# DEPARTMENT OF COMPUTER SCIENCE
# RAJAGIRI COLLEGE OF SOCIAL SCIENCES
# (Autonomous)



# M.Sc. COMPUTER SCIENCE
# (Data Analytics)

# DATA MINING LAB
# CSDA 207
# LAB RECORD

**NAME**         : **BALU S UNNY**

**SEMESTER**     : **SECOND**

**REGISTER NO**    : **2217013**

# DEPARTMENT OF COMPUTER SCIENCE
# RAJAGIRI COLLEGE OF SOCIAL SCIENCES
# (Autonomous)

## M.Sc. COMPUTER SCIENCE (Data Analytics)

# CERTIFICATE

**NAME** : **BALU S UNNY**

**SEMESTER** : **SECOND**

**REGISTER NO** : **2217013**

*Certified that this is a bonafide record of work done by **BALU S UNNY** MSCCS2211 in the Software Laboratory of Rajagiri Department of Computer Science, Kalamassery.*

Sunu Mary Abraham                              Dr. Bindiya M Varghese
Faculty in Charge                              Dean, Computer Science

Internal Examiner                              External Examiner

Place  : Kalamassery
Date   :

# Table of Contents

| | | |
|---|---|---|
| 1 | Give a visualization of the distribution of Iris dataset w.r.t all the features | 78 |
| 2 | Display the plot matrix for the Iris data set | 79 |
| **Section 3 A** | **Association Rule Mining** | |
| 1 | Generate association rules using Market Basket Data set in Orange Tool. Compare the different measures to assess the quality of rules. | 80 |
| 2 | Generate association rules using Food mart Data set in Orange Tool. Compare the different measures to assess the quality of rules. | 81 |
| 3 | Generate association rules using supermarket Data set in WEKA using Apriori algorithm. | 84 |
| 4 | Generate association rules using supermarket Data set in WEKA using FP – growth Algorithm | 85 |
| **Section 3 B** | **Association Rule Mining** | |
| 1 | Generate association rules using Lenses Data set in Orange Tool | 86 |
| 2 | Generate association rules using Lenses Data set in WEKA using Apriori algorithm. | 88 |
| 3 | Generate association rules using Lenses Data set in WEKA using FP – growth Algorithm | 90 |
| **Section 4** | **Classification** | |
| 1 | Generate a classifier in Orange Tool from Iris dataset using Decision Tree. | 91 |
| 2 | Generate a classifier in Orange Tool from Titanic dataset using Decision Tree | 92 |
| 3 | Generate a classifier in Weka Tool from Pima- Diabetes dataset using Decision Tree. | 93 |
| 4 | Generate a classifier in Weka Tool from contact lenses dataset using Decision Tree. | 94 |
| 5 | Generate a classifier in WEKA using Housing Dataset Decision using Decision Tree and Naïve Bayesian Classifier and compare the results. | 96 |
| 6 | Generate a classifier in Orange Tool from Iris dataset using K-Nearest Neighbour and SVM Classification. Compare the models. | 97 |
| 7 | Generate a classifier in WEKA using Iris dataset with K-Nearest Neighbour Classification and SVM Classification. Compare the models. | 98 |
| 8 | Generate a classifier in Orange Tool /WEKA for diabetes dataset using Linear Regression. | 99 |
| 9 | Generate a classifier in orange Tool/WEKA for heart disease dataset using Logistic Regression. | 100 |