# DATA MINING
## ASSIGNMENT-1

# VISUALIZATION AND ANALYSIS OF THE GIVEN
# DATASET USING WEKA

**Prof: Dr Elizabeth D Diaz**

## TEAM MEMBERS

Harshini Kandimalla – 1001960046

Pratik Antoni Patekar – 1001937948

Pratik Dhanraj Chavan – 1001963580

# Introduction:

Weka (Waikato Environment for Knowledge Analysis) is a collection of machine learning algorithms for data mining tasks. It contains tools for data preparation, classification, regression, clustering, association rules mining, and visualization. Weka is open-source software issued under the **GNU General Public License**.
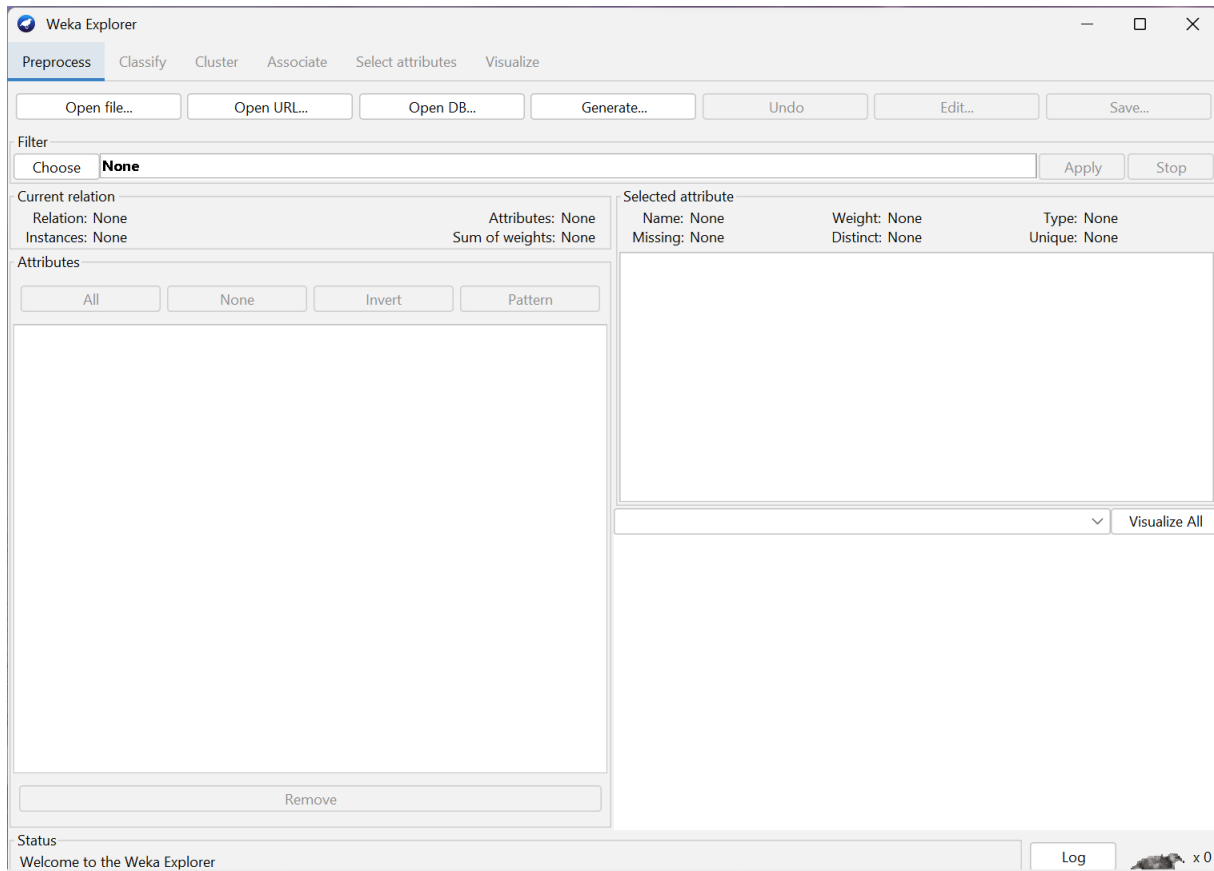
Weka supports **deep learning**!

In this, we have implemented visualization technique on Housing dataset.

# Retrieving Data:

Data must be loaded into weka. Dataset can be either csv file or arfftype. If we csv file, it can be uploaded directly or can be converted into arff type and then load and likewise.

1. Start the weka GUI chooser
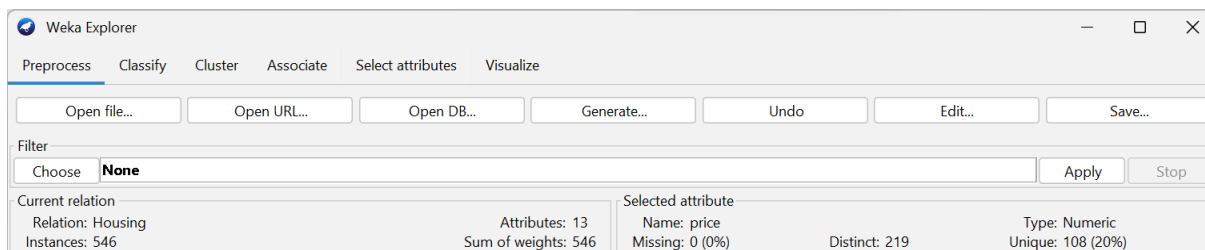
2. Launch the weeka explorer by clicking on explorer button

3. Now click on the open file to and it navigates to your directory and you can select the dataset to be used either in csv form or arff type.

# Glimpse of Data:

Once the file is opened in weka, we can check the data by clicking on the edit button in preprocessing tab.

From Current relation in above figure, we can get the details of attributes in the file and no of values it contains.

This is the dataset of Housing we are performing exploratory analysis in .csv file format.



Now, we are converting the .csv file to arff file and we can load data into weka

The arff type dataset is opened in the weka explorer where we can analyse each attribute in detail. For example, below we can see the details of attribute price such as the minimum value, maximum value, mean and standard deviation.



The following graph is about attribute bedrooms:

## Missing Data:

Once we load the data into explorer, on the right side in the selected attribute we can see missing from this we will know if there is any data missing and if there are any missing values, we click the button filter then click choose in that go to filter then unsupervised and then attribute in that we will find replace missing values you need to select that.

But in the dataset, we are using there is no missing data.



## Visualization of all attributes:

The data set on an attribute with respect to other attributes are show below. This can be seen by clicking on the tab visualize all.

- The X-axis and Y-axis represent the attribute

Visualization

1. The following graph shows bedroom over X-axis and Price over Y-axis
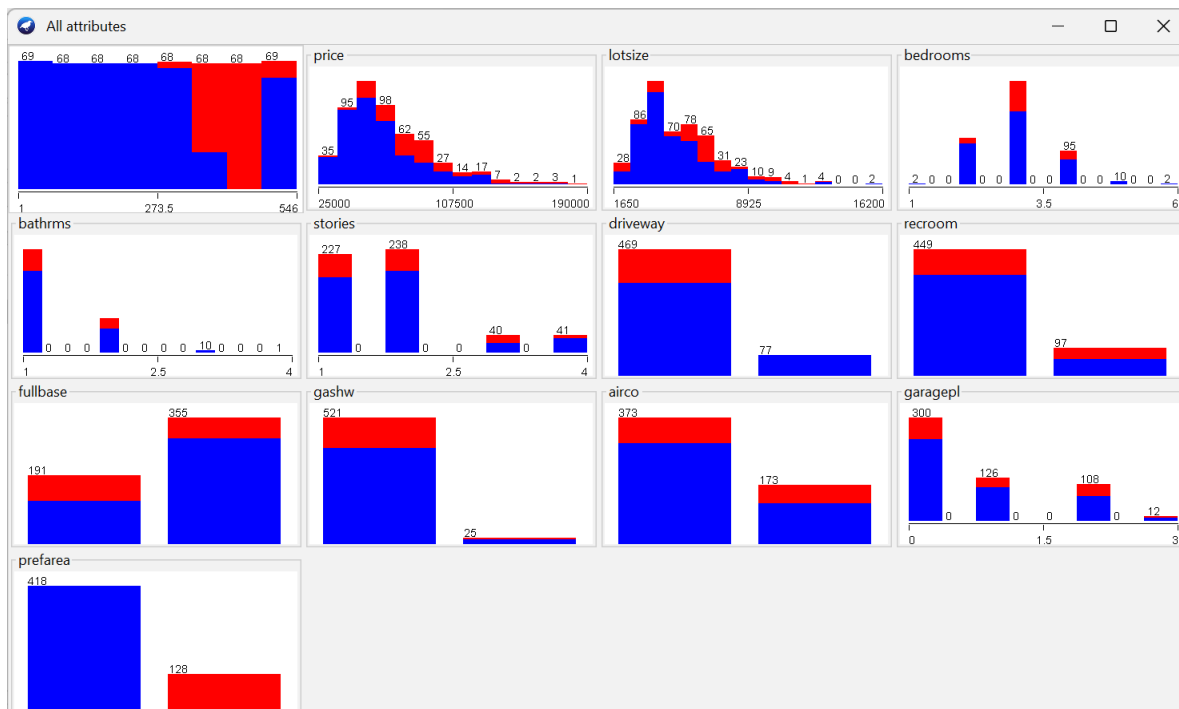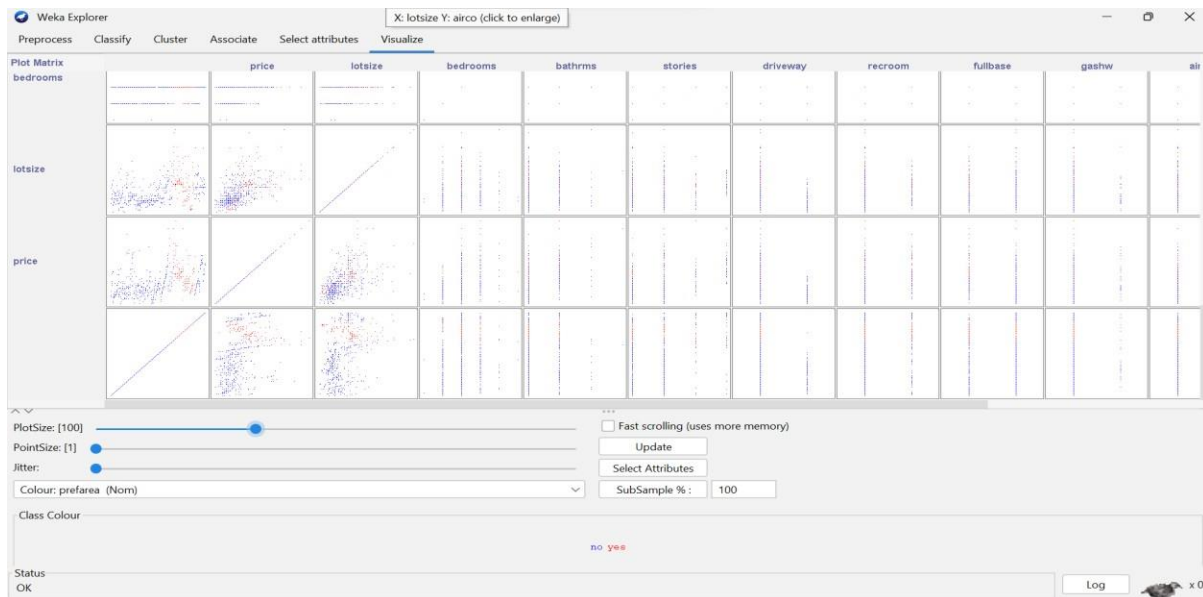


Price vs Bedrooms

The above plot shows no of bedrooms for a certain price. There are a lot of flats which can have over the pricing range between 25000 and 107500 with 2, 3, 4 and 5 rooms, but we have only two six bedrooms within this price. All the 2 bedrooms and most of 3 bedrooms are below 107500. From the plot we can say that we can get same no of bedrooms for all prices, but the size of each bedroom may not be same they can differ upon pricing. If we spend more money the size of the room might be large

1. Click on the instance represented by 'x' in the plot. It will give the instance details.
2. The X and Y-axis attributes can be changed in panel. We can view different plots by changing attributes in X and Y-axis.
3. Sometimes the points may overlap. The darker the spots we can say that they have multiple instances.



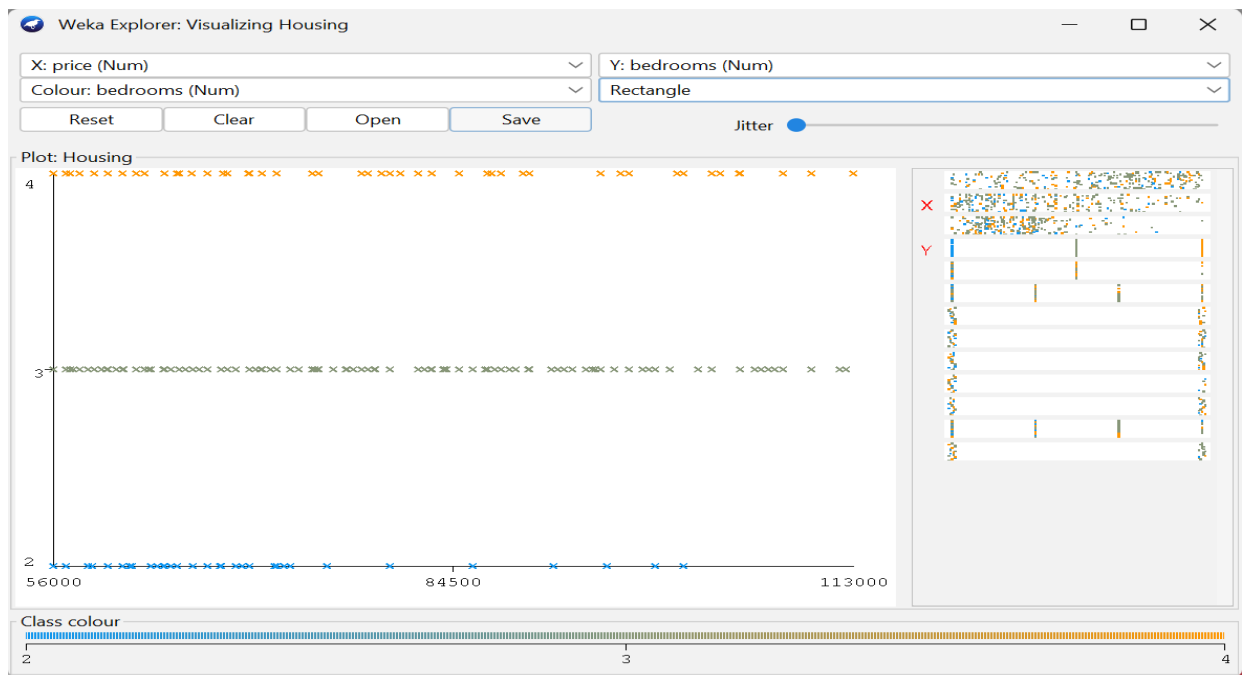To get the clearer view of the data, we can select an instance from dropdown. From below plot you can see we selected the rectangle, and we can get clear view of the data present in rectangle.
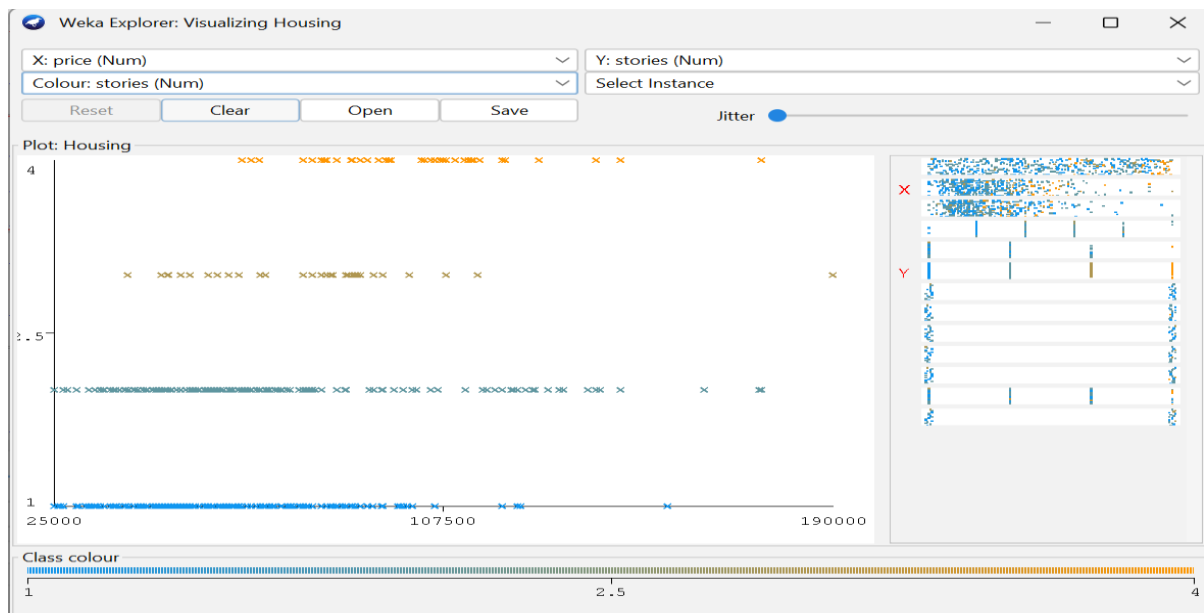
From above graph after selecting rectangle, you can click submit, then only the selected data will be displayed and the rest is cleared.

This is the elaborated graph of points selected in the rectangular box



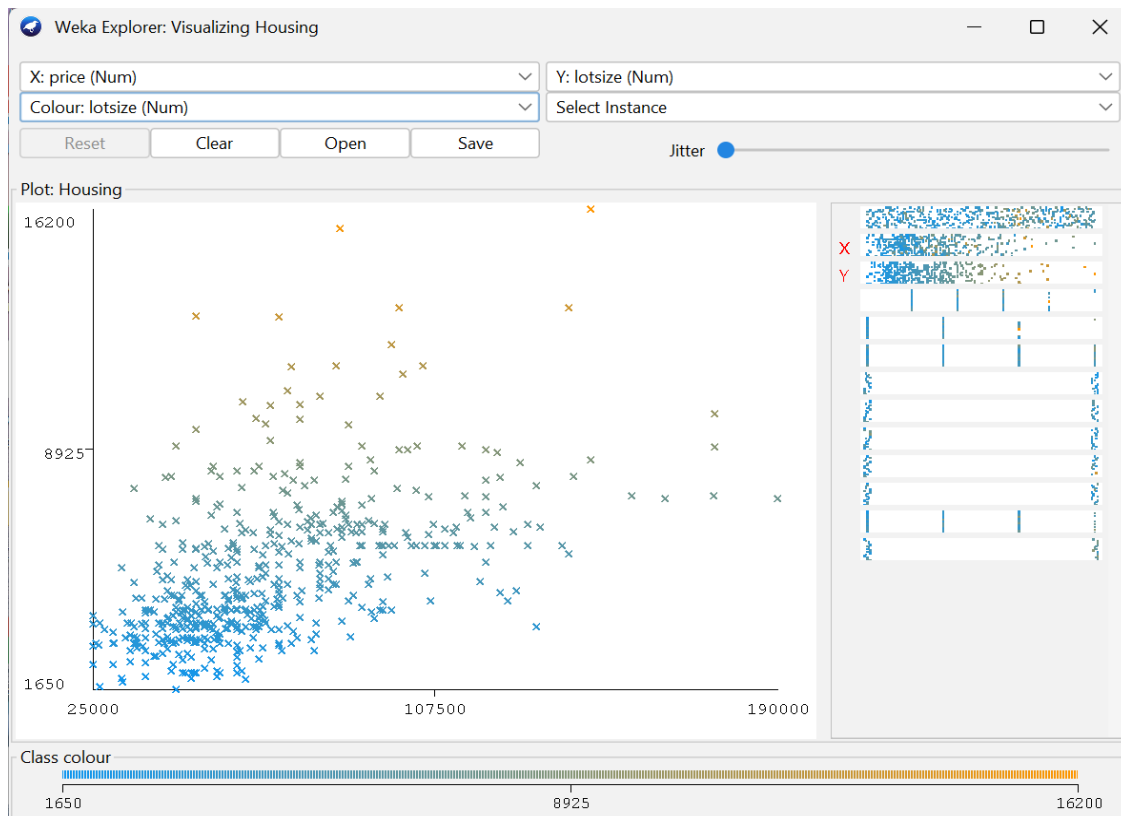2.The following graph shows plot between the price and stories.



Plot vs Stories

From the above plot we have price over X-axis and stories on Y-axis. There are most of the 2,3 and 4 stories building in the range between 25000 and 107500, and very few in the range of 190000. Even though we get more stories for the same price, the size of the every building may not be similar.

3. The below graph is price over X-axis and lotsize over Y-axis:

From the plot we can observe that the highest price i.e., 190000 has a lotsize around 7400 and the highest lotesize i.e., 16200 is around 142000. The most of lotsize between 1650 and 8925 has the pricing range between 25000 and 107500. From this we can say that all the Lotsize with in the same range have similar price and we get only few bigger lotsize with that similar price.



Price vs Lotsize

References:
- https://www.cs.waikato.ac.nz/ml/weka/
- https://cs.ccsu.edu/~markov/weka-tutorial.pdf
- https://www.softwaretestinghelp.com/weka-explorer-tutorial/

**Observation of other Teammates:**

Weka is a very efficient tool to perform data mining tasks and can experiment with new datasets. The interface and functionalities are very easy to use and have different analyses to choose. Almost every data can be loaded directly which reduces our work and worry about coding. Even if there is any missing data it can be seen in interface and can be replaced directly from the interface. It has many built-in features which can be modified as necessary.