# Ammanamanchi Sai Karthik

# B150310CS

## 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.
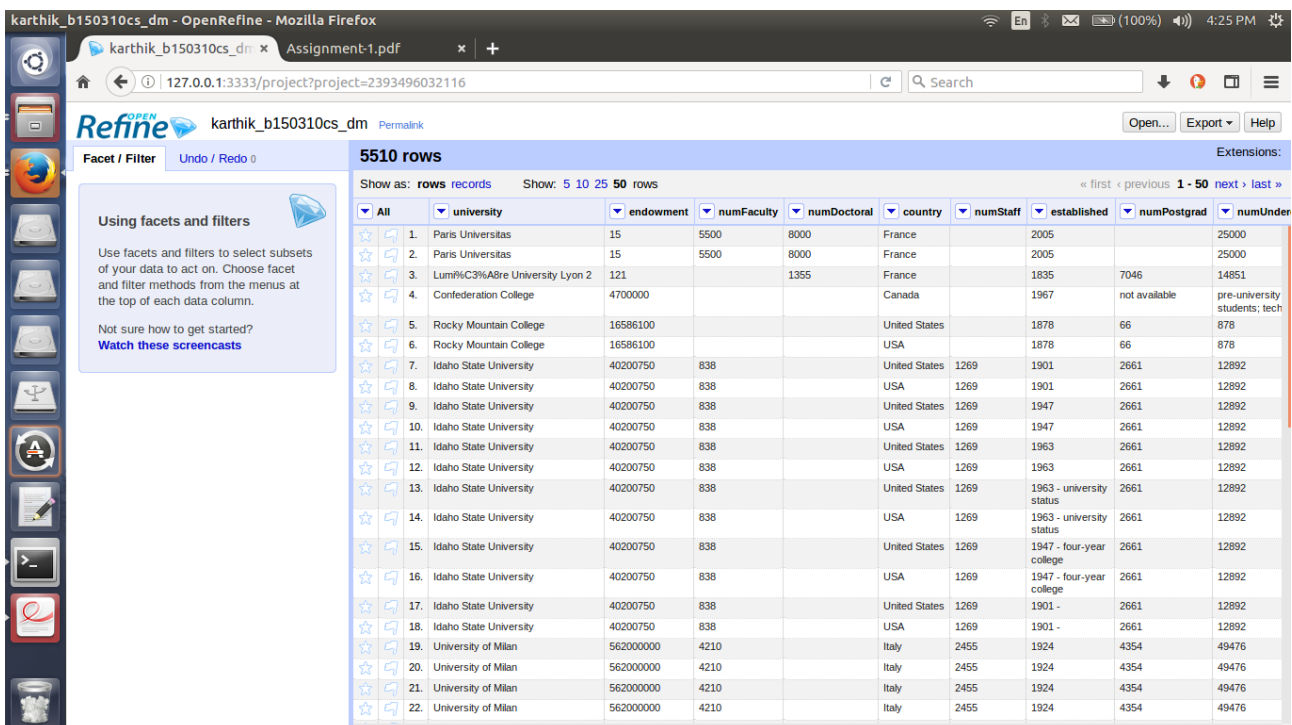
**Raw Dataset**: universityData.csv

    **Link**: https://drive.google.com/open?id=0B5bh5JNC6R13SWlRS2puMWZHZFk

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

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

## DataSet

The dataset consists of University data with over 5000 rows and 10 attributes.



## Attributes

University(Nominal), Endowment(Numeric), numFaculty(Numeric), numDoctoral(Numeric), Country(Ordinal),numStaff(Numeric),Established(Ordinal),numPostgrad(Numeric), numUndergrad(Numeric).

# Clean up country names

To sort out the issue of country names like USA, U.S.A, U.S we can use **Edit cells->Cluster and edit** on the country column.
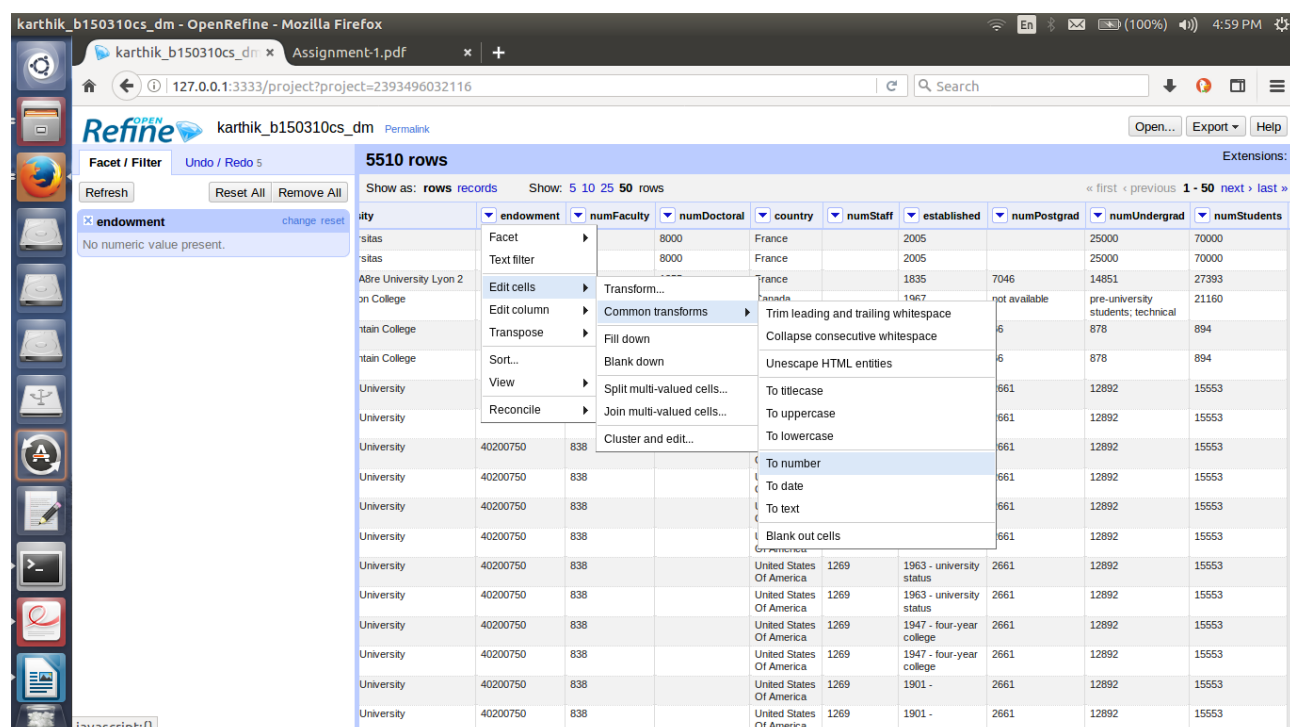


We use an inbuilt clustering algorithm with key function fingerprint and rename the country name to United States Of America.

# Converting Values such as $123 million to 123000000

We will focus on the attribute endowments.

Firstly, we will have to convert the strings to numbers using the common transformation tonumber.

We now change the US$ and US $ to ""(empty string) using value.replace("US $","").replace("US$", "") in tthe transform section of endowments attribute.
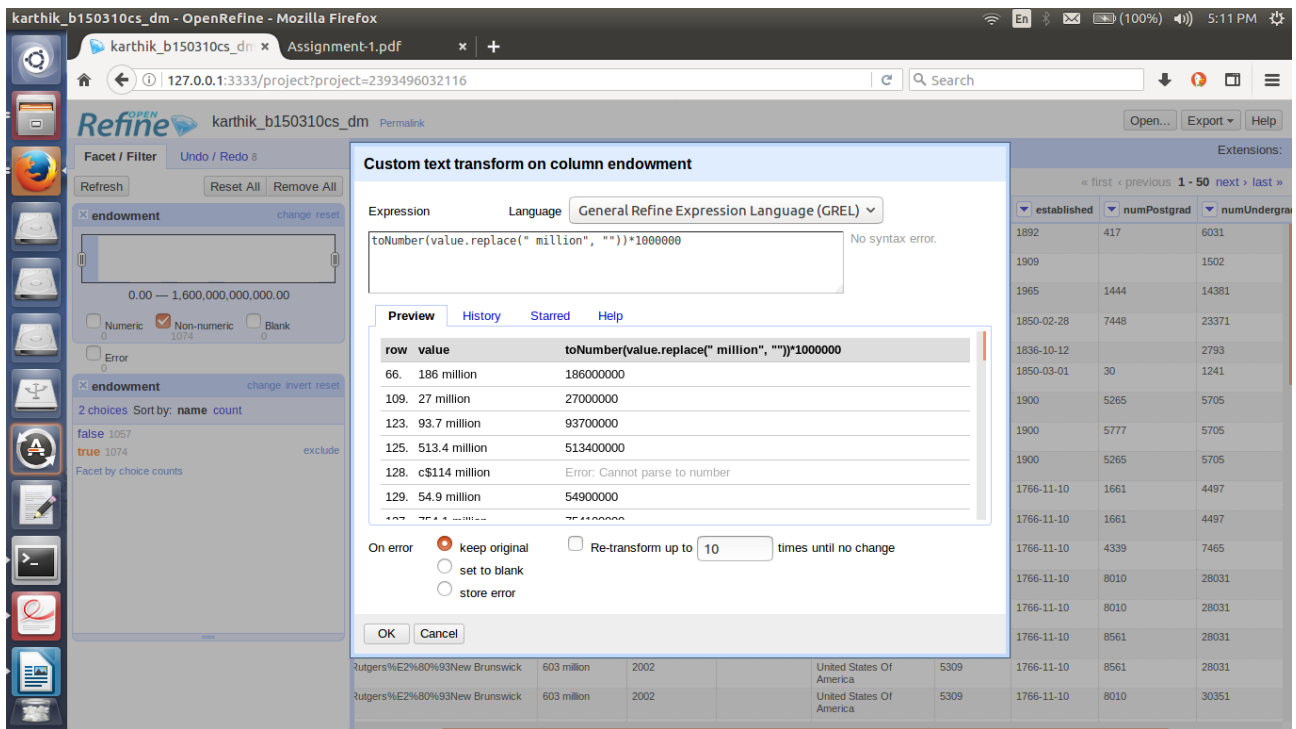


Now we convert the "millions" and "Millions" to the lowercase "millions".
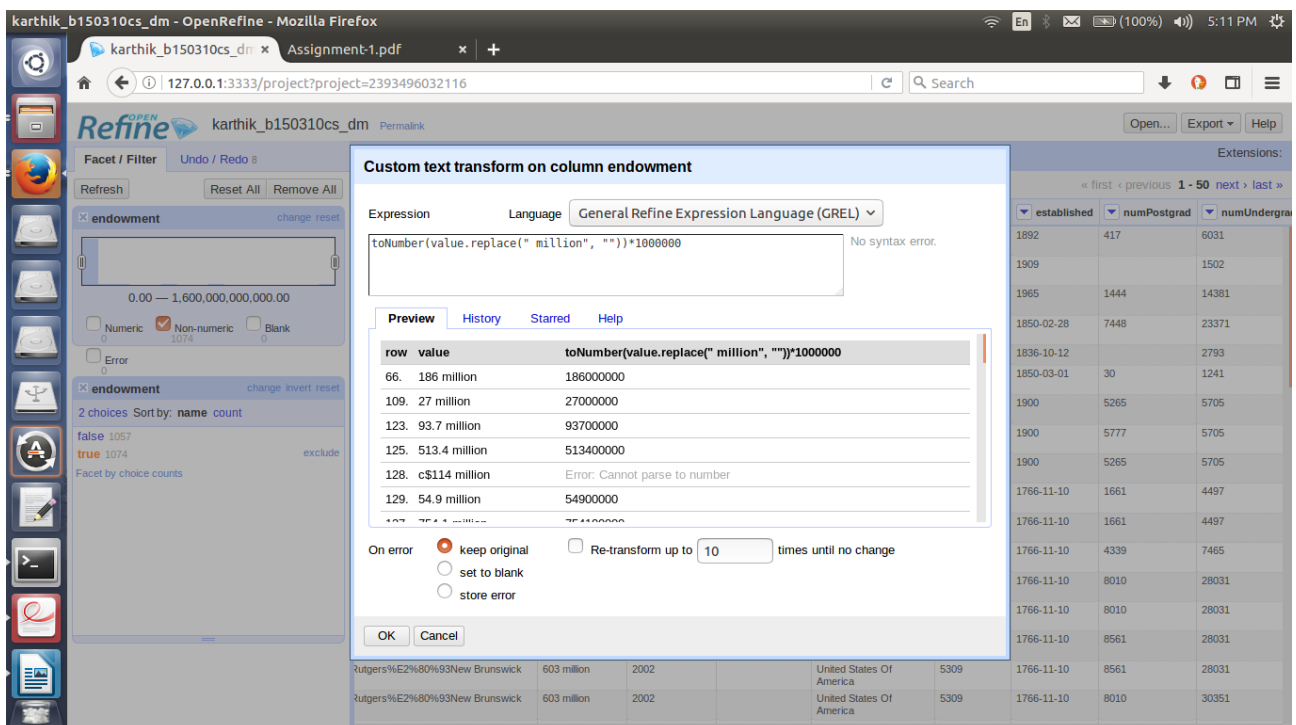


Create a custom text facet to select the rows containing "million".
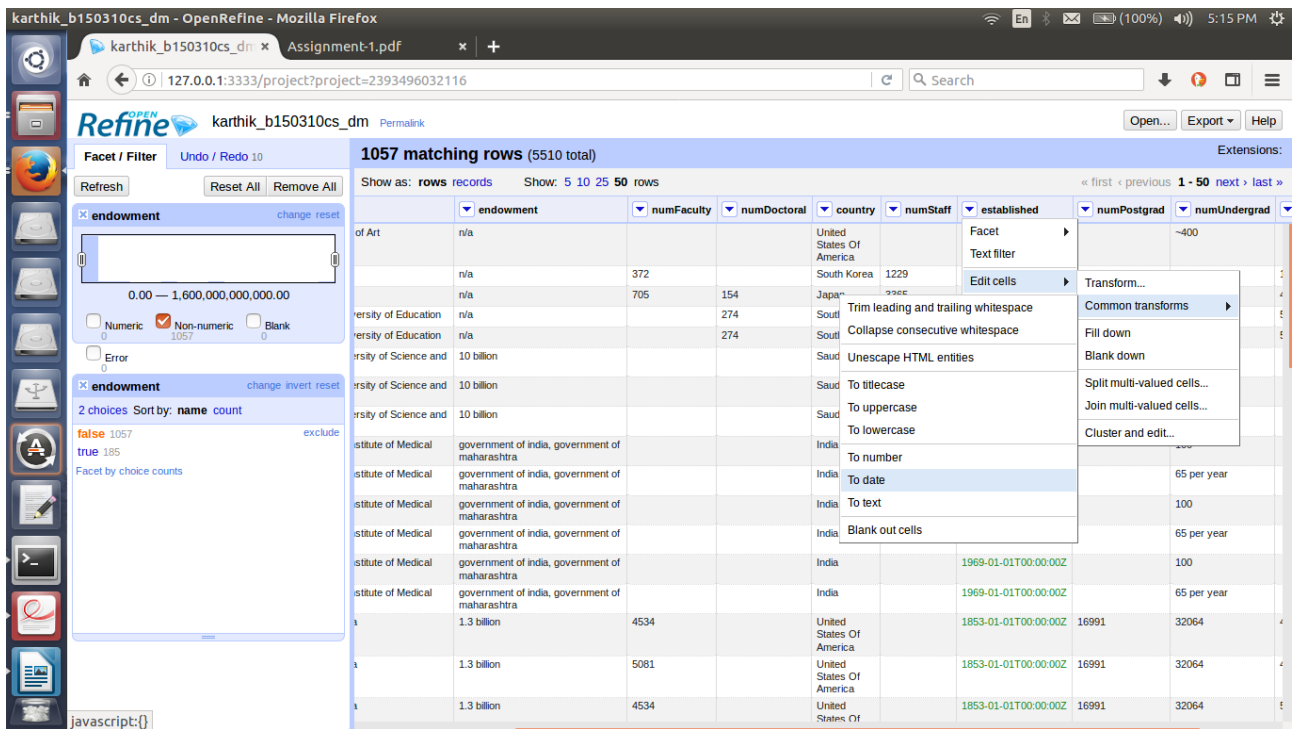
```
value.contains("million")
```

We now replace the million with a ""(empty string), convert the resulting string to number and multiply by 1000000.

```
toNumber(value.replace(" million", ""))*1000000
```



# Cleaning up dates

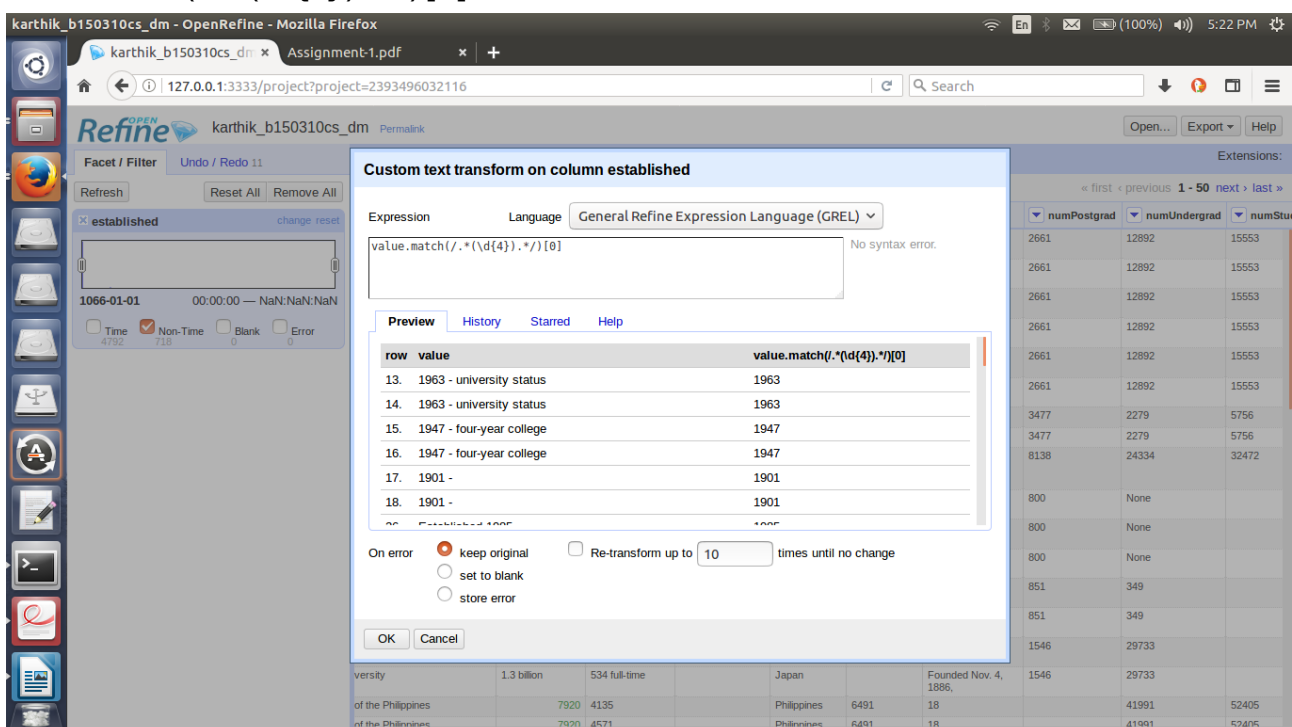First we will use a common transformation to Date on established column.

Now we select the non date tuples in the column by using the timeline facet and selecting the non-time dialog box.

**Facet -> Timeline facet**

Now we use a regular expression to select the first four characters in the string and adain use the to Date common transformation to convert them dates.

```
value.match(/.*(\d{4}).*/)[0]
```

Use transform on the column and apply value.toString('yyyy').



# Deduplicate Entries

We will try to deduplicate entries by tackling the university names column.

We first sort the rows and select reorder rows permanently.

We now edit the cells by **Edit cells -> Blank down** property which blanks out the row below another row if they are equal.
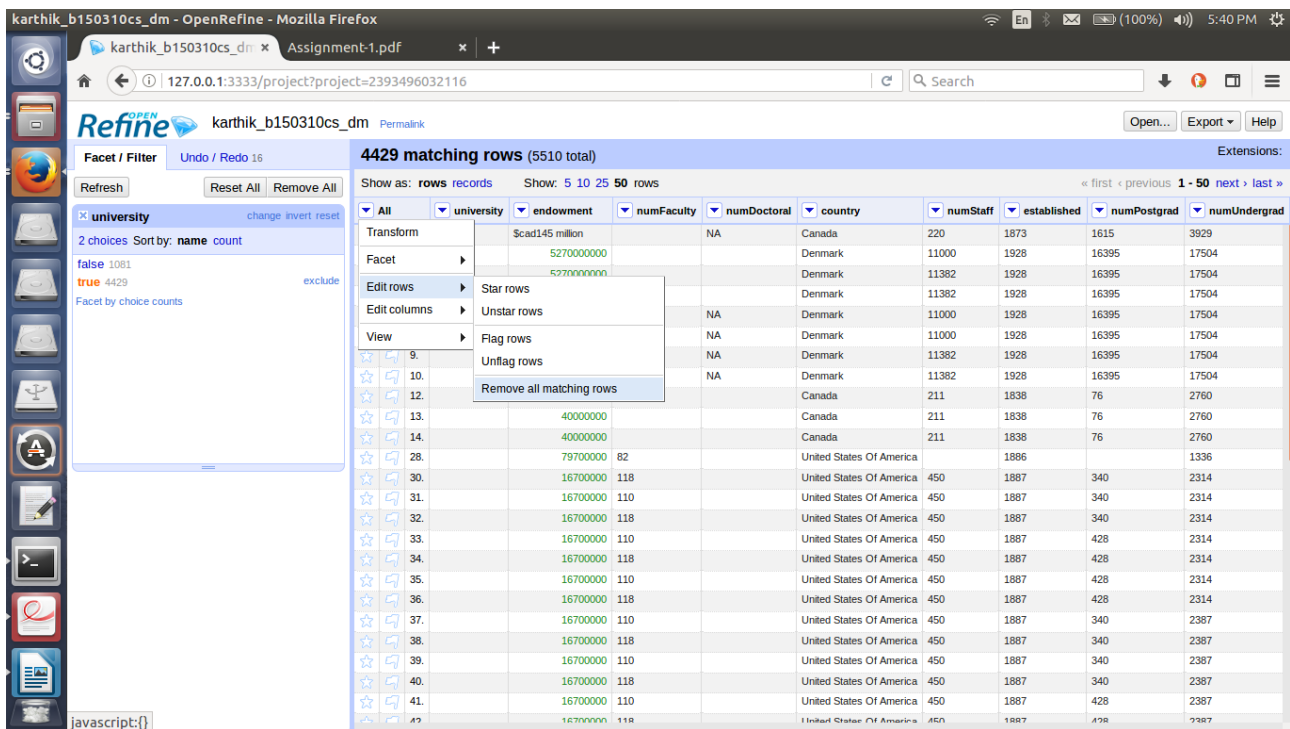


Now we use a customized facet Facet By Blank proprty(**Facet -> Customized facets -> Facet by blank)** and select all the blank rows and remove all the blank rows.

# Exploring the data with scatter plots

Click on the "endowment" column, **Facet -> Scatterplot facet**.

Use this scatterplot on all the attributes and generate the scatter plot.
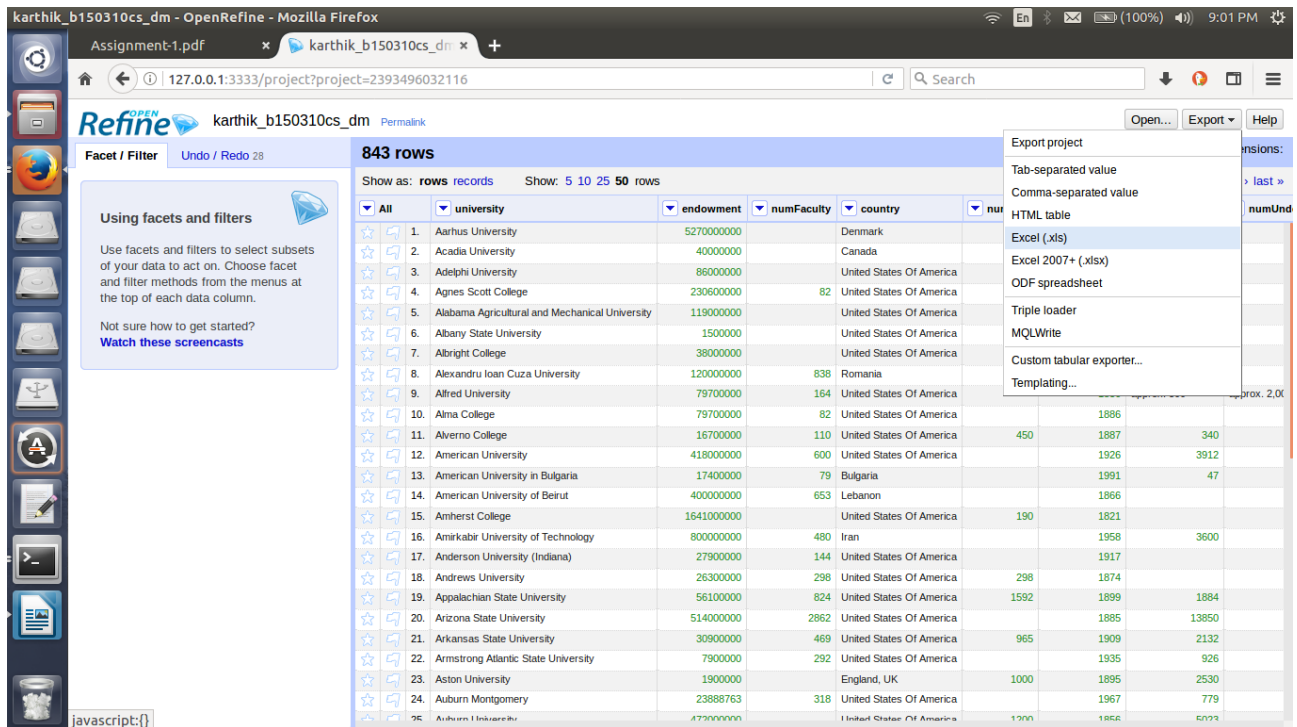
# Missing values

Open Refine does not have an option to fill the missing values whereas weka has.

So we have to fill out the missing values manually in OpenRefine.

# Exporting data to an Excel sheet

Export the final dataset in .xls format.



# Final Data