

Handling Numerical Data

Q How to Split the Data Set - I?

SNo	User_Rating	Downloaded
1	3.5	Yes
2	4.6	Yes
3	2.2	No
4	1.6	Yes
5	4.1	No
6	3.9	No
7	3.2	No
8	2.9	Yes
9	4.8	Yes
10	3.5	No

(Before) DataSet - I

Ans As, earlier we have Categorical Data. So It is easier to categorise the data. But, Now we have the Numerical Data. It is difficult to split as per the Numeric Values...

Q To Split the data when there is Numerical Data, what should be the steps?
 Ans Step 1:- Sort the DataSet as per the Numerical Column.

SNo	User_Rating	Downloaded
1	1.6	Yes
2	2.2	No
3	2.9	Yes
4	3.2	No
5	3.3	No
6	3.5	Yes
7	3.9	No
8	4.1	No
9	4.6	Yes
10	4.8	Yes

(After)

DataSet - I

Spiral

Date.....

Step 2:- Split the entire data as per the individual user-rating.
For example:-

User_Rating ≥ 1.6

1: \rightarrow User_Rating ≤ 1.6

SNo	User_Rating	Downloaded
1	1.6	Yes

Dataset_ First

SNo	User_Rating	Downloaded
2	2.2	No
3	2.9	Yes
4	3.2	No
5	3.3	No
6	3.5	Yes
7	3.9	No
8	4.1	No
9	4.6	Yes
10	4.8	Yes

Dataset_Second

2: User_Rating ≤ 2.2

SNo	User_Rating	Downloaded
1	1.6	Yes
2	2.2	No

Rest Complete Table..

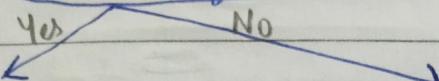
3: Splitting with all the values, till \rightarrow User_Rating ≤ 4.8

Complete table

Date.....

Step 3:- For every DataSet we get after Step 2 From each User_Rating Value
For example:-

1. $(User_Rating \leq 1.6)$



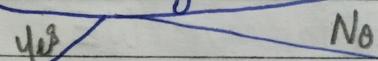
Dataset_First $\Rightarrow E(D_{-First})$ Dataset_Second $\Rightarrow E(D_{-Second})$

2. $(User_Rating \leq 2.2)$



Dataset_First $\Rightarrow E(D_{-First})$ Dataset_Second $\Rightarrow E(D_{-Second})$

3. $(User_Rating \leq 4.8)$



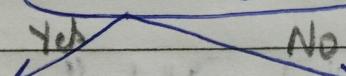
Dataset_First $\Rightarrow E(D_{-First})$

Dataset_Second $\Rightarrow E(D_{-Second})$

Step 4:- Calculate the Weighted Entropy of each of the Dataset.

For example:-

1. $(User_Rating \leq 1.6)$



Dataset_First



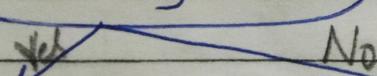
$E(D_{-First})$

Dataset_Second



WE_1

2. $(User_Rating \leq 2.2)$



$E(D_{-First})$

$E(D_{-Second})$

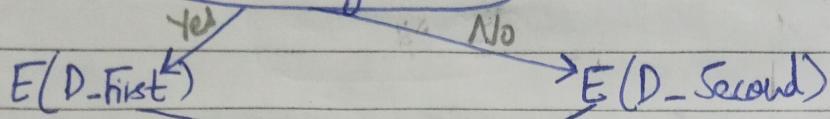
WE_2

Date.....

Step 5:- Calculate the Information Gain of every dataset.

For example:-

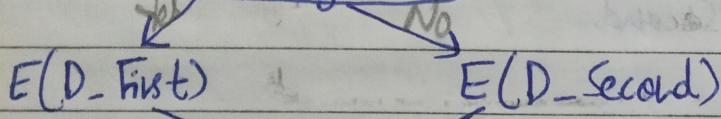
1. $(User_Rating \leq 1.6)$



$$WE_1 \implies \text{Information} = E(\text{Parent}) - WE_1$$

Gain 1

2. $(User_Rating \leq 2.2)$



$$WE_2 \implies \text{Information} = E(\text{Parent}) - WE_2$$

Gain 2

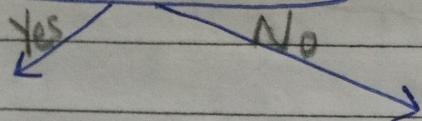
Step 6:- Find the Maximum Information Gain Value..

For example:-

$$\text{MAX}\{IG_1, IG_2, IG_3, \dots, IG_n\}$$

Assuming that " IG_3 " is the Maximum Value..

$User_Rating = IG_3$



"Till you get all the Leaf Nodes."