# A REPORT ON

# DRUG CONSUMPTION

## JULY 2018

**Submitted to:**

**Dr. Ankush Mittal**

**Director**

**Raman Classes**

**Submitted By:**

**Rajat Tyagi**

**B.tech(CSE)**

**GEU**

# CANDIDATE'S DECLARATION

We hereby declare that I own the full responsibility for the information, results etc. provided in the projects in this report

submitted to **RWX Technology** for the award of certificate of **Certificate of Summer Internship** done during June-July 2018. We have taken care in all respect to honor the intellectual property right and have acknowledged the contribution of others for using them in academic purpose and further declare that in case of any violation of intellectual property right or copyright we shall be fully responsible for the same. Our supervisor should not be held responsible for full or partial violation of copyright or intellectual property right.

**NAME :** Rajat Tyagi

**DATE :** 30-7-18

**PLACE :** Roorkee

# Abstract

Machine learning (ML) is a category of algorithm that allows software applications to become  more accurate in predicting outcomes without being explicitly programmed.

The Weka machine learning workbench is a modern platform for applied machine learning. Weka is an acronym which stands for Waikato Environment for Knowledge Analysis. It is also the name of a New Zealand bird the Weka.The basic premise of machine learning is to build algorithms that can receive input data and use statistical analysis to predict an output while updating outputs as the new data becomes available. In this mini project, various machine learning techniques have been applied on the data set taken from [1]. Among the operations performed upon the data set, Classification and selection of attributes or more precisely, feature selection are the major ones. The Weka 3.8 software is used for the analysis of the data set and all the operations mentioned above. Various graphs have been plotted for different attributes and parameters. The data set taken is  drug consumption datase quantified data set donated in year 2016.

# TABLE OF CONTENT

**CLASSIFICATION OF DRUG CONSUMPTION DATASE:**

# INTRODUCTION

**MACHINE LEARNING**
Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. Machine learning focuses on the development of computer programs that can access data and use it learn for themselves.

**SOME MACHINE LEARNING METHODS**
Machine learning algorithms are often categorized as supervised or unsupervised.

- **Supervised machine learning algorithms** can apply what has been learned in the past to new data using labeled examples to predict future events. Starting from the analysis of a known training dataset, the learning algorithm produces an inferred function to make predictions about the output values. The system is able to provide targets for any new input after sufficient training. The learning algorithm can also compare its output with the correct, intended output and find errors in order to modify the model accordingly.

- In contrast, **unsupervised machine learning algorithms** are used when the information used to train is neither classified nor labeled. Unsupervised learning studies how systems can infer a function to describe a hidden structure from unlabeled data. The system doesn't figure out the right output, but it explores the data and can draw inferences from datasets to describe hidden structures from unlabeled data.

- **Semi-supervised machine learning algorithms** fall somewhere in between supervised and unsupervised learning, since they use both labeled and unlabeled data for training – typically a small amount of labeled data and a large amount of unlabeled data. The systems that use this method are able to considerably improve learning accuracy. Usually, semi-supervised learning is chosen when the acquired labeled data requires skilled and relevant resources in order to train it / learn from it. Otherwise, acquiringunlabeled data generally doesn't require additional resources.

- **Reinforcement machine learning algorithms** is a learning method that interacts with its environment by producing actions and discovers errors or rewards. Trial and error search and delayed reward are the most relevant characteristics of reinforcement learning. This method allows machines and software agents to automatically determine the ideal behavior within a specific context in order to maximize its performance. Simple reward feedback is required for the agent to learn which action is best; this is known as the reinforcement signal.

**WEKA (3.8)**

**What is Weka?** Weka is a collection of machine learning algorithms for data mining tasks. The algorithms can either be applied directly to a dataset or called from your own Java code. Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization. It is also well-suited for developing new machine learning schemes.
**Datasets in Weka** Each entry in a dataset is an instance of the java class: − weka.core.Instance Each instance consists of a number of attributes.

# DATA SET DESCRIPTION

Database contains records for 1885 respondents. For each respondent 12 attributes are known: Personality measurements which include NEO-FFI-R (neuroticism, extraversion, openness to experience, agreeableness, and conscientiousness), BIS-11 (impulsivity), and ImpSS (sensation seeking), level of education, age, gender, country of residence and ethnicity. All input attributes are originally categorical and are quantified. After quantification values of all input features can be considered as real-valued. In addition, participants were questioned concerning their use of 18 legal and illegal drugs (alcohol, amphetamines, amyl nitrite, benzodiazepine, cannabis, chocolate, cocaine, caffeine, crack, ecstasy, heroin, ketamine, legal highs, LSD, methadone, mushrooms, nicotine and volatile substance abuse and one fictitious drug (Semeron) which was introduced to identify over-claimers. For each drug they have to select one of the answers: never used the drug, used it over a decade ago, or in the last decade, year, month, week, or day. Database contains 18 classification problems. Each of independent label variables contains seven classes: "Never Used", "Used over a Decade Ago", "Used in Last Decade", "Used in Last Year", "Used in Last Month", "Used in Last Week", and "Used in Last Day".

Problem which can be solved:
* Seven class classifications for each drug separately.
* Problem can be transformed to binary classification by union of part of classes into one new class. For example, "Never Used", "Used over a Decade Ago" form class "Non-user" and all other classes form class "User".
* The best binarization of classes for each attribute.
* Evaluation of risk to be drug consumer for each drug.

**ATTRIBUTES**
1. ID is number of record in original database. Cannot be related to participant. It can be used for reference only.
2. Age (Real) is age of participant
3. Gender (Real) is gender of participant:
4. Education (Real) is level of education of participant
5. Country (Real) is country of current residence of participant
6. Ethnicity (Real) is ethnicity of participant
7. Nscore (Real) is NEO-FFI-R Neuroticism.
8. Escore (Real) is NEO-FFI-R Extraversion.
9. Oscore (Real) is NEO-FFI-R Openness to experience
10. Ascore (Real) is NEO-FFI-R Agreeableness.
11. Cscore (Real) is NEO-FFI-R Conscientiousness
12. Impulsive (Real) is impulsiveness measured by BIS-11.
13. SS (Real) is sensation seeing measured by ImpSS
14. Alcohol is class of alcohol consumption.
15. Amphet is class of amphetamines consumption.

16. Amyl is class of amyl nitrite consumption
17. Benzos is class of benzodiazepine consumption.
18. Caff is class of caffeine consumption
19. Cannabis is class of cannabis consumption.
20. Choc is class of chocolate consumption.
21. Coke is class of cocaine consumption.
22. Crack is class of crack consumption.
23. Ecstasy is class of ecstasy consumption.
24. Heroin is class of heroin consumption.
25. Ketamine is class of ketamine consumption.
26. Legalh is class of legal highs consumption.
27. LSD is class of alcohol consumption
28. Meth is class of methadone consumption.
29. Mushrooms is class of magic mushrooms consumption.
30. Nicotine is class of nicotine consumption.
31. Semer is class of fictitious drug Semeron consumption.
32. VSA is class of volatile substance abuse consumption.

# TECHNIQUES USED

**J48 DECISION TREE**

Classification is the process of building a model of classes from a set of records that contain class labels. Decision Tree Algorithm is to find out the way the attributes-vector behaves for a number of instances. Also on the bases of the training instances the classes for the newly generated instances are being found . This algorithm generates the rules for the prediction of the target variable. With the help of tree classification algorithm the critical distribution of the data is easily understandable . J48 is an extension of ID3. The additional features of J48 are accounting for missing values, decision trees pruning, continuous attribute value ranges, derivation of rules, etc. In the WEKA data mining tool, J48 is an open source Java implementation of the C4.5 algorithm. The WEKA tool provides a number of options associated with tree pruning. In case of potential over fitting pruning can be used as a tool for précising. In other algorithms the classification is performed recursively till every single leaf is pure, that is the classification of the data should be as perfect as possible. This algorithm it generates the rules from which particular identity of that data is generated. The objective is progressively generalization of a decision tree until it gains equilibrium of flexibility and accuracy.

**Basic Steps in the Algorithm**

(i) In case the instances belong to the same class the tree represents a leaf so the leaf is returned by  labeling with the same class.

(ii) The potential information is calculated for every attribute, given by a test on the attribute. Then the gain in information is calculated that would result from a test on the attribute.

(iii) Then the best attribute is found on the basis of the present selection criterion and that attribute selected for branching.

**The steps performed for data mining in WEKA are:**

• Data pre-processing and visualization
 • Attribute selection
• Classification (Decision trees)
• Prediction (Nearest neighbour)
• Model evaluation
• Clustering (Cobweb, K-means) • Association rules

Test Options included using training set, supplying test set, cross validation and percentage splitting. Accuracy of the machine was tested through the data in many methods: by selecting all the original attributes and then training the machine, then by deselecting some features and then calculating the accuracy. For feature selection the attribute evaluator used is **CfsSubsetEvaluator**and the method used is **BestFirst method.**

# RESULTS AND DISCUSSION

## LOADING DATASET OR PRE –PROCESSING

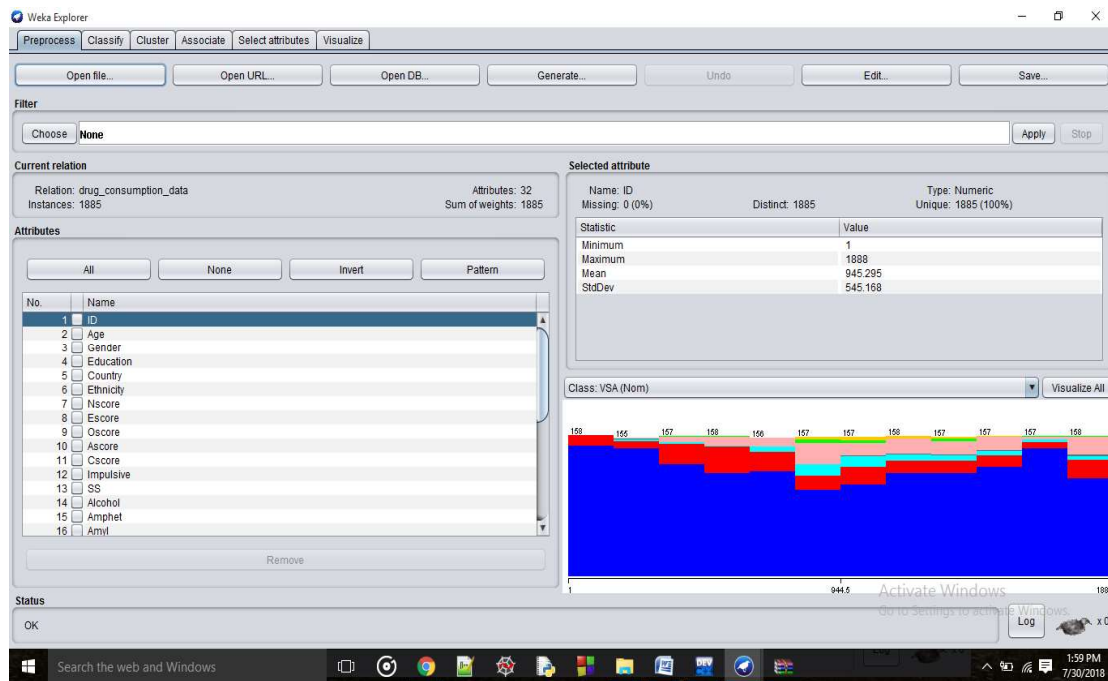This dataset loaded in comma separated value (.csv) format



Figure 1: Pre-processing Window where the data is loaded.

# ATTRIBUTE SELECTION

- **Selected attributes** : 2,5,13,15,16,17,21,22,24,26,27,28,30,31 : 14

   Age
   Country
   SS
   Amphet
   Amyl
   Benzos
   Coke
   Crack
   heroin
   Legalh
   LSD
   Meth
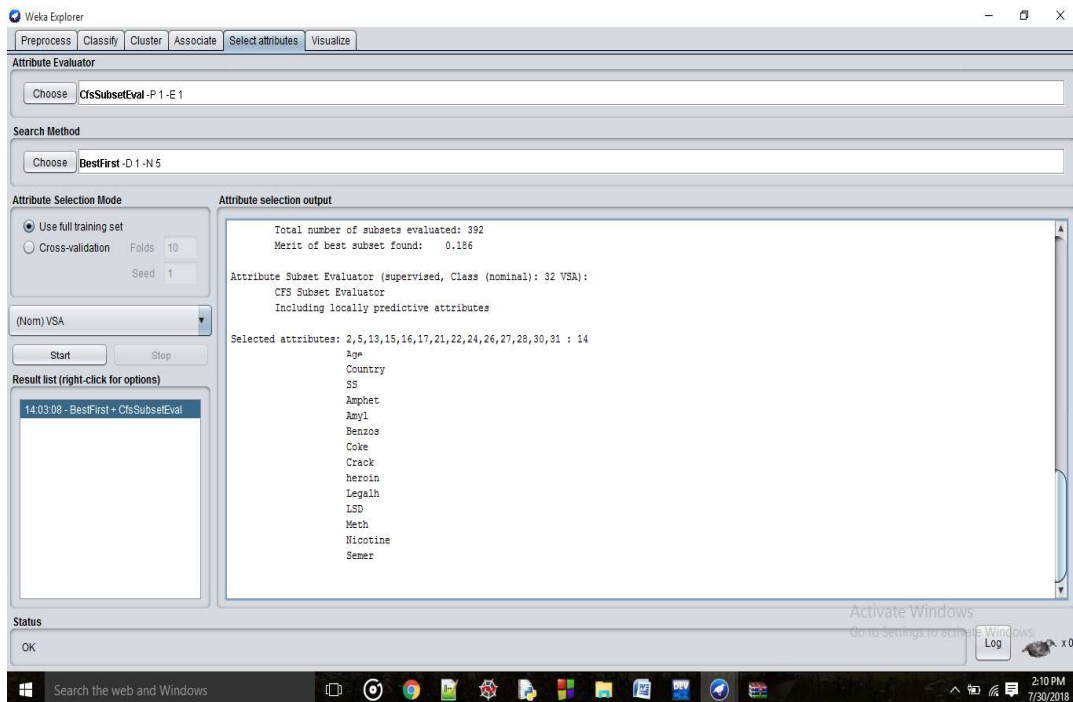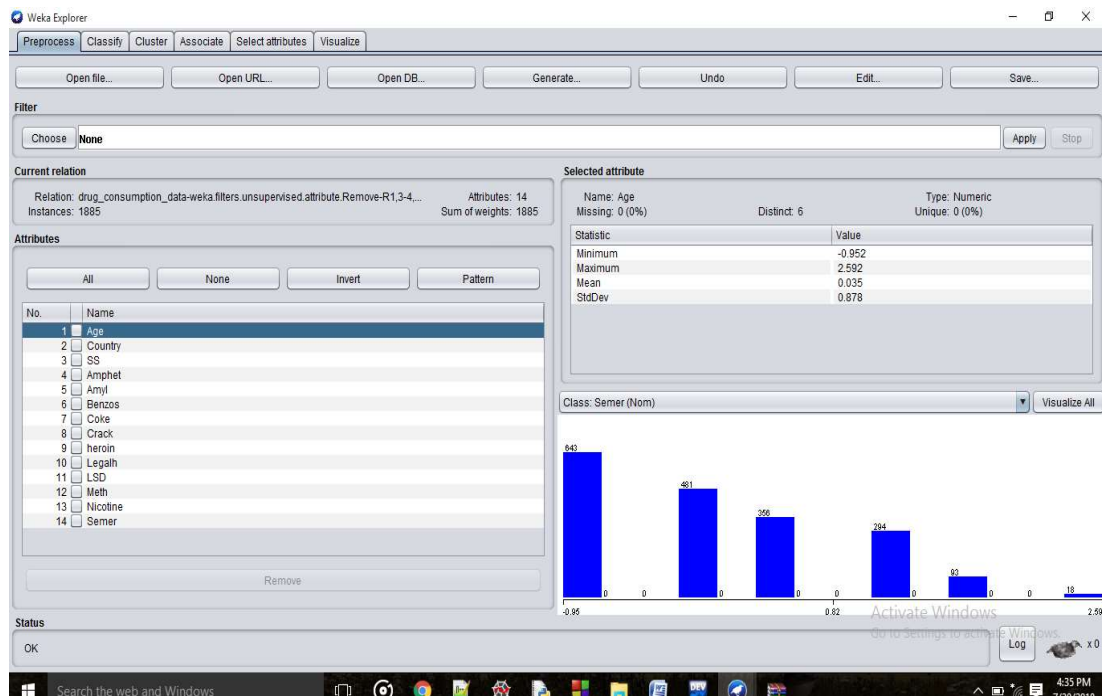   Nicotine
   Semer



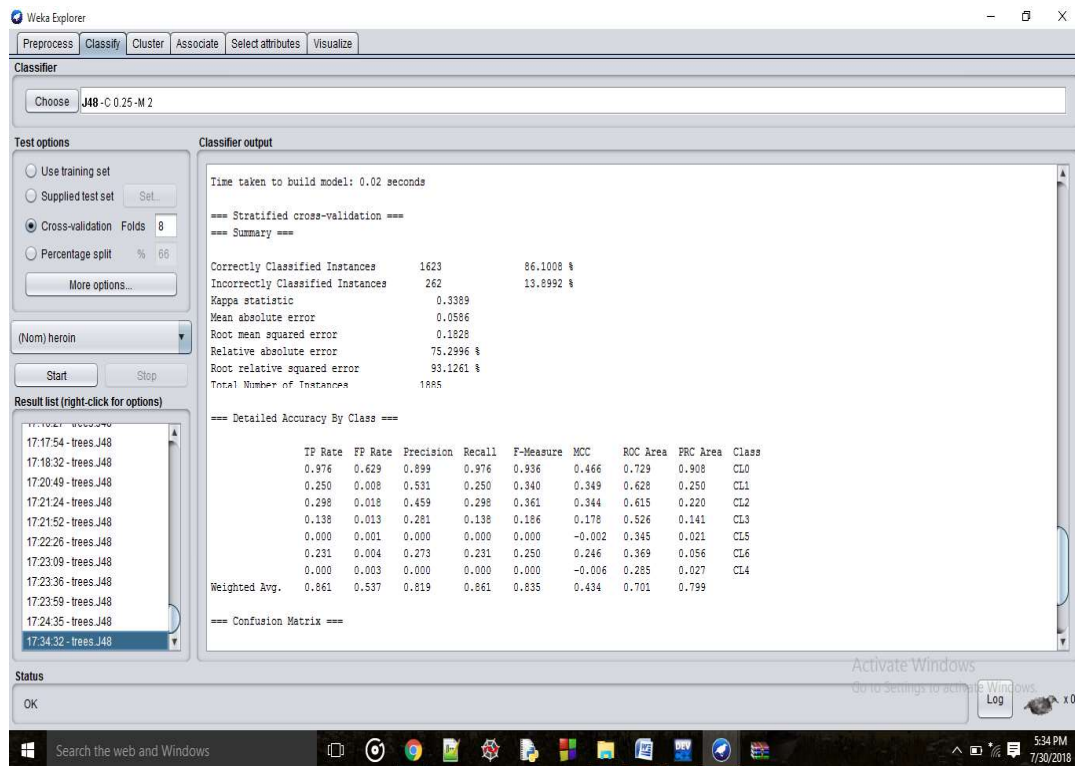Figure 2: Features selection process

•     Figure 3: Remaining Attributes after Feature selection

**SELECTED  ATTRIBUTES TAKEN, CLASSIFICATION APPLIED, CROSS VALIDATION USED, ABOUT HEROIN ATTRIBUTE**

| No. of Cross Validation Folds | Accuracy obtained | Error observed |
|---|---|---|
| 10 | 85.7825 | 14.2175 |
| 9 | 85.4642 | 14.5358 |
| 8 | 86.1008 | 13.8992 |
| 7 | 85.2520 | 14.7480 |
| 6 | 85.7294 | 14.2706 |
| 5 | 85.9947 | 14.0053 |
| 4 | 85.4111 | 14.5889 |
| 3 | 85.8886 | 14.1114 |
| 2 | 85.6764 | 14.3236 |
| 11 | 86.1008 | 13.8992 |
| 12 | 85.6233 | 14.3767 |
| 13 | 85.7294 | 14.2706 |
| 14 | 85.7294 | 14.2706 |
| 15 | 85.5703 | 14.4297 |

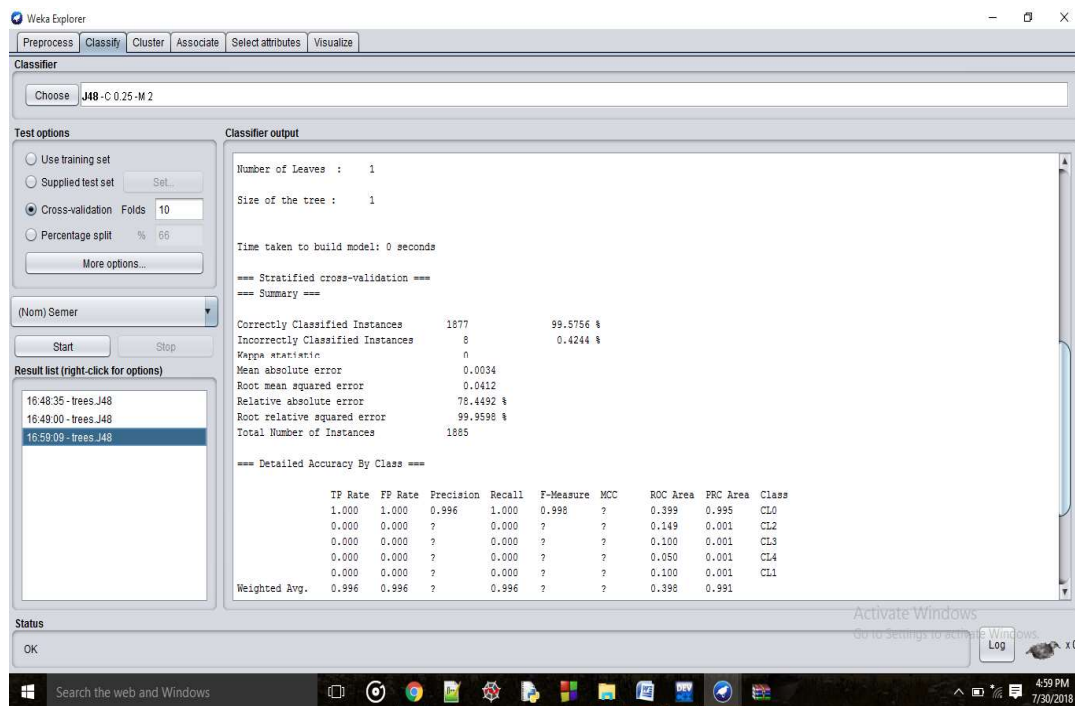| 16 | 85.8355 | 14.1645 |
|---|---|---|
| 17 | 85.7294 | 14.2706 |
| 18 | 85.3581 | 14.6419 |
| 19 | 85.5703 | 14.4297 |
| 20 | 85.6764 | 14.3236 |



- **Discussion for heroin attribute**

As per the above table , it is clear that the maximum accuracy is 86.0001% when we take heroin as root attribute and cross-validation folds are 8 .
This highlighted the importance of Feature Selection in classification. Attribute selection is done in the project using CfsSubsetEval evaluator.

**SELECTED ATTRIBUTES TAKEN, CLASSIFICATION APPLIED, CROSS VALIDATION USED, ABOUT SEMER ATTRIBUTE**

| No. of Cross Validation Folds | Accuracy obtained | Error observed |
|---|---|---|
| 10 | 99.5756 | 0.4244 |
| 9 | 99.5756 | 0.4244 |
| 8 | 99.5756 | 0.4244 |
| 7 | 99.5756 | 0.4244 |
| 6 | 99.5756 | 0.4244 |
| 5 | 99.5756 | 0.4244 |
| 4 | 99.5756 | 0.4244 |
| 3 | 99.5756 | 0.4244 |
| 2 | 99.5756 | 0.4244 |
| 11 | 99.5756 | 0.4244 |
| 12 | 99.5756 | 0.4244 |
| 13 | 99.5756 | 0.4244 |
| 14 | 99.5756 | 0.4244 |
| 15 | 99.5756 | 0.4244 |

| 16 | 99.5756 | 0.4244 |
|---|---|---|
| 17 | 99.5756 | 0.4244 |
| 18 | 99.5756 | 0.4244 |
| 19 | 99.5756 | 0.4244 |
| 20 | 99.5756 | 0.4244 |



- **Discussion for semer attribute**

As per the above table , it is clear that the accuracy is 99.5756% when we take semer as root attribute .
This highlighted the importance of Feature Selection in classification. Attribute selection is done in the project using CfsSubsetEval evaluator.

# CONCLUSION

In this project we learn how to select important features of a dataset and reduce the computational work.
We also learn how attributes play important role for accuracy.
Example :  When we take heroin as root attribute we got maximum 86.0001% accuracy but in case of semer we got 99.5756% accuracy.
Above example shows how attribute selection is important in weka.