# Empirical Software Engineering (SE-404)

# LAB A1-G2

# Laboratory Manual



# Department of Software Engineering

# DELHI TECHNOLOGICAL UNIVERSITY(DTU)

Shahbad Daulatpur, Bawana Road, Delhi-110042

**Submitted to: -**                                          **Submitted by:-**

Mr. Sanjay Patidar                                    Name: ASHISH KUMAR

Roll number: 2K18/SE/041

# INDEX

| S.No. | EXPERIMENT | DATE | REMARKS |
|---|---|---|---|
| **10.** | Perform a comparison of the following data analysis tools. WEKA, KEEL, SPSS, MATLAB, R. | 04-01-2022 | |
| **1.** | Consider any empirical study of your choice (Experiments, Survey Research, Systematic Review, Postmortem analysis and case study). Identify the following components for an empirical study:<br>a. Identify parametric and nonparametric tests<br>b. Identify Independent, dependent and confounding variables<br>c. Is it Within-company and cross-company analysis?<br>d. What type of dataset is used? Proprietary and open-source software | 18-01-2022 | |
| **2.** | Defect detection activities like reviews and testing help in identifying the defects in the artifacts (deliverables). These defects must be classified into various buckets before carrying out the root cause analysis. Following are some of the defect categories: Logical, User interface, Maintainability, and Standards. In the context of the above defect categories, classify the following statements under the defect categories. | 25-01-2022 | |
| **3.** | Consider any prediction model of your choice.<br>a. Analyze the dataset that is given as a input to the prediction model<br>b. Find out the quartiles for the used dataset<br>c. Analyze the performance of a model using various performance metrics. | 25-01-2022 | |
| **8.** | Why is version control important? How many types of version control systems are there? Demonstrate how version control is used in a proper sequence (stepwise). | 01-02-2022 | |
| **9.** | Demonstrate how Git can be used to perform version control? | 01-02-2022 | |
| **11.** | Validate the results obtained in experiment 3 using 10-cross validation, hold out validation or leave one out cross-validation. | 15-02-2022 | |
| **4.** | Consider defect dataset and perform following feature reduction techniques using Weka tool. Validate the dataset using 10-cross validation.<br>a. Correlation based feature evaluation<br>b. Relief Attribute feature evaluation<br>c. Information gain feature evaluation<br>d. Principle Component | 23-02-2022 | |

# Empirical Software Engineering LAB – A1 G2 EXPERIMENT 4

**-** ASHISH KUMAR
- 2K18/SE/041

## Experiment Objective:-

Consider defect dataset and perform following feature reduction techniques using WEKA tool. Validate the dataset using 10-cross validation.

a. Correlation based feature evaluation
b. Relief Attribute feature evaluation
c. Information gain feature evaluation
d. Principle Component

## Introduction:-

**DATASET USED:**
- KC2 Dataset in .arff format.
- Author: Mike Chapman, NASA
- Link to dataset: https://datahub.io/machine-learning/kc2#readme

One of the NASA Metrics Data Program defect data sets. Data from software for science data processing. Data comes from McCabe and Halstead features extractors of source code. These features were defined in the 70s in an attempt to objectively characterize code features that are associated with software quality.

**Attribute Information**

1. loc: numeric  McCabe's line count of code
2. $v(g)$ : numeric  McCabe "cyclomatic complexity"
3. $ev(g)$ : numeric  McCabe "essential complexity"
4. $iv(g)$ : numeric  McCabe "design complexity"
5. n: numeric  Halstead total operators + operands
6. v: numeric  Halstead "volume"
7. l: numeric  Halstead "program length"
8. d: numeric  Halstead "difficulty"
9. i: numeric  Halstead "intelligence"
10. e: numeric  Halstead "effort"
11. b: numeric  Halstead
12. t: numeric  Halstead's time estimator
13. lOCode : numeric % Halstead's line count
14. lOComment : numeric % Halstead's count of lines of comments
15. lOBlank : numeric % Halstead's count of blank lines
16. lOCodeAndComment: numeric
17. uniq_Op : numeric  unique operators

18. uniq_Opnd : numeric % unique operands
19. total_Op : numeric % total operators
20. total_Opnd : numeric % total operands
21. branchCount : numeric % Branch Count of the flow graph
22. problems: {false,true} % module has/has not one or more reported defects

## Procedure:

1. Open the Weka GUI Chooser.
2. Click the "Explorer" button to launch the Explorer.
3. Open the dataset.
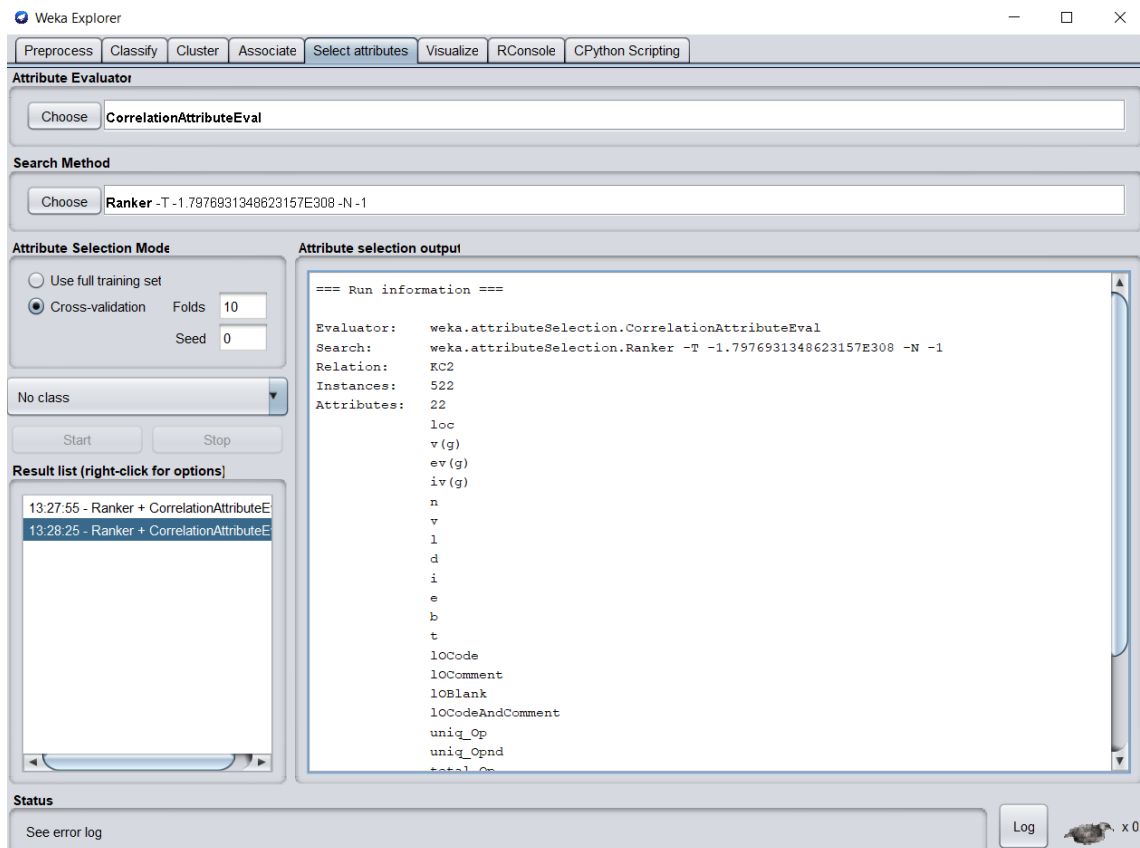4. Click the "Select attributes" tab to access the feature selection methods
Feature selection is divided into two parts:

- Attribute Evaluator
- Search Method.

## Result:-

**a) Correlation Based Feature Selection:**

A popular technique for selecting the most relevant attributes in your dataset is to use correlation. Correlation is more formally referred to as Pearson's correlation coefficient in statistics. Weka supports correlation based feature selection with the CorrelationAttributeEval technique that requires use of a Ranker search method.

## b) Relief Attribute feature evaluation:

Weka supports correlation based feature selection with ReliefFAttributeEval, the technique that requires use of a Ranker search method.

## ReliefFAttributeEval :

Evaluates the worth of an attribute by repeatedly sampling an instance and considering the value of the given attribute for the nearest instance of the same and different class.

## c) Information gain feature evaluation:

Another popular feature selection technique is to calculate the information gain. You can calculate the information gain (also called entropy) for each attribute for the output variable. Entry values vary from 0 (no information) to 1 (maximum information).
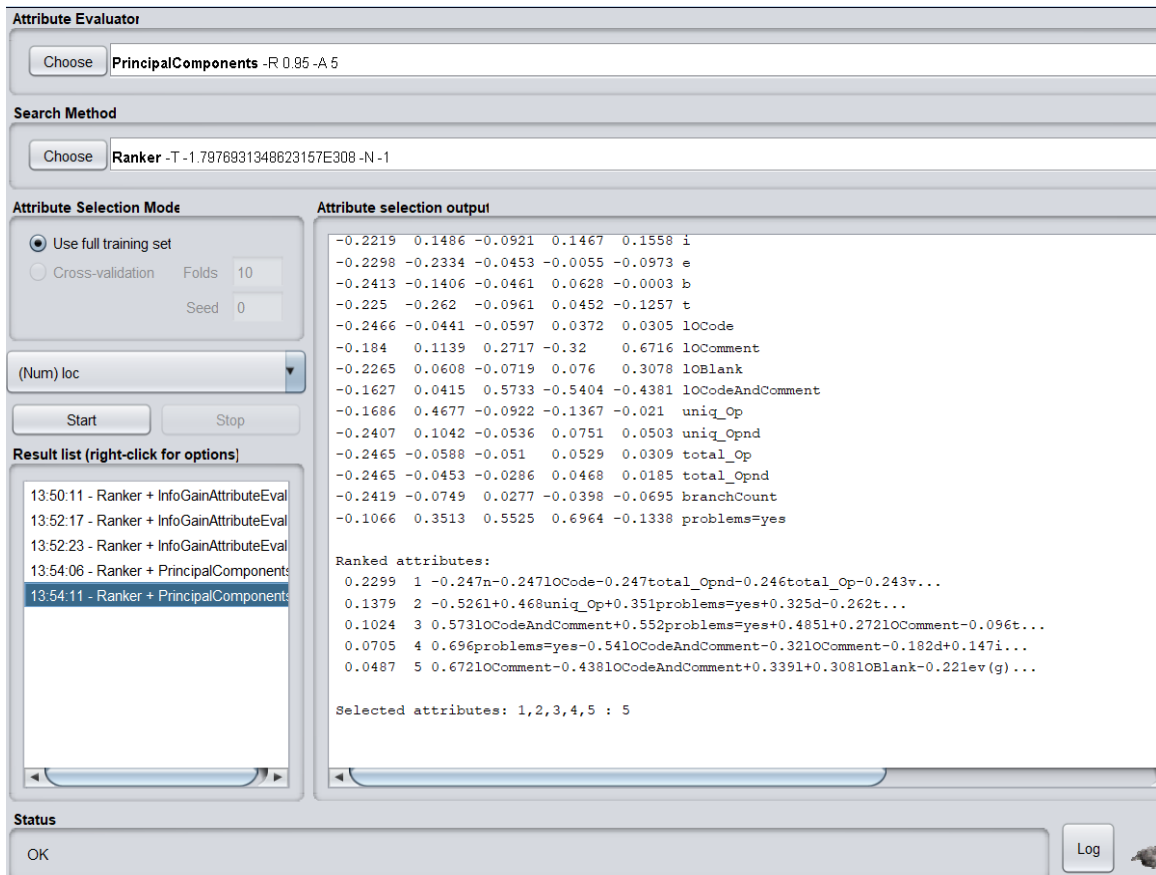
Weka supports feature selection via information gain using the InfoGainAttributeEval Attribute Evaluator. Like the correlation technique above, the Ranker Search Method must be used.

```
=== Attribute selection 10 fold cross-validation (stratified), seed: 0

average merit        average rank  attribute
 0.23   +- 0.016      1.7 +- 0.64    20 total_Opnd
 0.228  +- 0.015      1.8 +- 1.47    18 uniq_Opnd
 0.215  +- 0.012      3.5 +- 0.92     1 loc
 0.21   +- 0.012      5.1 +- 2.47    11 b
 0.212  +- 0.015      5.1 +- 2.17    19 total_Op
 0.209  +- 0.011      5.5 +- 1.2      6 v
 0.207  +- 0.011      6.3 +- 1        5 n
 0.198  +- 0.017      9.6 +- 2.33    12 t
 0.198  +- 0.014      9.9 +- 3.21     9 i
 0.198  +- 0.017     10   +- 2.49    10 e
 0.192  +- 0.018     11   +- 2.24    13 lOCode
 0.187  +- 0.016     12.1 +- 2.3      2 v(g)
 0.185  +- 0.012     13   +- 1.34    15 lOBlank
 0.185  +- 0.013     13   +- 1.61     8 d
 0.175  +- 0.013     15.3 +- 1.35     7 l
 0.175  +- 0.013     15.6 +- 1.62     4 iv(g)
 0.171  +- 0.015     16   +- 2.49    17 uniq_Op
 0.17   +- 0.017     16.5 +- 1.86    21 branchCount
 0.141  +- 0.009     19   +- 0        3 ev(g)
 0.126  +- 0.012     20   +- 0       14 lOComment
 0.051  +- 0.006     21   +- 0       16 lOCodeAndComment
```

**d) Principle Component:**

Weka Explorer can be used to perform principal components analysis and transformation of the data. It is used in conjunction with a Ranker search.



**<u>Learning from experiment:</u>**- We have successfully learned about WEKA and Correlation based feature evaluation, Relief Attribute feature evaluation, Information gain feature evaluation and Principle Component. We successfully have used the select attributes feature of WEKA for validating a dataset using 10-cross validation.