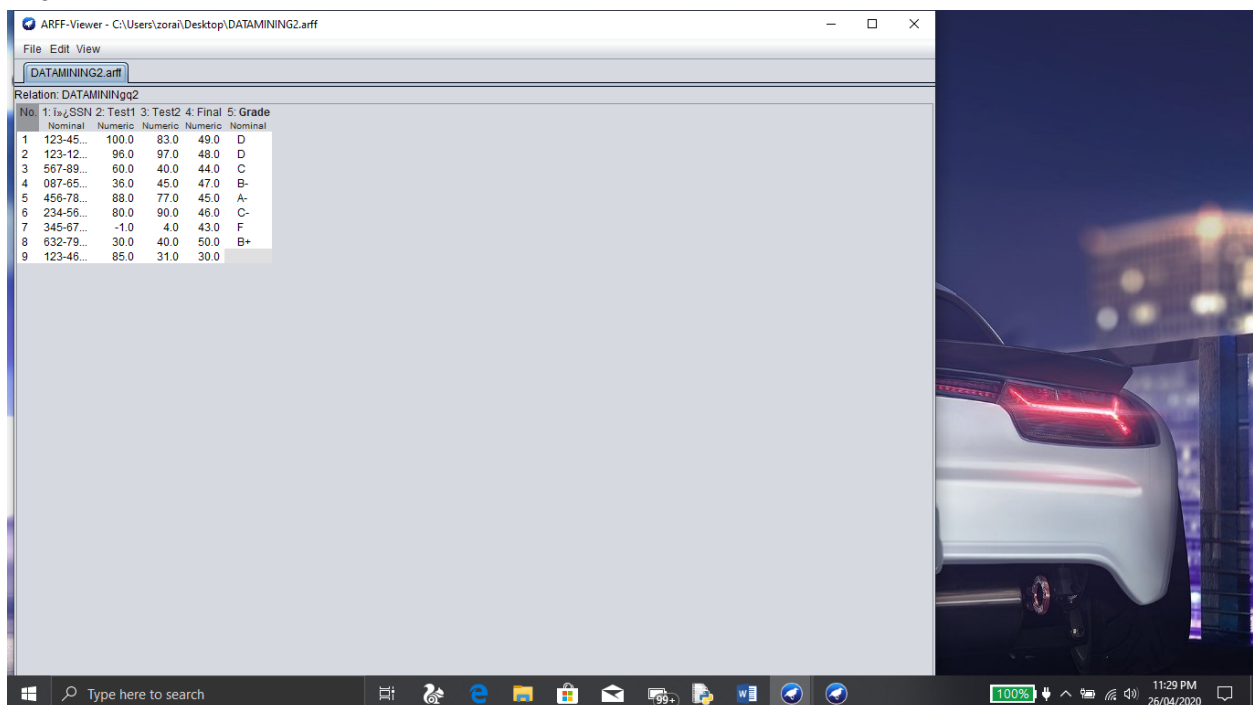
Regression Scenario Example:

Suppose from your past data (train data) you come to know that your best friend likes the above movies. You also know how many times each particular movie seen by your friend. Now one new movie (test data) released. Now you are going to find how many times this newly released movie will your friend watch. It could be 5 times, 6 times,10 times etc

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

If you clearly observe the problem is about finding the count, sometimes we can say this as predicting the value. Keep in mind, here we are forecasting a value (prediction) and the other thing this prediction also belongs to supervised learning. This is because you are learning this from you train data.

Q2

Ans:



| No. | 1: ï»¿SSN | 2: Test1 | 3: Test2 | 4: Final | 5: Grade |
|-----|-----------|----------|----------|----------|----------|
|     | Nominal   | Numeric  | Numeric  | Numeric  | Nominal  |
| 1   | 123-45... | 100.0    | 83.0     | 49.0     | D        |
| 2   | 123-12... | 96.0     | 97.0     | 48.0     | D        |
| 3   | 567-89... | 60.0     | 40.0     | 44.0     | C        |
| 4   | 087-65... | 36.0     | 45.0     | 47.0     | B-       |
| 5   | 456-78... | 88.0     | 77.0     | 45.0     | A-       |
| 6   | 234-56... | 80.0     | 90.0     | 46.0     | C-       |
| 7   | 345-67... | -1.0     | 4.0      | 43.0     | F        |
| 8   | 632-79... | 30.0     | 40.0     | 50.0     | B+       |
| 9   | 123-46... | 85.0     | 31.0     | 30.0     |          |

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat



Q3

Ans:

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat

Name M Zoraiz Ali
Id 14413
Data mining
Sir Zain Shaukat