# DATA PREPROCESSING (Replace Missing Value).

**Aim:** To implement data preprocessing for missing numerical value and finding the error using point estimation method.

## Definition:

### Point estimation:

- estimate a population parameter.
- May be made by calculating the parameter for a simple.
- may be used to predict value for missing data.

### Estimation Error:

Bias difference between expected value and actual value.

$$bias = \epsilon(\theta^\wedge) - \theta.$$

### Mean Squared Error:

expected value of the ssuared difference between the estimate and the actual value

$$MSE = \epsilon(\theta^\wedge - \theta)2.$$

## Algorithm:

Step 1: open the Excel sheet

step 2: create 10 records of student mark list.

step 3: Delete some data in that list and save as missing data.csv

step 4: open the weka tool and open the csv file.

step 5: choose filter - unsupervised in that attribute select replacmissing values then click apply

step 6: Then it will replace the missing values

step 7: Then using the formula for find the prediction accuracy and error accuracy

* prediction accuracy = $\dfrac{\text{observed value}}{\text{prediced value} * 100}$

* error Accuracy = error / predicted value * 100

### Calculation:

⇒ prediction accuracy = observed value / predicted value * 100

$$= 419.4 / 396 * 100$$

$$= 94.4$$

⇒ error Accuracy = error / predicted value * 100

$$= -23.4 / 419.4 * 100$$

$$= 5.6$$

Good

Result: The program executed successfully.

# Experiment - 2

## Data PREPROCESSING FILTERS APPLIED FOR unsupervised ATTRIBUTE

**Aim:** To convert Nominal to Binary values.

**Algorithm:**

**step1:** open weka explorer window

**step2:** click on preprocess tab

**step3:** open the file or dataset (eg. weather . arff) which is available in your computer after installing weka tool.

**step4:** observe the datatype of the features or an attributes on your selected dataset. click edit and see the file content

**step5:** click the data preprocessing and choose the filter unspervised folder, in that choose the attributes, then choose the nominal to binary option.

**step6:** APPly the converted work and save the file.

**step7:** click the edit and check the changed attribute with nominal to binary.

**Result:** Thus the conversion of data type from nomial to binary values is implemented successfully.

# Experiment-3

## Data Processing – Add Expression.

**Aim:** To Add the expressions x and y using the weka tool.

**Algorithm:**

**step1:** open excel and create table with x & y variables and save in csv format.

**step2:** open weka tool Explorer and open the csv file.

**step3:** And Click choose filter and click on the unsupervised then attribute and open the select the Add Expression.

**step4:** Then apply filter then go to show properties then give name z and $(a^2+b^2)*2$.

**step5:** Then you get the z value

**step6:** To see the z value open the edit option.

**Result:**

The program executed and Added the expression successfully.

# Experiment - 4

## Attribute selection

**Aim :** To select the particular attribute best first

**procedure :**

**step1 :** open weka explorer window

**step2 :** Go to open file and goto this PC then click on program files select the weka tool select the iris.

**step3 :** Now choose filter and supervised and click on the Attribute selection.

**step4 :** Then show properties select the best first

**step5 :** The best are shown on the screee (Attribute)

**Result :**

The program executed successfully.

9/5/2023

## Linear Regression

**Aim:** To find a linear regression equation and predict the salary college graduates whose experience is 10 years.

**calculation:**

| x | y | $x_i - \bar{x}$ | $y_i - \bar{y}$ | $(x_i - \bar{x})^2$ | $(y_i - \bar{y})^2$ | $(x_i - \bar{x}) * (y_i - \bar{y})$ |
|---|---|---|---|---|---|---|
| 3 | 80 | −53 | −560 | 2809 | 313600 | 29680 |
| 4 | 45 | −52 | −595 | 2704 | 354025 | 30940 |
| 5 | 60 | −51 | −580 | 2601 | 336400 | 29580 |
| 7 | 75 | −49 | −565 | 2401 | 319225 | 27685 |
| 6 | 55 | −50 | −585 | 2500 | 342225 | 29250 |
| 2 | 20 | −54 | −620 | 2916 | 384400 | 33650 |
| 1 | 10 | −55 | −630 | 3025 | 396900 | 34650 |
| 8 | 85 | −48 | −555 | 2304 | 308025 | 26648 |
| 9 | 95 | −47 | −545 | 2209 | 297025 | 25615 |
| 11 | 115 | −45 | −525 | 2205 | 275625 | 23625 |

**Algorithm:**

**Step 1:** open exel sheet and assign x & y

**Step 2:** Next give some values of table in exel sheets

**Step 3:** Do manual calculation in exel sheet for x mean, y mean $x_o \times y_i - y_i$  $x_i - x^2$, $y_i - m =$ and $(x_1 - x) \times (y_i - y)$

step4: check the Answers

step5: open a New file in csv format and the enter
values in excel sheet

step6: save as salam in Desktop

step7: open wekatool and open file and choose classify
and choose linear Regression.

step8: By click use training set and start

step9: check the Answers stop the program

$$b_1 = \sum [(x_i - x) - (y_i - y)] \left[ \varepsilon (x_1 - x)^2 \right]$$

$$= 97 / 1924$$

$$= 10.5087$$

$$b_2 = y - b_1 * y$$

$$= 59 - 10.5087 * 5.6$$

$$= 0.1515$$

$$y = 0.1515 + 10.5087 \times 10$$

$$= 105.24$$

Result:

Thus the program executed successfully.

cal
10/5/2023

# Experiment - 6

## Classification using Naive Bayessian Classifier Algorithm - (multi class classification problem).

**Aim:** To find the classification prediction using naive bayesian classifier algorithm the new person will buy a comper or not.

**Algorithm:**

step 1: create a table in the excel sheet and save as a csv file.

step 2: Then later open the weka tool and open the computer .csv file

step 3: we need save as computer arff

step 4: later goto edit isee the viewer table iselected all the rows and delete it.

step 5: And then Now click add instance. again open the computer .arff file then open classify in the menu bar.

step 6: click bayes and click naive bayes click cross validation and click start.

step 7: we will get Answer of correctly classified Instances and incorrectly classified instances.

| Age | Income | student | credit-rating | Buys-computer |
|---|---|---|---|---|
| <=30 | High | NO | Fair | No |
| <=30 | High | NO | encuient | yes |
| 31...40 | High | NO | Fair | yes |
| >40 | medinum | NO | Fair | yes |
| >40 | Low | yes | Fair | yes |
| >40 | Low | yes | exclerent | yes No |
| 31-40 | Low | yes | exclunt | yes |
| <=30 | Medimun | NO | Fair | NO |
| <=30 | Low | yes | Fair | NO |
| >40 | medium | yes | Fair | yes |
| <=30 | Medium | yes | excellent | yes |
| 31-40 | Medium | NO | excellent | yes |
| 31-40 | High | yes | Fair | yes |
| >40 | Medinum | NO | Excelunt | NO |

E = age <=30, income =medium, student = yes, credit rating =fair

E1 = age <=30

e2 = Income = Medium

e3 = student = yes

E4 = credit rating = fair

$$P(yes/e) \frac{P(E_1/yes) P(E_2/yes) P(E_3/yes) P(E_4/yes)}{P(E)}$$

P(yes) = 9/14 = 0.643   |   P(no) = 5/14 = 0.357

P(E1/yes) = 2/9 = 0.222   |   P(E1/No) = 3/5 = 0.6

P(E2/yes) = 4/9 = 0.444   |   P(E2/No) = 2/5 = 0.4

P(E3/yes) = 6/9 = 0.667   |   P(E3/No) = 1/5 = 0.2

P(E4/yes) = 6/9 = 0.667   |   P(E4/No) = 2/5 = 0.4

$$P(yes/e) = \frac{0.222 * 0.444 * 0.6667 * 0.667 * 0.643}{P(E)} \quad \frac{0.0028}{P(E)}$$

$$P(no/e) = \frac{0.6 * 0.4 * 0.2 * 0.4 * 0.357}{P(e)} \quad \frac{0.006}{P(E)}$$

# Experiment - 7

## Naive Bayes Classifier Algorithm

Working of Naive Bayer's classifier (single class classification)

**Aim:** To find the classification prediction using naive bayesian classifier Algorithm on weather conditions whether player should play or not.

**Algorithm:**

**step1:** Create a table in xshell sheet and save file as csv file.

**step2:** Then open weka tool, select weather nominal data set

**step3:** keep only outlook and play attribute, remove other attributes save the file as player.

**step4:** goto edit, in the viewer delete all the records.

**step5:** And click add, instance, then add only sunny, after save the new file as playertest.

**step6:** Again the open player file, click classify, again click bayes then naivebayes click the cross folder then click start.

**step7:** Note the corrected and incorrected values.

**step8:** Now click the supply test set, open the playertest file, then click more options setup output predictions as on system

**step9:** Then click start, find out the prediction result with use of right click.

| S.NO | outlook | Play |
|------|---------|------|
| 1 | Sunny | no |
| 2 | Sunny | no |
| 3 | overcast | yes |
| 4 | rainy | yes |
| 5 | rainy | yes |
| 6 | rainy | no |
| 7 | overcast | yes |
| 8 | sunny | no |
| 9 | sunny | yes |
| 10 | rainy | yes |
| 11 | sunny | yes |
| 12 | overcast | yes |
| 13 | overcast | yes |
| 14 | rainy | no |

weather conditions

| weather | Yes | NO |
|---------|-----|-----|
| overcast | 4 | 0 |
| Rainy | 3 | 2 |
| Sunny | 2 | 3 |
| Total. | 9 | 5 |

Applying Bayes theorem.

$P(Yes| sunny) = P(Sunny|Yes) * P(Yes/P(Sunny)$

$P(Sunny) = 2/9 = 0.22$

$P(Sunny) = 0.36$

$P(Yes) = 0.64$

$P(Yes|Sunny)$

$= 0.22 * 0.64/0.36$

$= 0.39$

$P(No|Sunny) = P(Sunny|No) * P(No)/P(Sunny)$

$P(Sunny|No) = 3/5 = 0.60$

$P(No) = 0.36$

$P(Sunny) = 0.36$

So $P(No| Sunny) = 0.60 * 0.36/0.36$

$= 0.60$

$P(A|B) = \dfrac{P(A \cap B)}{P(B)}$

Result: Thus the program executed successfully, hence on sunny day, player cannot play the game.

# Experiment-8

## DESISION TREE ALGORITHM

**Aim:** Given the training data, build a decision tree and predict the class of the following new. ex: age<=30, income=medium, student=yes

## Algorithm:

**step1:** Create csv file with name of computerbuy with above mentioned.

**step2:** Save as arff file in weka tool.

**step3:** open classify and choose rules in that choose decision tree with name of J48.

**step4:** start the program.

**step5:** print the visualize tree.

**step6:** Apply supply set (open computerbuy, go to edit, delete all instance, add new instance age<=30, income= medium, student=yes, buys-computer empty) save it another name.

**step7:** Run with supply set and find the answer.

| Age | Income | student | Buys-computer |
|---|---|---|---|
| <=30 | High | No | NO |
| <=30 | High | NO | NO |
| 31~40 | High | NO | Yes |
| >40 | medium | NO | Yes |
| >40 | Low | Yes | Yes |
| >40 | Low | Yes | NO |
| 31~40 | Low | Yes | yes |
| <=30 | Medimun | NO | NO |
| <=30 | Low medium | Yes | yes |
| >40 | medium | Yes | yes |
| <=30 | medium | Yes | |

| | | | |
|---|---|---|---|
| 31....40 | medium | NO | yes |
| 31--40 | High | yes | yes |
| >40 | medium | NO | NO |

★ Gain $(D, A) = Entropy (D) - \sum_{j=1}^{y} \frac{|D_j|}{|D|} Entropy (D_j)$

$I(syes, 5No) = I(9, 5) = -9/14 \log_2 (9/14) - 5/14 \log_2 (5/14) = 0.94$

Entropy (age) $= 8/14$

$= 0.6935$

Gain (age) = income low (3yes 1 no)

entropy (income) =

$= 0.285714 + 0.393428 + 0.231714$

$= 0.9108$

Gain (income) =

$0.94 - 0.9108$

$= 0.0292$



age

<=30          >40

| Income | student | class |
|---|---|---|
| high | nu | NO |
| high | h | NO |
| medium | no | yes |
| low | yes | yes |
| medium | yes | yes |

| Income | student | class |
|---|---|---|
| medium | no | Yes |
| low | yes | ves |
| low | yes | NO |
| medium | yes | yes |
| medium | no | NO |

| Income | student | class |
|---|---|---|
| high | no | Yes |
| low | yes | Yes |
| medium | no | Yes |
| high | yes | yes |

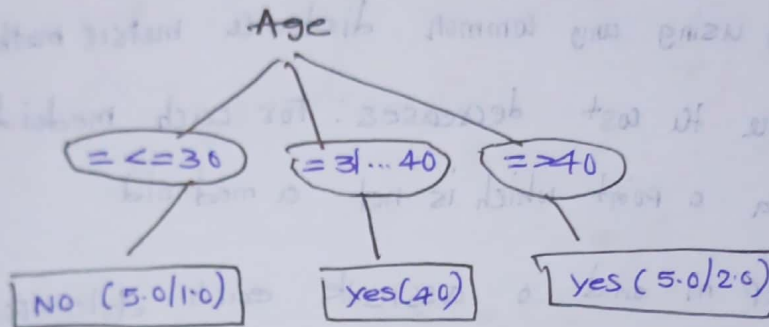Again the same process is needed for the other branch of age :

Entropy (income)
$$= 3/5 (0.9182) + 2/5 (1) = 0.55 + 0.4 = 0.95$$

Gain (income) $= 0.97 - 0.95 = 0.02$

Entropy (student) $= 0.95$

Gain (student) $= 0.97 - 0.95$
$$= 0.02$$

Age

= <= 30    (= 31...40)    = >40

NO (5.0/1.0)    yes(40)    yes (5.0/2.0)

Result : Thus the program executed successfully , Hence person will not buy a computer

Octal
13/6/2023

# Experiment-9

## K-Medoids Algorithm.

**Aim:** To prove the k-medoids Algorithm using weka tool.

**Algorithm:**

**Step1:** initialize select k random points out of the n data points as the medoids.

**Step2:** Associate each data point to the closet medoid by using any common distance metric methods.

**Step3:** While the cost decreases: For each medoid m, for each data o point which is not a medoid.

**Step4:** swap m and o associate each data point to the closet medoid, and recompute the cost.

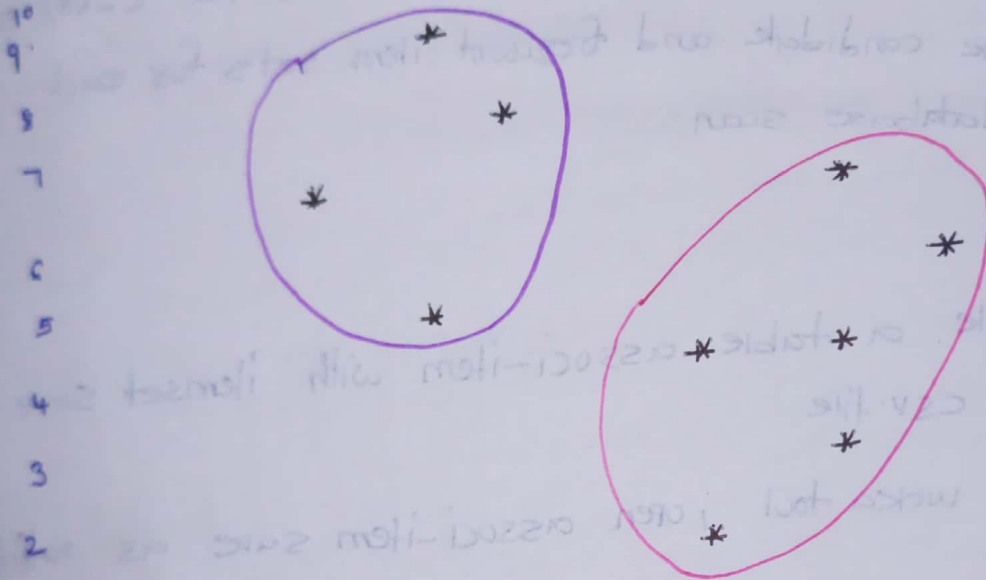**Step5:** save the csv in enecel sheet and enecute in the wek tool.

| X | Y | Dissimidarity C1 | Dissimilarity C2 |
|---|---|---|---|
| 8 | 7 | 6 | 2 |
| 3 | 7 | 3 | 7 |
| 4 | 9 | 4 | 8 |
| 9 | 6 | 6 | 2 |
| 8 | 5 | - | - |
| 5 | 8 | 4 | 6 |
| 7 | 3 | 5 | 3 |
| 8 | 4 | 5 | 1 |
| 7 | 5 | 3 | 1 |
| 4 | 5 | - | - |

The cost in k-medoids algorithm is given as

$$c = \sum_{c_i} \sum_{p_i \in c_i} |p_i - c_i|$$

That formula tell that $\text{Distance} = |x_1 - x_2| + |y_1 - y_2|$.

* $O(k * (n-k)^2)$

# Experiment — 10

## ASSOCIATION RULE MINING
## APRIORI ALGORITHM.

**Aim:** Trace the result of using the apriori algorithm on the grocery store example, with support threshold s = 33·34% and confidence threshold c = 60% show the candidate and frequent item sets for each database scan.

**Algorithm:**

**step1:** create a table associ-item with itemset save as csv·file

**step2:** open weka tool, open associ-item save as arff.file.

**step3:** open association, click apriori algorithm.

**step4:** click show properties, change support 0·33 and confidence 0·60

**step5:** start the program

**step6:** result will display.

| Transaction ID | Items |
|---|---|
| T1 | HotDogs , Buns, Ketchup |
| T2 | Hotdogs , Buns |
| T3 | HotDogs , coke , chips |
| T4 | chips , coke |
| T5 | chips , ketchup |
| T6 | HotDogs, coke , chips |

confidence: The confidence of a rule is

$$\text{conf}(x \to y) = \text{supp}(x \cup y)/\text{supp}(x)$$

| Transcation ID | hotdogs | buns | ketchup | coke | chips |
|---|---|---|---|---|---|
| T1 | T | T | T | | |
| T2 | T | + | | | |
| T3 | + | | | T | T |
| T4 | | | | + | T |
| T5 | | | | T | T |
| T6 | T | | | T | T |

calculation:

support threshold = 33.34%.

⟹ threshod is at least 2 transcations

confidence = 0.60

Association rules:

| Itemset | support | confidence |
|---|---|---|
| Hot Dogs, Buns | 2/6 = 33.33 | 2/4 = 50 |
| Buns, Hot Dogs | 2/6 = 33.33 | 2/2 = 1000 |
| Hotdogs, coke | 2/6 = 33.33 | 2/4 = 50 |
| Coke, Hot Dogs | 2/6 = 33.33 | 2/3 = 66.66 |
| Hot Dogs, chips | 2/6 = 33.33 | 2/4 = 50 |
| chips, Hot Dogs | 2/6 = 33.33 | 2/4 = 50 |
| Coke, chips | 3/6 = 50 | 3/3 = 100 |
| chips, coke | 3/6 = 50 | 3/4 = 75 |
| Hotdogs → coke n chips | 2/6 = 33.33 | 2/4 = 50 |
| coke → chips n Hot dogs | 2/6 = 33.33 | 2/3 = 66.66 |

chips -> Hotdogs n coke       2/6 = 33.33       2/4 = 50

Hotdogs n coke -> chips       2/6 = 33.33       2/2 = 1

chips n Hotdogs -> coke       2/6 = 33.33       2/2 = 1

coke n chips -> Hotdogs       2/6 = 33.37       2/3 = 66.66

**Result:** Thus the program have been executed successfully through Association rule mining.

Generated sets of large itemsets:

→ size of set of large items L(1) : 11

size of set of large itemsets L(2) : 22

size of set of large itemsets L(3) : 14

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

Best rules found:

coke = T    3 ==> chips = T 3 con

buns = T    2 ==> hotdogs = T 2 conf : (1)
Mint

hotdogs = T chips = T 2 ==> coke = T 2 conf : (1)

Minimum support : 0.25 (1 instances)

Minimum metric < confidence >: 0.9

Number of cycles performed 15

# Experiment - 11

Classic agglomerative, hierarchical clustering methods using with linkage criteria.

**Aim:** To find the cluster using hierarchical Agglomerative cluster with linkage criteria and find dendrogram.

**Algorithm:**

step1: open contact lens file in the weka tool.

step2: click cluster and choose hierarchical algorithm. change the property link type with single

step3: start the program. find the results and dendrogram.

step4: change link type completed and start the program, find the results and dendrogram.

step5: change link type average and start the program find the results and dendrogram.

single linkage =

$$L(r,s) = min (D(x_{ri}, x_{sj}))$$

Average linkage =

$$L(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$

complete linkage:

$$L(r,s) = max (D(x_{ri}, x_{sj}))$$

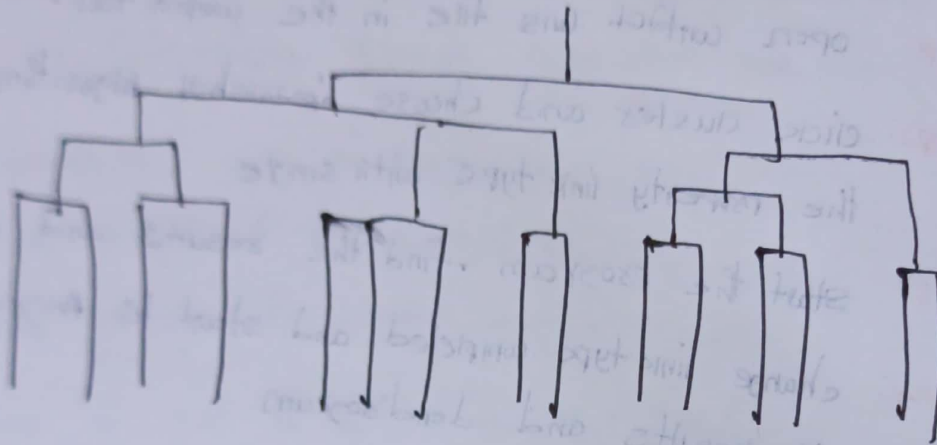Euclidean Distance.

$$X = (a, b) \text{ and } Y = (c, d)$$

The euclidean distance between x and y
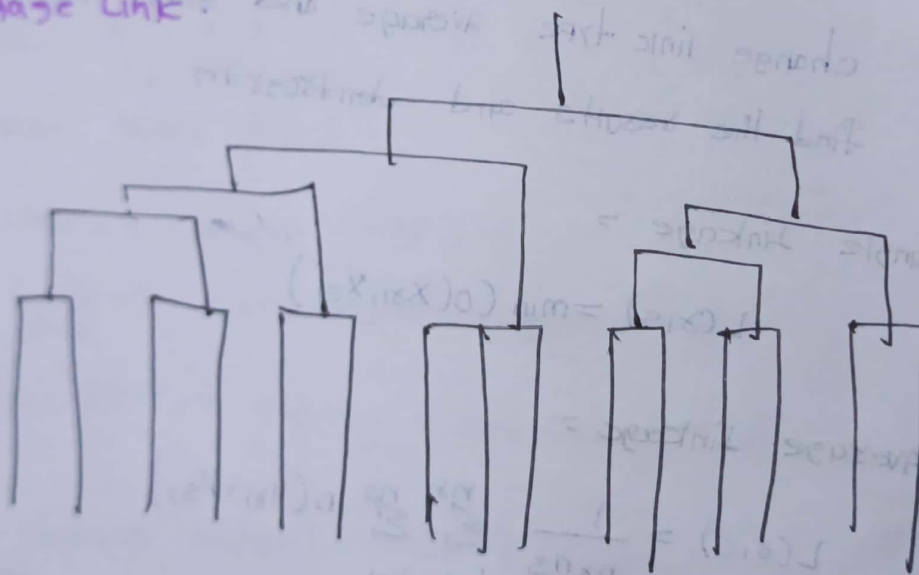
$$\sqrt{(a-c)^2 + (b-d)^2}$$

comparision of link.

| Link | Tot recordz+ | cluster 0 | cluster 1 |
|------|------|------|------|
| single | 20-4 | 83 | 4 |
| completed | 15-9 | 23 | 37 |
| Average | 15-9 | 63 | 37 |
|  |  |  | 37 |

complete Link:



Average Link:



Result: Thus the program executed successfuly
through classic Agglomerative hierarchical
clustering,

Croof
16/05/2023

# Experiment -12

## comparison various algorithms in classification

**Aim:** To briefly described three algorithms in terms of how it works, key algorith parameters will be hightened and the algorithm will be demostrated in the weka explorer interfac.

→ Navie Bayes
→ Decision Tree
→ k-Nearest Neighbors.

**Algorithm:**

step1: open the weka GUI chooser

step2: click the explorer button to open the weka explorer.

step3: Load the Ionosphere dataset from data/ionosphere.arff file.

step4: click "classify to open the classify tab.

### Decision tree algorithm

step1: click the "choose" button and select "REPTree" under the "tree" group

step2: click on the name of the algorithm to riview the algorithm configuration.

step3: click "ok" to close the algorithm configuration.

step4: click the "start" button to run the algorithm on the Ionospha dataset.

### k-Nearest Neighbors algorithm.

step1: click the "choose" button and slect "IBk" under the "Iazy'l group

step2: click on the name of the algorithm to review the agorithm confisuratior.

step3: click "ok" to cose the algorithm configuration.

step4: click the "start" button to run the algorithm on the Ionosphere dataset.

| Algorithm | Accuracy |
|-----------|----------|
| Navie bayes | 82 |
| decision tree | 89 |
| k NN | 86 |

**Result:** Thus the program was executed successfully using weka tool.

Croal
16/5/2022

Experiment - 13

FP GROWTH ALGORITHM USING WEKA.

Aim: To briefly describe about the FP Growth Algorithm using weka and enplosed in the weka tool.

Algorithm:

Step1: open the data file in weka Enpluses.

Step2: It is presumed that the resuired data fields have been discretized. In this example it is age attoibute.

Step3: clicking on the associate. tab will bring up the interface for association rule algorithm

Step4: we will use FP - growth algorithm. This is the default algorithm.

Step5: Inorder to change the parameters for the run (example, support, confidence etc).

Step6: we click on the text box immediately to the right of the choose buton.

Data set:
shopping. arff
@ relation shopping
@ attoibute milk {yes,no}
@ attribute bread {yes,no}
@ attribute honey {yes,no}
@ attribute ghee {yes,no}
@ attribute jam {yes,no}
@ data
yes, yes, no, no, yes
yes, no, yes, no

Yes, yes, no, yes, no

yes, no, yes, no, no

no, yes, yes, no, no

yes, no, yes, no, no

yes, yes, yes, no, yes

yes, yes, yes, no, no

**Result:** Thus the program was executed successfully using weka tool.

Crook-
18/5/2023