



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 1

- Write a R program for the intervals and corresponding frequencies are as follows. age frequency
 1-5. 300
 5-15 550
 15-20 500
 20-50 1500
 50-80 800
 80-110 74

Compute an approximate median value for the data

- Write a R program to calculate co variance and corelation Children of three ages are asked to indicate their preference for three photographs of adults. Do the data suggest that there is a significant relationship between age and photograph preference? What is wrong with this study?

Photograph:

| Age of child | A | B | C |
|--------------|----|----|----|
| 5-6 years: | 18 | 22 | 20 |
| 7-8 years: | 2 | 28 | 40 |
| 9-10 years: | 20 | 10 | 40 |

- Use cov() to calculate the sample covariance between B and C.
 - Use another call to cov() to calculate the sample covariance matrix for the preferences.
 - Use cor() to calculate the sample correlation between B and C.
 - Use another call to cor() to calculate the sample correlation matrix for the preferences.
- Implement using WEKA for the given the data set and perform the Apriori Algorithm and FP algorithm support:3 and confidence=50%

| Customer ID | Transaction ID | Items Bought |
|-------------|----------------|--------------|
| 1 | 0001 | {a, d, e} |
| 1 | 0024 | {a, b, c, e} |
| 2 | 0012 | {a, b, d, e} |
| 2 | 0031 | {a, c, d, e} |
| 3 | 0015 | {b, c, e} |
| 3 | 0022 | {b, d, e} |
| 4 | 0029 | {c, d} |
| 4 | 0040 | {a, b, c} |
| 5 | 0033 | {a, d, e} |
| 5 | 0038 | {a, b, e} |

- Implement using WEKA for the given the data set and perform Bayes classification and descion tree (using training and test data)

| RID | age | income | student | credit_rating | Class: buys_computer |
|-----|-------------|--------|---------|---------------|----------------------|
| 1 | <=30 | high | no | fair | no |
| 2 | <=30 | high | no | excellent | no |
| 3 | 31 . . . 40 | high | no | fair | yes |
| 4 | >40 | medium | no | fair | yes |
| 5 | >40 | low | yes | fair | yes |
| 6 | >40 | low | yes | excellent | no |
| 7 | 31 . . . 40 | low | yes | excellent | yes |
| 8 | <=30 | medium | no | fair | no |
| 9 | <=30 | low | yes | fair | yes |
| 10 | >40 | medium | yes | fair | yes |
| 11 | <=30 | medium | yes | excellent | yes |
| 12 | 31 . . . 40 | medium | no | excellent | yes |
| 13 | 31 . . . 40 | high | yes | fair | yes |
| 14 | >40 | medium | no | excellent | no |



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 2

1. Write a R program for the data to analysis the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.

(a) What is the mean of the data? What is the median?

(b) What is the mode of the data? Comment on the data's modality (i.e., bimodal, trimodal, etc.).

(c) What is the midrange of the data?

(d) Can you find (roughly) the first quartile (Q1) and the third quartile (Q3) of the data?

2. Write a R program for the data that you have selected data from the All Electronics data warehouse for analysis. The data set will be huge! The following data are a list of All Electronics prices for commonly sold items (rounded to the nearest dollar). The numbers have been sorted:

1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.

(i) dataset using an equal-frequency partitioning method with bin equal to 3

(ii) apply data smoothing using bin means and bin boundary.

(iii) Plot Histogram for the above frequency division

3. Implement using WEKA for the given the data set and perform the Apriori Algorithm and FP algorithm support:3 and confidence=50%

Consider the market basket transactions shown in the above table.

(a) What is the maximum number of association rules that can be extracted from this data (including rules that have zero support)?

(b) What is the maximum size of frequent itemsets that can be extracted (assuming minsup > 0)?

| Transaction ID | Items Bought |
|----------------|--------------------------------|
| 1 | {Milk, Beer, Diapers} |
| 2 | {Bread, Butter, Milk} |
| 3 | {Milk, Diapers, Cookies} |
| 4 | {Bread, Butter, Cookies} |
| 5 | {Beer, Cookies, Diapers} |
| 6 | {Milk, Diapers, Bread, Butter} |
| 7 | {Bread, Butter, Diapers} |
| 8 | {Beer, Diapers} |
| 9 | {Milk, Diapers, Bread, Butter} |
| 10 | {Beer, Cookies} |

4. Create the following dataset using CSV file format. To perform cluster analysis using K-Means in WEKA. To change the cluster size and plot the graph and illustrate the visualization of cluster.

| EmployeeID | Gender | Age | Salary | Credit |
|------------|--------|-----|--------|--------|
| 111 | Male | 28 | 150000 | 39 |
| 222 | Male | 25 | 150000 | 27 |
| 333 | Female | 26 | 160000 | 42 |
| 444 | Female | 25 | 160000 | 40 |
| 555 | Female | 30 | 170000 | 64 |
| 666 | Male | 29 | 200000 | 72 |



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 3

1. Write a R program for the data, Use the two methods below to normalize the following group of data: 200, 300, 400, 600, 1000 (a) min-max normalization by setting min = 0 and max = 1 (b) z-score normalization

2. Write a R program for Two Maths teachers are comparing how their Year 9 classes performed in the end of year exams. Their results are as follows:

Class A: 76, 35, 47, 64, 95, 66, 89, 36, 84, 76, 35, 47, 64, 95, 66, 89, 36, 84

Class B: 51, 56, 84, 60, 59, 70, 63, 66, 50, 51, 56, 84, 60, 59, 70, 63, 66, 50

(i) Find which class had scored higher mean, median and range.

(ii) Plot above in boxplot and give the inferences

3. Implement using WEKA for the given Suppose a database has five transactions. Let min sup= 50%(2) and min con f = 80%.

| Transactions | Items |
|--------------|--------------------|
| T1 | (M, O, N, K, E, Y) |
| T2 | (D, O, N, K, E, Y) |
| T3 | (M, A, K, E) |
| T4 | (M, U, C, K, Y) |
| T5 | (C, O, O, K, I, E) |

Find all frequent item sets using Apriori algorithm

4. Write a R program for the following list of persons with vegetarian or not details given in the table. How will you find out how many of them are vegetarian and how many of them are non-vegetarian? Which type of the person total count is greater value?

| Person | Gopu | Babu | Baby | Gopal | Krishna | Jai | Dev | Malini | Hema | Anu |
|-------------------|------|------|------|-------|---------|-----|-----|--------|------|-----|
| Vegetarian | yes | yes | yes | no | yes | no | no | yes | yes | yes |



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 4

1. Write a R program for the Data: 11,13,13,15,15,16,19,20,20,20,21,21,22,23,24,30,40,45,45,45,71,72,73,75 and find

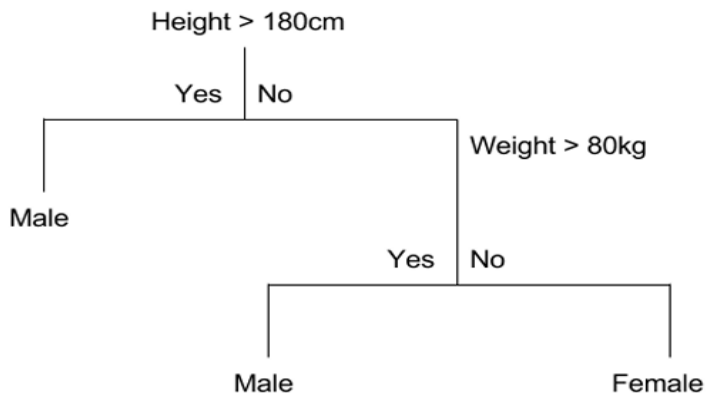
- a) Smoothing by bin mean
- b) Smoothing by bin median
- c) Smoothing by bin boundaries

2. Write a R program to calculate the minimum and maximum values for the feature F are \$50,000 and \$100,000 correspondingly. It needs to range F from 0 to 1. In accordance with min-max normalization, $v =$ \$80,

3. Prediction of Categorical Data using Rule base classification and decision tree classification through WEKA using any datasets. Compare the accuracy using two algorithm and plot the graph

4. Consider this Decision tree :

- a) create the data set for the below tree using ARFF format and calculate accuracy and decision for the same
- b) Using this decision tree generate the rules based on rule based induction.
- c) Compare both the algorithms and plot the confusion matrix.





SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 5

1. Write a R program for the hospital data tested the age and body fat data for 18 randomly selected adults with the following results:

| | | | | | | | | | |
|------|------|------|------|------|------|------|------|------|------|
| age | 23 | 23 | 27 | 27 | 39 | 41 | 47 | 49 | 50 |
| %fat | 9.5 | 26.5 | 7.8 | 17.8 | 31.4 | 25.9 | 27.4 | 27.2 | 31.2 |
| age | 52 | 54 | 54 | 56 | 57 | 58 | 58 | 60 | 61 |
| %fat | 34.6 | 42.5 | 28.8 | 33.4 | 30.2 | 34.1 | 32.9 | 41.2 | 35.7 |

- (a) Calculate the mean, median, and standard deviation of age and %fat.
- (b) Draw the boxplots for age and %fat.
- (c) Draw a scatter plot and a q-q plot based on these two variables.

2. Write a R program Use the below methods to normalize the following group of data: 200, 300, 400, 600, 1000

- (a) min-max normalization by setting min = 0 and max = 1
- (b) z-score normalization

3. Create the dataset using ARFF file format:

| Transaction ID | Items |
|----------------|-------------------------|
| T1 | Hot Dogs, Buns, Ketchup |
| T2 | Hot Dogs, Buns |
| T3 | Hot Dogs, Coke, Chips |
| T4 | Chips, Coke |
| T5 | Chips, Ketchup |
| T6 | Hot Dogs, Coke, Chips |

- a. Find the **frequent itemsets** and generate **association rules** on this. Assume that minimum support threshold ($s = 33.33\%$) and minimum confident threshold ($c = 60\%$).
- b. List the various rule generated by apriori and FP tree algorithm, mention whether accepted or rejected.

4. Implement of the R script using marks scored by a student in his model exam has been sorted as follows: 55, 60, 71, 63, 55, 65, 50, 55, 58, 59, 61, 63, 65, 67, 71, 72, 75. Partition them into three bins by each of the following methods. Plot the data points using histogram.

- (a) equal-frequency (equi-depth) partitioning
- (b) equal-width partitioning
- (c) clustering



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 6

1. Suppose that a hospital tested the age and body fat data for 18 randomly selected adults with the following results:

| | | | | | | | | | |
|------|------|------|------|------|------|------|------|------|------|
| age | 23 | 23 | 27 | 27 | 39 | 41 | 47 | 49 | 50 |
| %fat | 9.5 | 26.5 | 7.8 | 17.8 | 31.4 | 25.9 | 27.4 | 27.2 | 31.2 |
| age | 52 | 54 | 54 | 56 | 57 | 58 | 58 | 60 | 61 |
| %fat | 34.6 | 42.5 | 28.8 | 33.4 | 30.2 | 34.1 | 32.9 | 41.2 | 35.7 |

- Use min-max normalization to transform the value 35 for age onto the range [0.0, 1.0].
- Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.
- Use normalization by decimal scaling to transform the value 35 for age. Perform the above functions using R – tool

2. Perform the below two methods using R – tool to normalize the following group of data: 200, 300, 400, 600, 1000

- min-max normalization by setting min = 0 and max = 1
- z-score normalization

3. Implement using WEKA the Bayes classification and decision tree (using training and test data)

| RID | age | income | student | credit_rating | Class: buys_computer |
|-----|-----------|--------|---------|---------------|----------------------|
| 1 | <=30 | high | no | fair | no |
| 2 | <=30 | high | no | excellent | no |
| 3 | 31 ... 40 | high | no | fair | yes |
| 4 | >40 | medium | no | fair | yes |
| 5 | >40 | low | yes | fair | yes |
| 6 | >40 | low | yes | excellent | no |
| 7 | 31 ... 40 | low | yes | excellent | yes |
| 8 | <=30 | medium | no | fair | no |
| 9 | <=30 | low | yes | fair | yes |
| 10 | >40 | medium | yes | fair | yes |
| 11 | <=30 | medium | yes | excellent | yes |
| 12 | 31 ... 40 | medium | no | excellent | yes |
| 13 | 31 ... 40 | high | yes | fair | yes |
| 14 | >40 | medium | no | excellent | no |

4. The given are the strike-rates scored by a batsman in season 1 in different tournaments. 100, 70, 60, 90, 90 . Perform the above functions using R – tool

- min-max normalization by setting min = 0 and max = 1
- z-score normalization
- z-score normalization using the mean absolute deviation instead of standard deviation

SET 7



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

1. Write a R program for the following values are the number of pencils available in the different boxes. Create a vector and find out the mean, median and mode values of set of pencils in the given data.

Box1 Box2 Box3 Box4 Box5 Box6 Box7 Box8 Box9 Box 10

9 25 23 12 11 6 7 8 9 10

2. Write a R program to the following table would be plotted as (x,y) points, with the first column being the x values as number of mobile phones sold and the second column being the y values as money. To use the scatter plot for how many mobile phones sold.

x :4 1 5 7 10 2 50 25 90 36

y :12 5 13 19 31 7 153 72 275 110

3. Implement using WEKA for the given Suppose a database has five transactions. Let min sup= 50%(2) and min con f = 80%.

| Transactions | Items |
|--------------|--------------------|
| T1 | (M, O, N, K, E, Y) |
| T2 | (D, O, N, K, E, Y) |
| T3 | (M, A, K, E) |
| T4 | (M, U, C, K, Y) |
| T5 | (C,O, O, K, I ,E) |

Find all frequent item sets using Apriori algorithm and FP-Growth Tree

4. Implement using WEKA and generate rules using FP growth algorithm using the given dataset which has the following transactions with items purchased: Consider the values as support=50% and confidence=75%.

| Transaction ID | Items Purchased |
|----------------|---------------------------|
| 1 | Bread, Cheese, Egg, Juice |
| 2 | Bread, Cheese, Juice |
| 3 | Bread, Milk, Yogurt |
| 4 | Bread, Juice, Milk |
| 5 | Cheese, Juice, Milk |



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

1. Implement of the R script using marks scored by a student in his model exam has been sorted as follows: 55, 60, 71, 63, 55, 65, 50, 55, 58, 59, 61, 63, 65, 67, 71, 72, 75. Partition them into three bins by each of the following methods. Plot the data points using histogram.

(a) equal-frequency (equi-depth) partitioning (b) equal-width partitioning

2. Assume the Tennis coach wants to determine if any of his team players are scoring outliers. To visualize the distribution of points scored by his players, then how can he decide to develop the box plot? Give suitable example using Boxplot visualization technique.

3. Implement using WEKA for the data set and perform the Apriori Algorithm and FP algorithm support:3 and confidence=50%

| Customer ID | Transaction ID | Items Bought |
|-------------|----------------|--------------|
| 1 | 0001 | {a, d, e} |
| 1 | 0024 | {a, b, c, e} |
| 2 | 0012 | {a, b, d, e} |
| 2 | 0031 | {a, c, d, e} |
| 3 | 0015 | {b, c, e} |
| 3 | 0022 | {b, d, e} |
| 4 | 0029 | {c, d} |
| 4 | 0040 | {a, b, c} |
| 5 | 0033 | {a, d, e} |
| 5 | 0038 | {a, b, e} |

4. Prediction of Diabetes Data using Decision tree classifier in WEKA. Compare it with Support Vector Machine classifier. Show the result accuracy and F1 measure calculation. Plot the graph and explain the summary of results



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



Course Code: CSA1614

Course Title: Data Warehousing and Data Mining for Healthcare

SET 9

1. Create an ARFF file for the table below and implement for the Apriori Algorithm and FP growth algorithm and compare the rules generated by both the algorithms. Identify the unique rules generated by the above algorithms.

NOTE: Assume Min_sup=2 and confidence= 50%

| T.ID | ITEMS |
|------|----------------------|
| T1 | SONY, BPL, LG |
| T2 | BPL, SAMSUNG |
| T3 | BPL, ONIDA |
| T4 | SONY, BPL, SAMSUNG |
| T5 | SONY, ONIDA |
| T6 | BPL, ONIDA |
| T7 | SONY, ONIDA |
| T8 | SONY, BPL, ONIDA, LG |
| T9 | SONY, BPL, ONIDA |

2. Implement of the R script using marks scored by a student in his model exam has been sorted as follows: 55, 60, 71, 63, 55, 65, 50, 55, 58, 59, 61, 63, 65, 67, 71, 72, 75. Partition them into three bins by each of the following methods. Plot the data points using histogram.

(a) equal-frequency (equi-depth) partitioning

(b) equal-width partitioning

(c) clustering

3. Generate rules using FP growth algorithm using the given dataset which has the following transactions with items purchased: Consider the values as support=50% and confidence=75%.

| Transaction ID | Items Purchased |
|----------------|---------------------------|
| 1 | Bread, Cheese, Egg, Juice |
| 2 | Bread, Cheese, Juice |
| 3 | Bread, Milk, Yogurt |
| 4 | Bread, Juice, Milk |
| 5 | Cheese, Juice, Milk |

4. Suppose some car is tested for the AvgSpeed and TotalTime data for 9 randomly selected car with the following result

| | | | | | | | | | |
|--------------------|----|----|----|----|----|----|----|----|----|
| AvgSpeed(in kph) | 78 | 81 | 82 | 74 | 83 | 82 | 77 | 80 | 70 |
| TotalTime(in mins) | 39 | 37 | 36 | 42 | 35 | 36 | 40 | 38 | 46 |

a) Calculate the standard deviation of AvgSpeed and TotalTime.

b) Calculate the Variance of AvgSpeed and TotalTime for the above dataset.