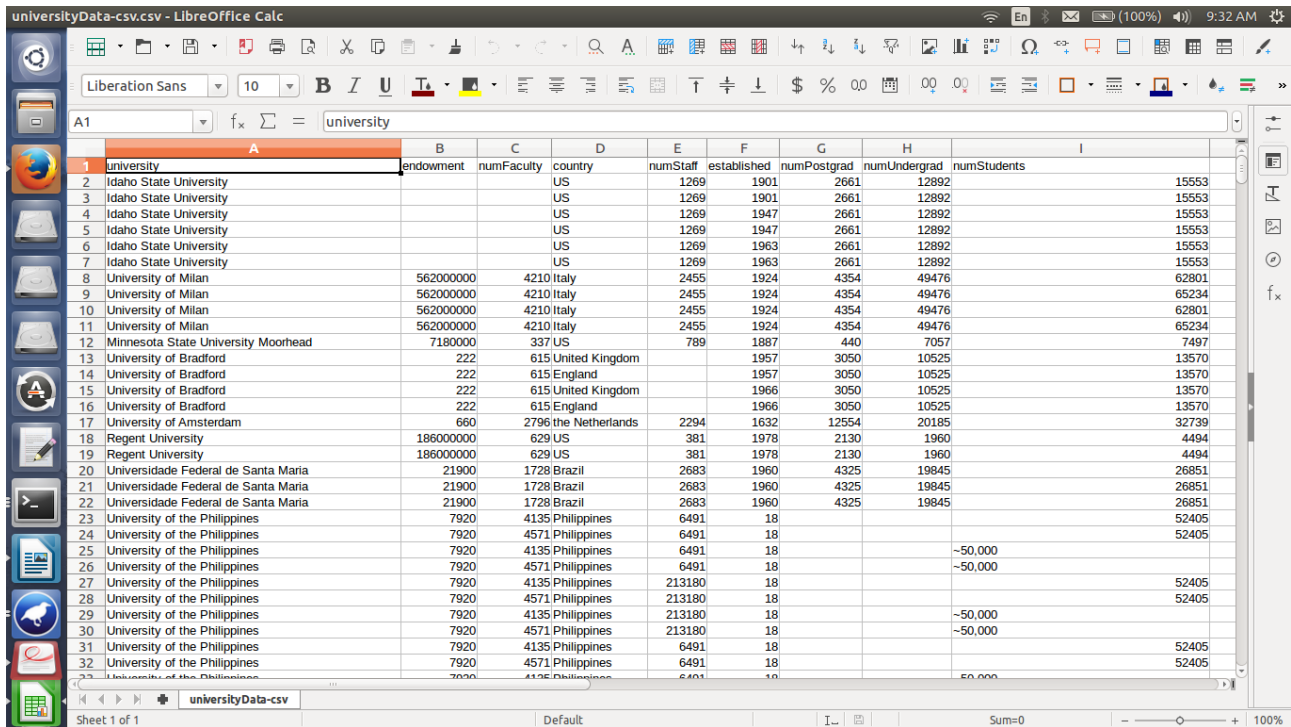# AMMANAMANCHI SAI KARTHIK

# B150310CS

## Data Set

The size of the data set has been reduced(10000 rows) and the data set is converted into proper CSV fomat and loaded into WEKA which uses ARFF format through the inbuilt converter.



## Format of Submission

The initial data set used is included as a .csv file and the final modified data set is availiable in .xls format.

**Raw Dataset**: universityData-csv.csv

> **Link**: https://drive.google.com/open?id=0B5bh5JNC6R13RktBR2pDOEJMX0U

**Pre-Processed Dataset**: karthik_b150310cs_dm.xls

> **Link**: https://drive.google.com/open?id=0B5bh5JNC6R13NGw4WFlfMldjQ00

# Missing values



We can replace all the missing values by going to **choose->unsupervised attribute filter->replace missing value** .

It automatically replaces the missing values in the data with the mean ot the mode of the dataset.

# CHANGING DATE FORMAT

In our data we have an attribute violation data which is of the form MM/dd/yyyy.The dataset  is clean regarding dates.

We can change date format by choosing ChangeDateFormat from unsupervised folder.



# REMOVING DUPLICATE DATA

To remove duplicate data we click on choose->filter->remove Duplicate Data under unsupervised category.

Before removing missing values.



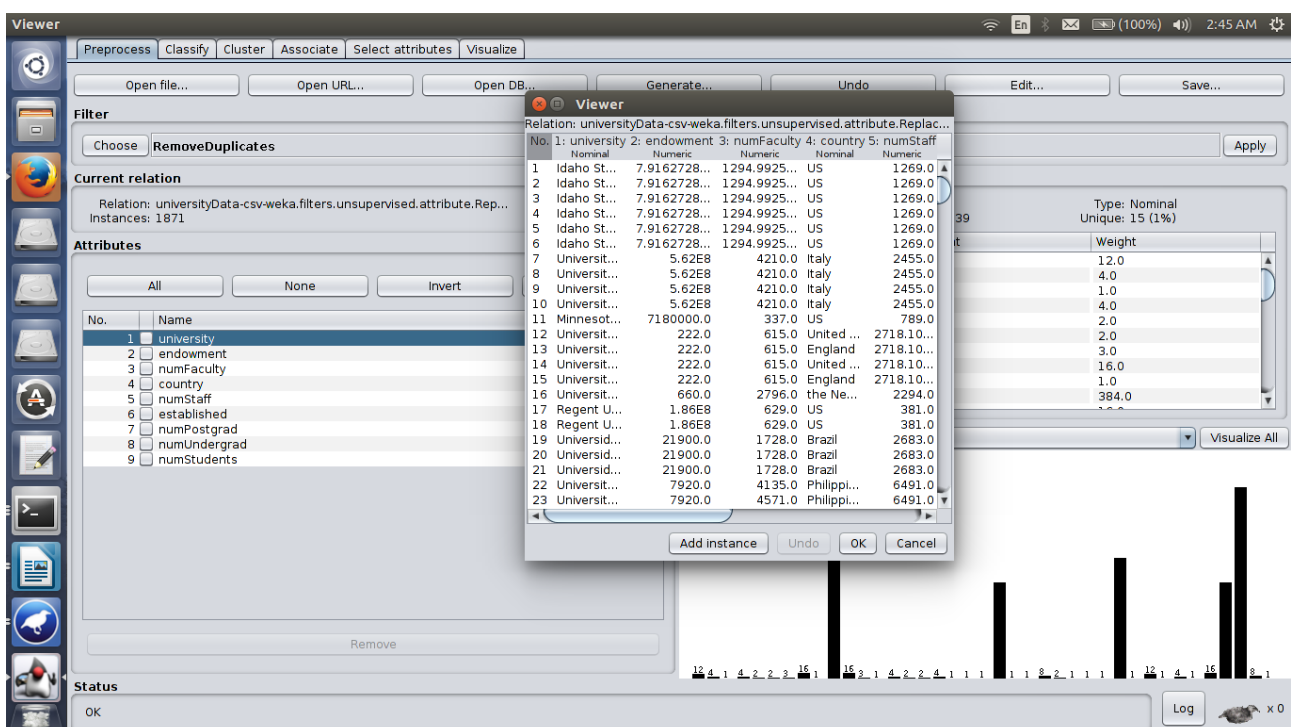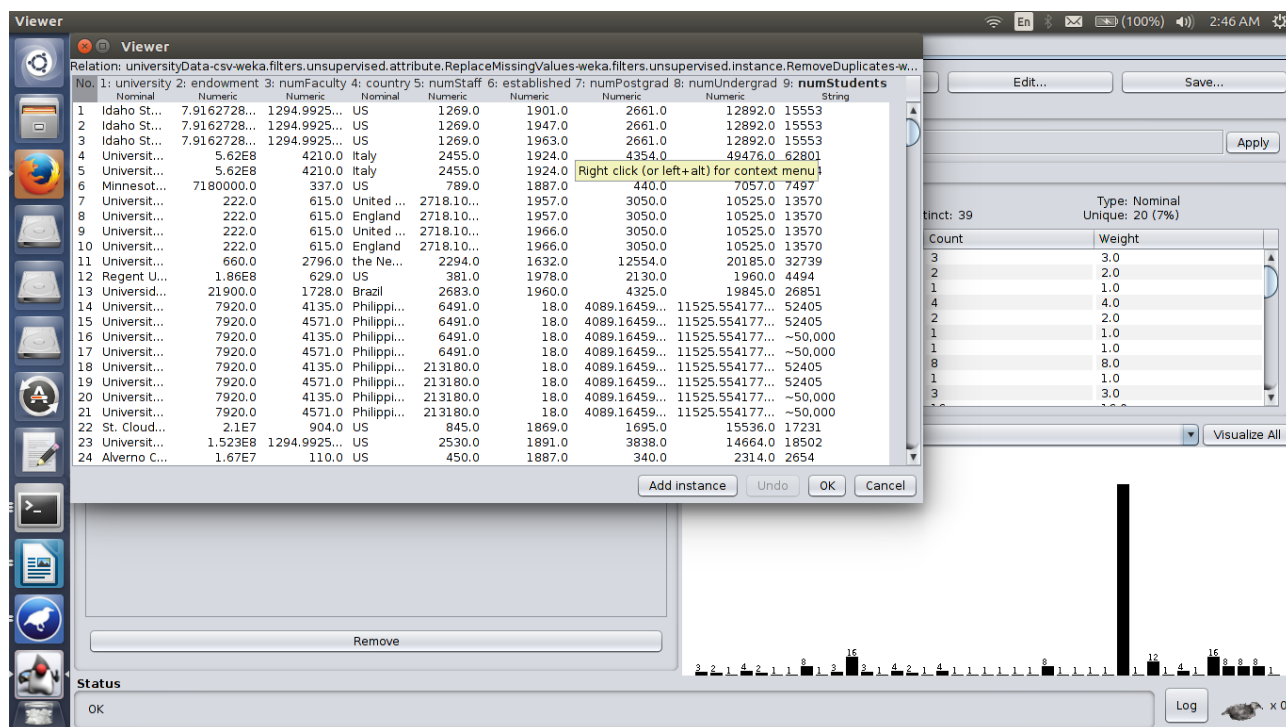After removing missing values.

Here two instances of a university are same as some of the mmissing values are replaced by means and mode which might come out to be different for differrent instances.

# SCATTERPLOT MATRIX

When attributes are numeric we can create a scatter plot of one attribute against another. This is useful as it can highlight any patterns in the relationship between the attributes, such as positive or negative correlations.So,Weka provides us Visualize tab for this purpose.

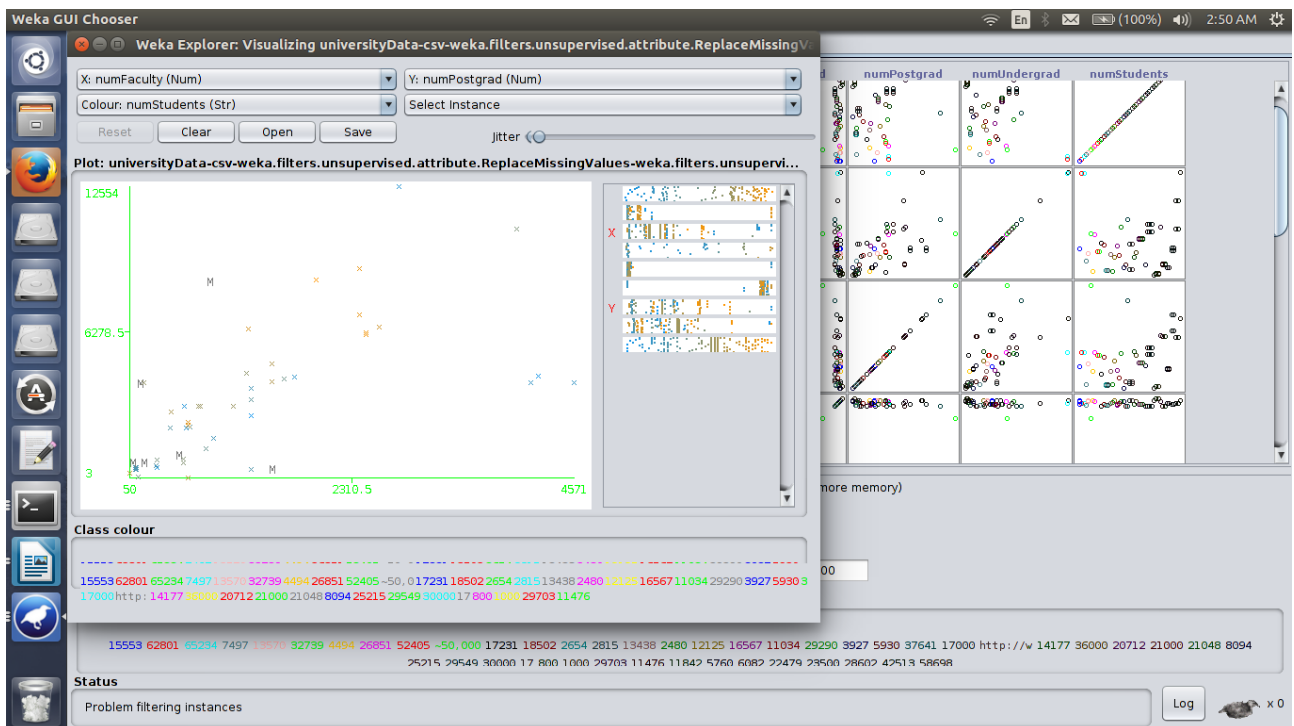The dots in the scatter plots are colored by their class value.Clicking on a plot will give you a new window with the plot .All combinations of attributes are plotted in a systematic way.
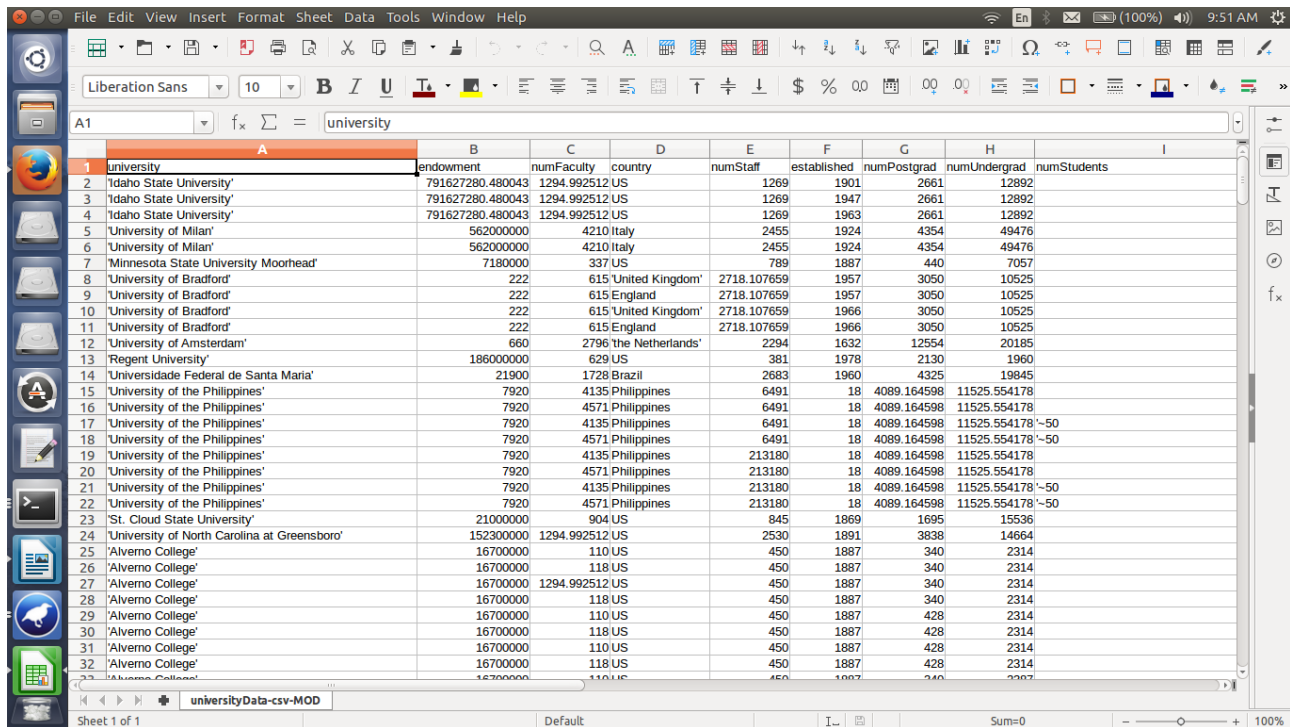
We can also see that each plot appears twice, first in the top left triangle and again in the bottom right triangle with the axes flipped. We can also see a series of plots starting in the bottom left and continuing to the top right where each attribute is plotted against itself.

Controls at the bottom of the screen. These let us increase the size of the plots, increase the size of the dots and add jitter.

# EXPLORE DATA

Right-click the model, choose "Visualize classifier errors" and Save the data (including the prediction).



OpenRefine does not have the option/pre-defined algorithm to fill in missing values.

## Advantages of open refine over weka :

1. In open refine we can edit a lot of instances at the same time by using the facet option but in weka it is not possible we can edit only one instance at a time.

2. Exporting the data from open refine to excel is very easy as it provided a direct option but in weka the file has to be saved in some format and then it has to be loaded onto excel.

3. quick, interactive, filter facets which allow for easy browsing of instances/rows which match a variety of filters

4.complete provenance/undo history of all modifications

5. wide variety of input & output formats including both file formats and online repositories like Google Spreadsheets & Fusion Tables.

## Advantages of weka over open refine:

1. Weka contains a large collection of data algorithms.

2. removal of attributes can be done easily in weka rather than open refine.

3. Weka provides a lot of attribute filter options such as discritisation,principal components,binarization,normalization etc..

4. Weka also provides a lot of attribute selection options such as correlation attribute evaluation,gainratio attribute evaluation,principal components etc..

5.Weka contains various data transformation techniques such as nominal to binary Nominal to string etc..

6. Missing values can be replaced for all attributes at once in weka but in open refine it has to be done separately for each attribute.