HO CHI MINH UNIVERSITY OF SCIENCE

VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY

# REPORT LAB 01

# DATA PREPROCESSING - DATA EXPLORATION

Course: CSC14004 – DATA MINING AND APPLICATIONS

| | |
|---|---|
| Lecturer: | Prof. Dr. Le Hoai Bac |
| Teaching Assistants: | Ms. Nguyen Thi Thu Hang |
| | Mr. Nguyen Bao Long |
| Group students: | Dang Tien Dat      20127458 |
| | Pham Thi Anh Phat    20127680 |

(2022 - 2023)

# CONTENTS

# I.    GROUP INFORMATION

| ID | Full Name | Contribution rate (%) |
|---|---|---|
| 20127458 | Dang Tien Dat | 50 |
| 20127680 | Pham Thi Anh Phat | 50 |

- We definitely completed all questions and requirements. 👏 👏
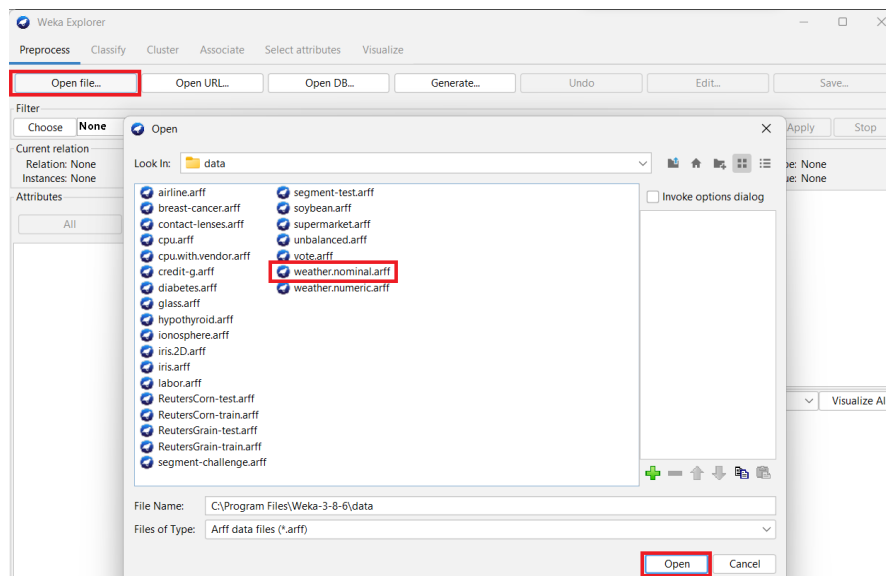
# II.    REQUIREMENTS

### a.  Install WEKA

a. Requirement 1: capture a screen that contains the "**Explorer**" function in your desktop background.



b. Requirement 2: Explain the meaning of **Current Relation**, **Attributes**, and **Selected Attribute** in **Preprocess** tag. Briefly explain the meaning of the other tags in WEKA Explorer.

➕ After load available dataset



➕ **The screen displays:**

1. **PREPROCESS TAB**



- **Current relation:** shows the name of the of the database currently loaded as well as number of instances, number of attributes and sum of weights.

- **Attributes:** displays all fields in the database. The weather database has 5 attributes: outlook, temperature, humidity, windy, play. When clicking one of these attributes, **more details** on the attribute itself are displayed on the **Selected attribute**. Further more, we can remove feature by click on **Remove** button.
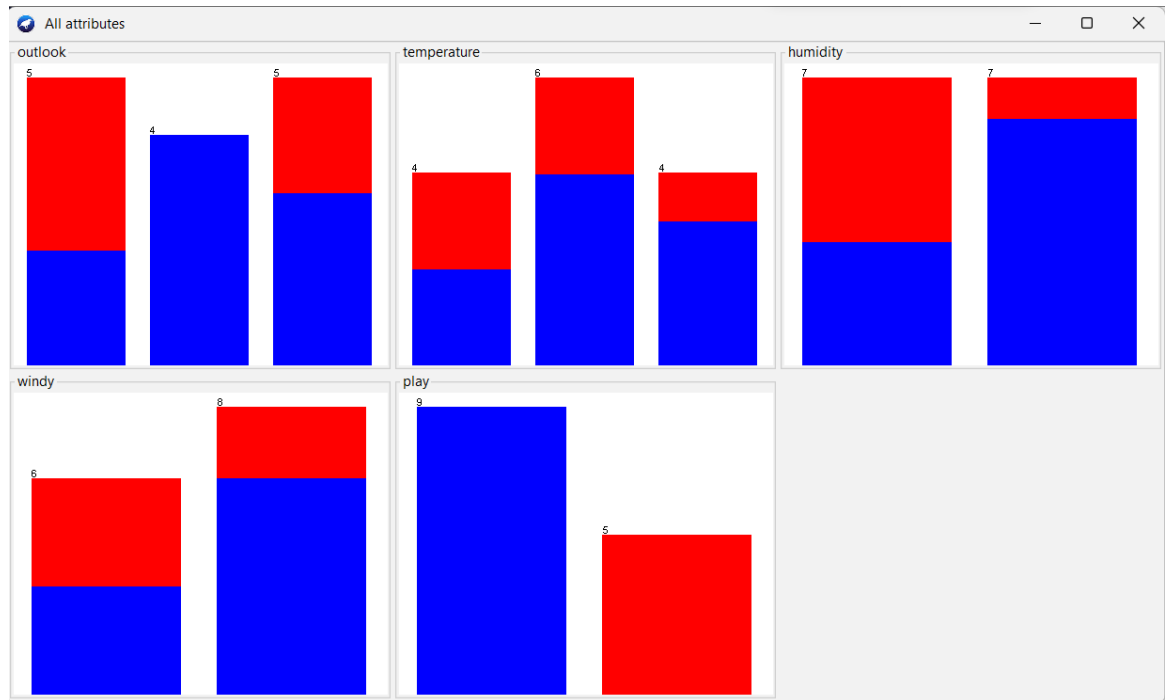


- **Selected attribute** displays:
    - c. Name and type of chosen attribute.
    - d. Number of missing value, distinct value and unique value of this attribute.
    - e. *For example:* type of humidity is Nominal, no missing value, 2 distinct values with no unique value.
    - f. The table below shows the nominal values for this humidity field as high and normal, it also shows count and weight for each value.
    - g. The bottom subwindow visualize representation of the class values.
    - h. If you click on **Visualize All**, all features in one single window as shown here:

5

2. **ANOTHER TABS IN WEKA EXPLORER:**
   a. **CLASSIFY TAB:** provides you several both supervised and unsupervised machine learning <u>algorithms for the classification</u> of your data, such as Linear Regression, Logistic Regression, Support Vector Machines, Decision Trees, RandomTree, RandomForest, NaiveBayes, and so on.
   b. **CLUSTER TAB:** provides several <u>clustering algorithms</u> such as SimpleKMeans, FilteredClusterer, HierarchicalClusterer, and so on.
   c. **ASSOCIATE TAB:** several <u>associate algorithms</u>: Apriori, FilteredAssociator and FPGrowth.
   d. **SELECT ATTRIBUTES:** allows you <u>feature selections</u> based on several algorithms such as ClassifierSubsetEval, PrinicipalComponents, etc.
   e. **VISUALIZE:** visualize your processed data for analysis**.**


b. **Getting Acquainted With WEKA**
   ❖ Exploring Breast Cancer data set: **breast cancer.arff**
     a. <u>How many instances does this data set have? How many attributes does this data set have?</u>

⇨ Data set has **286 instances** with **10 attributes**. The information is in the red box.

b.  <u>Which attribute is used for the label? Can it be changed? How?</u>

⇨ The **label** is **Class** attribute.



We see that the **highlighted Class attribute** tells us that it is used as the label.

Absolutely yes, we **can change label** by following these steps:

- **Step 1:** Click on the "Edit" button in the top-left corner of the screen.



- **Step 2:** Select the attribute that want to use as the label by clicking on its name in the **Attributes** section. And right-click to attribute that we want.
- **Step 3:** Click on the **"Nominal"** button to change the attribute type to nominal if it is not already a nominal attribute. Click on the "**Set as class**" button to set the selected attribute as the label attribute.
- **Step 4:** Save your modified dataset using the "**OK**" button.



c. What is the meaning of each attribute?

- **Class** - the target variable or label that indicates whether the tumor is benign or malignant.

- **Age** - the age of the patient at the time of diagnosis.
- **Menopause** - whether the patient has gone through menopause or not.
- **Tumor size** - the size of the tumor in millimeters.
- **Inv-nodes** - the number of axillary lymph nodes that contain metastatic cancer cells.
- **Node-caps** - whether the cancer cells have spread to the lymph node capsule or not.
- **Deg-malig** - the degree of malignancy of the tumor on a scale of 1 to 3.
- **Breast** - which breast the tumor was found in (left or right).
- **Breast-quad** - the quadrant of the breast in which the tumor was found.
- **Irradiat** - whether the patient received radiation therapy after surgery or not.
    d. <u>Let's investigate the missing value status in each attribute and describe in general ways to solve the problem of missing values</u>
- There are 2 attributes have missing values: node-caps attribute(3%), breast-quad (1 instance)

| Selected attribute | | |
|---|---|---|
| Name: node-caps | | Type: Nominal |
| Missing: 8 (3%) | Distinct: 2 | Unique: 0 (0%) |

| Selected attribute | | |
|---|---|---|
| Name: breast-quad | | Type: Nominal |
| Missing: 1 (0%) | Distinct: 5 | Unique: 0 (0%) |

- To solve the problem of missing values in a dataset, there are several general approaches that you can take using Weka:
1. **Delete the instances or attributes that contain missing values**. This approach is simple but can result in a loss of information if there are many missing values in the dataset.
2. **Replace the missing values with a fixed value**, such as the mean, median, or mode of the attribute values. This approach can work well if the missing values are random and not correlated with the target variable, but can introduce bias if the missing values are related to other variables in the dataset.
3. **Use an imputation algorithm to estimate the missing values** based on the values of other attributes in the dataset. Weka provides several built-in imputation algorithms, such as k-nearest neighbors imputation and expectation-maximization imputation. These algorithms can be useful for

handling missing values when there is a complex relationship between the missing values and other variables in the dataset.

    e. <u>Let's propose solutions to the problem of missing values in the specific attribute.</u>

- There is a solution for **missing values in node-caps** attribute.

**Solution: Filter → unsupervised → instance → RemoveWithValues**



Click **filter button**, the screen shows:



Click **matchMissingValues → True → OK → Apply**

➢ Result:

Selected attribute
　Name: node-caps　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　Type: Nominal
　Missing: 0 (0%)　　　　　　　　　　　　　　　Distinct: 0　　　　　　　　　　　　Unique: 0 (0%)

| No. | Label | Count | Weight |
|-----|-------|-------|--------|
| 1 | yes | 0 | 0 |
| 2 | no | 0 | 0 |

    f. <u>Let's explain the meaning of the chart in the WEKA Explorer. Setting the title for it and describing its legend</u>.



- We can quickly get an overview of the distribution of all attributes in the dataset and the breakdown of distributions by class, red color is **recurrence-events** label, blue is **no-recurrence-events** label.

- *For example*, in **breast** chart. There are 152 instances having left breast and 134 instances having right breast.



11

❖ Exploring Weather data set: **weather.numeric.arff**
  a. <u>How many attributes does this data set have? How many samples?</u>
     <u>Which attributes have data type categorical? Which attributes have</u>
     <u>a data type that is numerical? Which attribute is used for the label?</u>



- Data set has **14 samples** with **5 attributes, 3 categorical attributes**
  (outlook, windy, play) and **2 numeric attributes** (temperature, humidity).
- **Play** attribute is the **label** with 2 values: yes, no.
  b. <u>Let's list **five-number summary** of two attributes temperature and</u>
     <u>humidity. Does WEKA provide these values?</u>
- Five-number summary of numeric attribute are: min, max, mean, StdDev
  and median.
- WEKA just provide four values from **Selected attribute** in **Preprocess** tab.

Selected attribute
Name: temperature        Type: Numeric
Missing: 0 (0%)    Distinct: 12    Unique: 10 (71%)

| Statistic | Value |
|-----------|-------|
| Minimum | 64 |
| Maximum | 85 |
| Mean | 73.571 |
| StdDev | 6.572 |

c. <u>Let's explain the meaning of all charts in the WEKA Explorer. Setting the title for it and describing its legend.</u>
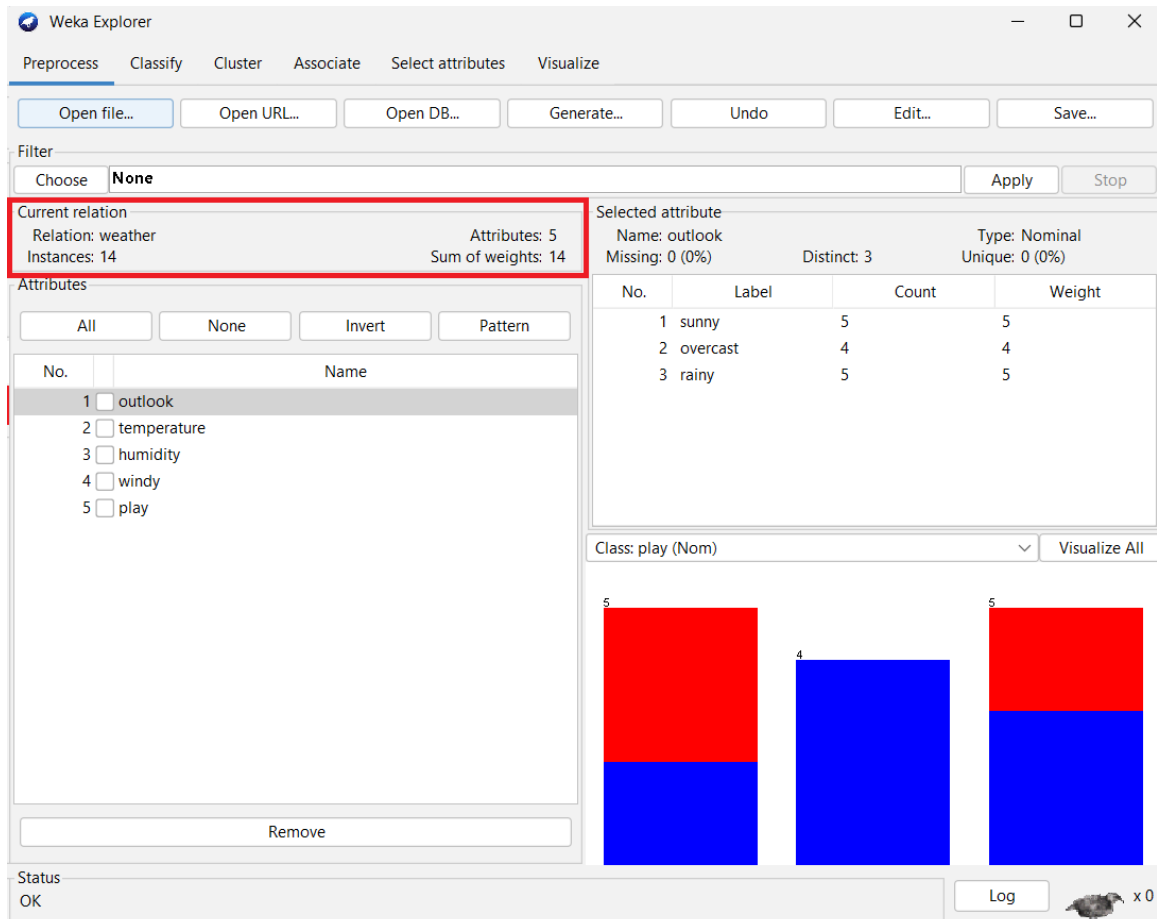


- We can quickly get an overview of the distribution of all attributes in the dataset and the breakdown of distributions by class, red color is **no** label, blue is **yes** label.
- *For example*, in **temperature** chart. There are 8 instances having temperatures in range[64, 74.5] and 6 instances having temperature in range[74.5, 85].



13

d. In **Visualize** tag. What's the name of this chart? Do you think there are any pairs of different attributes that have correlated?



- This is plot matrix.
- Of course, there are pairs of different attributes that have correlated. And to know more detail how correlative between them, you can click here.

❖ Exploring Credit in Germany data set: **credit-g.arff**
  a. What is the **content of the comments section** in credit-g.arff (when opened with any text editor) about? How many samples does the data set have? How many attributes? Describe any five attributes (must have both discrete and continuous attributes).

- Content of comments have many sections:
  1. Title
  2. Source Information
  3. Number of instances
  4. Number of attributes german
  5. Number of attributes german.numer
  6. Attribute description for german
  7. Cost matrix

```
% Description of the German credit dataset.
%
% 1. Title: German Credit data
%
% 2. Source Information
%
% Professor Dr. Hans Hofmann
% Institut f"ur Statistik und "Okonometrie
% Universit"at Hamburg
% FB Wirtschaftswissenschaften
% Von-Melle-Park 5
% 2000 Hamburg 13
%
% 3. Number of Instances:  1000
%
% Two datasets are provided.  the original dataset, in the form
provided
% by Prof. Hofmann, contains categorical/symbolic attributes and
% is in the file "german.data".
%
% For algorithms that need numerical attributes, Strathclyde
University
% produced the file "german.data-numeric".  This file has been
edited
% and several indicator variables added to make it suitable for
% algorithms which cannot cope with categorical variables.
Several
% attributes that are ordered categorical (such as attribute 17)
have
% been coded as integer.    This was the form used by StatLog.
%
%
% 6. Number of Attributes german: 20 (7 numerical, 13
categorical)
%    Number of Attributes german.numer: 24 (24 numerical)
%
%
% 7.  Attribute description for german
```

- Data set has **1000 samples** with **21 attributes.**
- Describe 5 attributes:
  - **Foreign_worker:** Binary attribute with 2 values yes or no. Means whether worker is foreign or not. In the dataset, it has 963 yes and 37 no.

Selected attribute

Name: foreign_worker  Type: Nominal
Missing: 0 (0%)    Distinct: 2    Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | yes | 963 | 963 |
| 2 | no | 37 | 37 |

o **Age**: Numeric attribute and has 53 distinct values, 1 unique value

Selected attribute

Name: age        Type: Numeric
Missing: 0 (0%)    Distinct: 53      Unique: 1 (0%)

| Statistic | Value |
|---|---|
| Minimum | 19 |
| Maximum | 75 |
| Mean | 35.546 |
| StdDev | 11.375 |

and no missing values. Min: 19, Max: 75, Mean: 35.546, StdDev: 11.375.

o **Housing:** Type of accommodation. This is Nominal attribute with 3 distinct values: rent, own or for free. No missing value. Own housing is majority.

Selected attribute

Name: housing        Type: Nominal
Missing: 0 (0%)    Distinct: 3      Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | rent | 179 | 179 |
| 2 | own | 713 | 713 |
| 3 | for free | 108 | 108 |

o **Savings_status:** status of savings with 5 distinct values, no missing or unique value.

  ▪ Saving < 100: 603 instances.
  ▪ 100 <= saving < 500: 103 instances.

Selected attribute

Name: savings_status        Type: Nominal
Missing: 0 (0%)    Distinct: 5      Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | <100 | 603 | 603 |
| 2 | 100<=X<500 | 103 | 103 |
| 3 | 500<=X<1000 | 63 | 63 |
| 4 | >=1000 | 48 | 48 |
| 5 | no known savings | 183 | 183 |

o **Personal_status:** male or female and div/sep or single. It has 4 distinct values and male single is majority.

b. Which attribute is used for the label?

- **Class is the label** with 2 distinct values: good (700 instances) and bad (300 instances).

c. Distribution of continuous attributes?(Left skewed or right skewed?)

- **age**: This attribute has a **right-skewed distribution**, with a more instances towards lower age (around 18-25) and a longer tail towards the higher age (around 40 - 70).

- **duration:** This attribute is slightly right-skewed, with most instances having shorter loan durations and a few instances having longer loan durations.



- **credit amount:** This attribute is highly right-skewed, with most instances having relatively low credit amounts and a few instances having very high credit amounts.

- **installment_commitment**: This attribute is left-skewed, with most instances having moderate to high installment commitments.



- **resident_since: sightly left-skewed**

d. Explain the meaning of all charts in the WEKA Explorer. Setting the title for it and describing its legend.



- We can quickly get an overview of the distribution of all attributes in the dataset and the breakdown of distributions by class, red color is **bad** label, blue is **good** label.
- For example, bar chart shows 596 instances have none own_telephone, 404 instances have own_telephone.



e. **Select attributes** tag. Describe all of the options for attribute selection.

- When click on Select Tab, you will see the next screen. It has **Attribute Evaluator, Search Method, Attribute Selection Mode,** you will find several options by click *Choose* button.

- **Attribute Evaluator.**

When click **Choose**, it shows all options for attribute selection.

If *left-click in CfsSubsetEval*, you can modify some parameters.

- **Search Method**

When click on **Choose**, it shows all search method.

*left-click on BestFirst* to modify parameters.





- **Attribute Selection Mode** has 2 options: use full training set or cross-validation.



When click **Start** button to process the dataset. You will see the following output.

- In the **Attribute selection output** subwindow, you will get result is the list of **Selected** attributes.

f. Which **options** should be used to **select the 5 attributes with the highest correlation**? (Step-by-step description, with step-by-step photos and final results).

- **Step 1:** To calculate correlation, we should choose **CorrelationAttributeEval** in attribute evaluator, the Alert will show as below, click **Yes** to select the **Ranker** search method.

- **Step 2:** Choose **Ranker** by click Yes on Alert, because it is recommendation if we use CorrelationAttributeEval.
- **Step 3:** Choose **Use full training set** in attribute selection mode.
- **Step 4:** Click on **Start** button to process the data.
- **Step 5:** Get the result from Attribute Selection Output as the figure. So five **attributes with highest correlation** are **checking_status, duration, credit_amount, savings_status, housing.**

   c. **Preprocessing Data in Python (5 points)**

**Directory structure:**

- We have a data folder (containing input and output data for programs) at the same level as the src folder (containing program files). In there:
- **data folder:**
  - house-prices.csv: file input after the program is ready.
  - test.csv: file using to test during programming.
- **src folder:**
  - The files corresponding to the question (question1.py, ...)
  - utils.py: This file contains functions that are reused many times.



  - run.sh: file contains scripts when we execute it will execute the commands contained in the file (Here this file will run all the questions in one go with bash run.sh command).

1. **Extract columns with missing values:**
   - **Describe**: The requirement for extracting columns with missing values is to identify which columns have missing data and determine the appropriate way to handle them.
   - **Input**: Read file **house-prices.csv**
   - **Run:** python3 question1.py house-prices.csv

- **Output:** List of columns with any missing values

```
● (min_ml-env) tiendat@TienDat57:/mnt/d/WORK/Lab-DataMining/Lab01/src$ bash run.sh
Question 1: List of columns with any missing values:
 ['LotFrontage', 'Alley', 'MasVnrType', 'MasVnrArea', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1
', 'BsmtFinType2', 'FireplaceQu', 'GarageType', 'GarageYrBlt', 'GarageFinish', 'GarageQual', 'GarageCond',
'PoolQC', 'Fence', 'MiscFeature']
----------------------------------
```

2. **Count the number of lines with missing data:**
   - **Describe:** Counting the number of lines with missing data involves scanning each row in a data table and identifying which rows have one or more missing values. We can do this by iterating over each row in the data table and checking for missing values, which can be represented in different ways. The count of lines with missing data can be used as a quality metric for data analysis, and can inform decisions on data cleaning, imputation, or removal.
   - **Input:** Read file **house-prices.csv**
   - **Run:** python3 question2.py house-prices.csv
   - **Output:** Number of rows with missing values

```
----------------------------------
Question 2: Number of rows with missing values: 1000 samples
----------------------------------
```

3. **Fill in the missing value using mean, median (for numeric properties) and mode (for the categorical attribute).**
   - **Describe:** Filling in missing values involves replacing missing values with estimated values based on the available data.
   - **Input:** Read file **house-prices.csv**
   - **Run:**

python3 question3.py house-prices.csv --method=mean --columns=LotFrontage --out=fill_nan_values.csv

   - **Output: After full nan values → data save to csv** data/fill_nan_values.csv

```
----------------------------------
Question 3: Fill missing values by method =  mean for column =  LotFrontage  successfully!
----------------------------------
```

4. **Deleting rows containing more than a particular number of missing values (Example: delete rows with the number of missing values is more than 50% of the number of attributes)**
   - **Input:** Read file **house-prices.csv**
   - **Run:**

```
python3 question4.py house-prices.csv --threshold=0.5 --out=del_rows.csv
```

- **Output:**

```
-----------------------------------
Question 4: Deleting rows containing more than a particular number of missing values with threshold 0.5
successfully!
-----------------------------------
```

5. **Deleting columns containing more than a particular number of missing values (Example: delete columns with the number of missing values is more than 50% of the number of samples).**
    - **Input:** Read file **house-prices.csv**
    - **Run:**

```
python3 question5.py house-prices.csv --threshold=0.5 --out=del_cols.csv
```

- **Output:**

```
-----------------------------------
Number of columns before deleting:  81
Deleting column:  Alley
Deleting column:  MasVnrType
Deleting column:  FireplaceQu
Deleting column:  PoolQC
Deleting column:  Fence
Deleting column:  MiscFeature
Number of columns after deleting:  75
Question 5: Delete columns with more than  0.5  missing values successfully!
-----------------------------------
```

6. **Delete duplicate samples**
    - **Input:** Read file **house-prices.csv**
    - **Run:** `python3 question6.py house-prices.csv --out=del_dup.csv`
    - **Output:**

```
-----------------------------------
Question 6: Delete duplicate samples successfully!
-----------------------------------
```

7. **Normalize a numeric attribute using min-max and Z-score methods**
    - **Input:** Read file **house-prices.csv**
    - **Run:**

```
python3 question7.py house-prices.csv --method=min-max --columns=LotFrontage --out=normalize.csv
```

- **Output:**

```
-----------------------------------
Question 7: Normalize data by min-max method for column LotFrontage successfully
-----------------------------------
```

8. **Performing addition, subtraction, multiplication, and division between two numerical attributes.**
- **Input:** Read file **house-prices.csv**
- **Run:**

```
python3 question8.py house-prices.csv --method=add --columns=LotFrontage,1stFlrSF --out=add.csv

python3 question8.py house-prices.csv --method=sub --columns=LotFrontage,1stFlrSF --out=sub.csv

python3 question8.py house-prices.csv --method=mul --columns=LotFrontage,1stFlrSF --out=mul.csv

python3 question8.py house-prices.csv --method=div --columns=LotFrontage,1stFlrSF --out=div.csv
```

- **Output:**

```
-------------------------------------
Question 8: Performing add between two numerical attributes successfully!
-------------------------------------
Question 8: Performing sub between two numerical attributes successfully!
-------------------------------------
Question 8: Performing mul between two numerical attributes successfully!
-------------------------------------
Question 8: Performing div between two numerical attributes successfully!
-------------------------------------
```

## III.    REFERENCES

[1]. Handle missing values

[2]. weka_feature_selection

[3]. skewed

[4]. Weka_tutorial

[5]. Lecture slides

[6]. passing args

[7]. discrete-vs-continuous-data

[8]. calculate correlation