

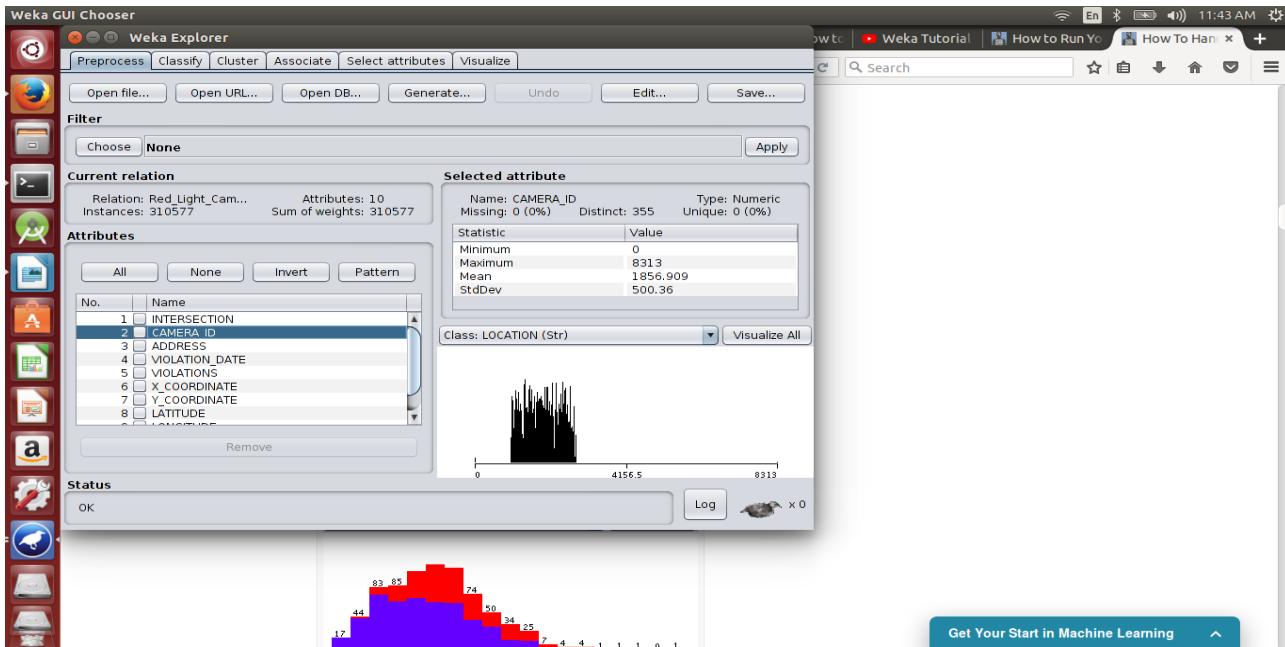
DM ASSIGNMENT

Question-2:

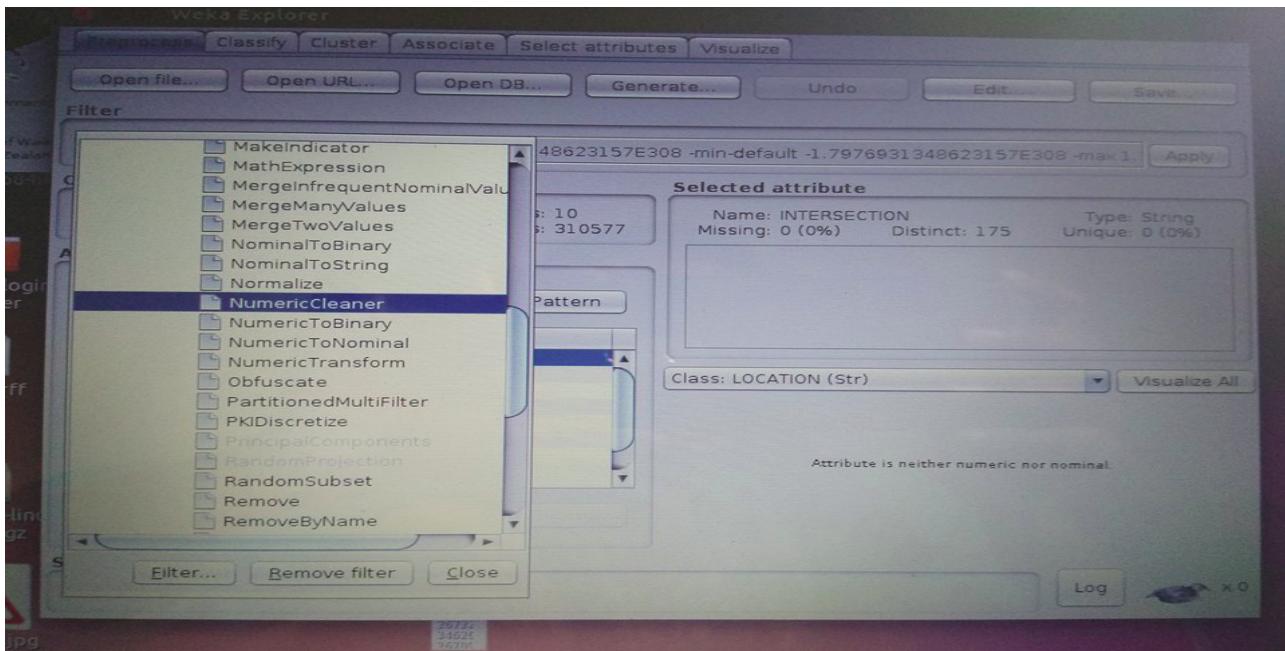
The dataset used is same as Question1.

1) Cleaning Non-numeric Camera Id's :

Open the Weka Explorer. Load the dataset.

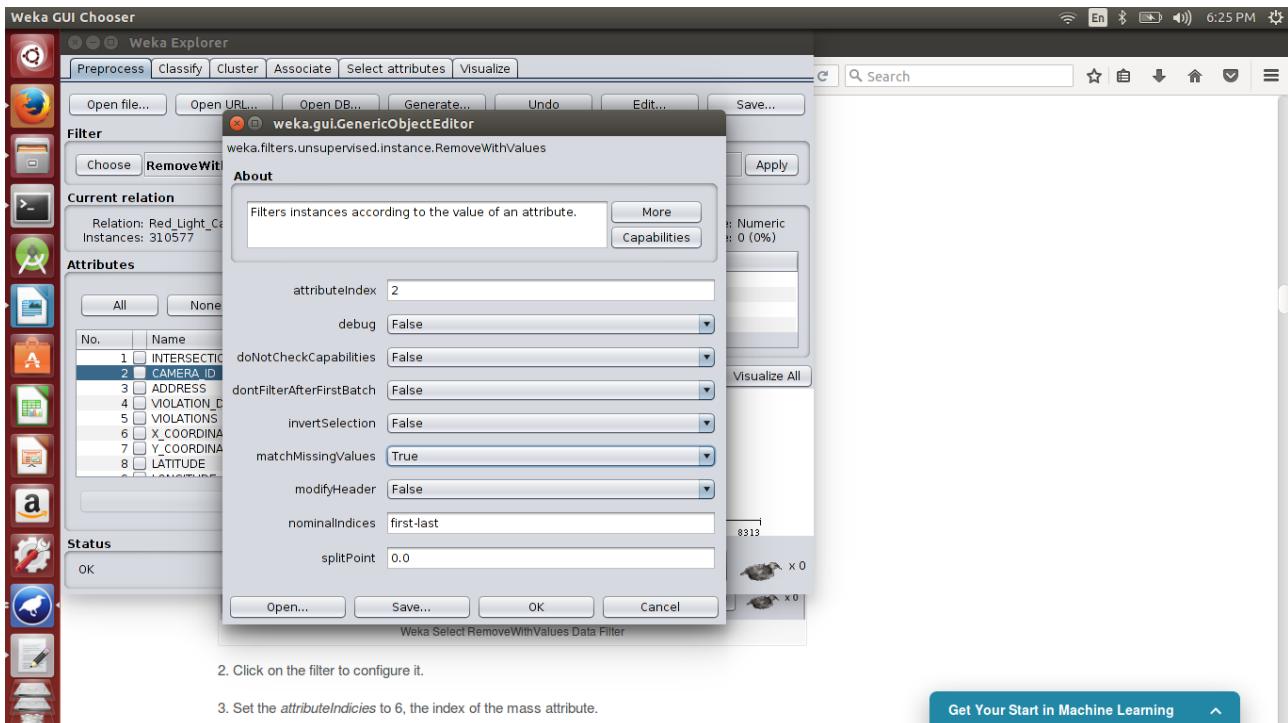


Click the “Choose” button for the Filter and select NumericalCleaner, it us under unsupervised attribute.



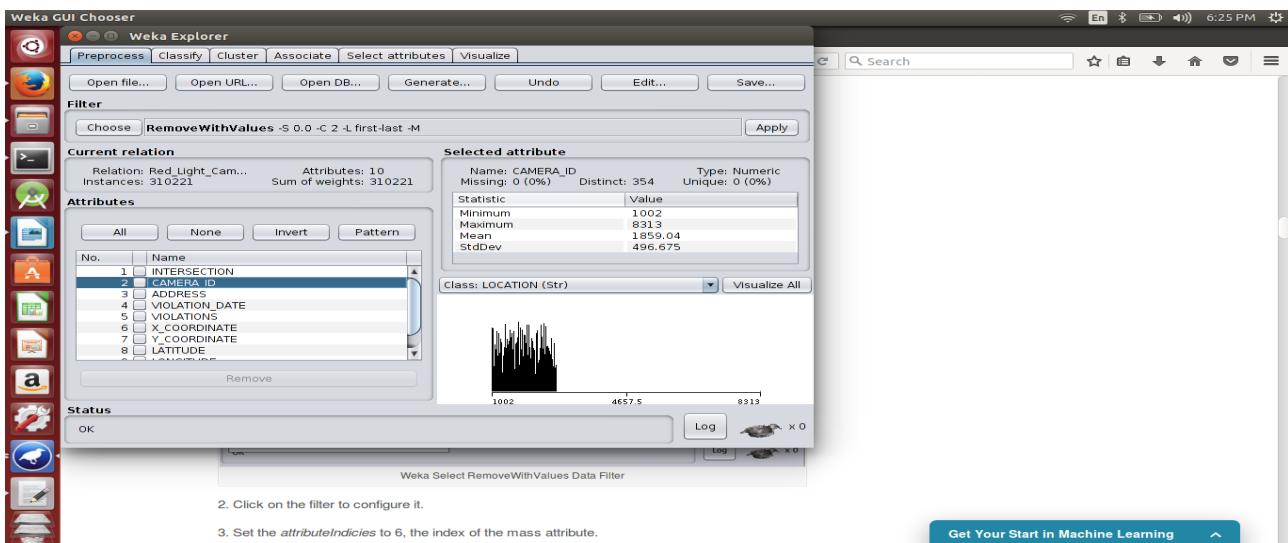
As Camera_ID is expected to be unique we remove nonnumeric and blank id's. Set *minThreshold* to 0.1E-8 (close to zero), which is the minimum value allowed for the attribute. Set *minDefault* to NaN, which is unknown

and will replace values below the threshold. So, this marks the missing values. Click the “Choose” button for the Filter and select RemoveWithValues, it's under unsupervised.instance.RemoveWithValues. Click on the filter to configure it. Set the *attributeIndices* to 2, the index of the Camera_id attribute. Set *matchMissingValues* to “True”.



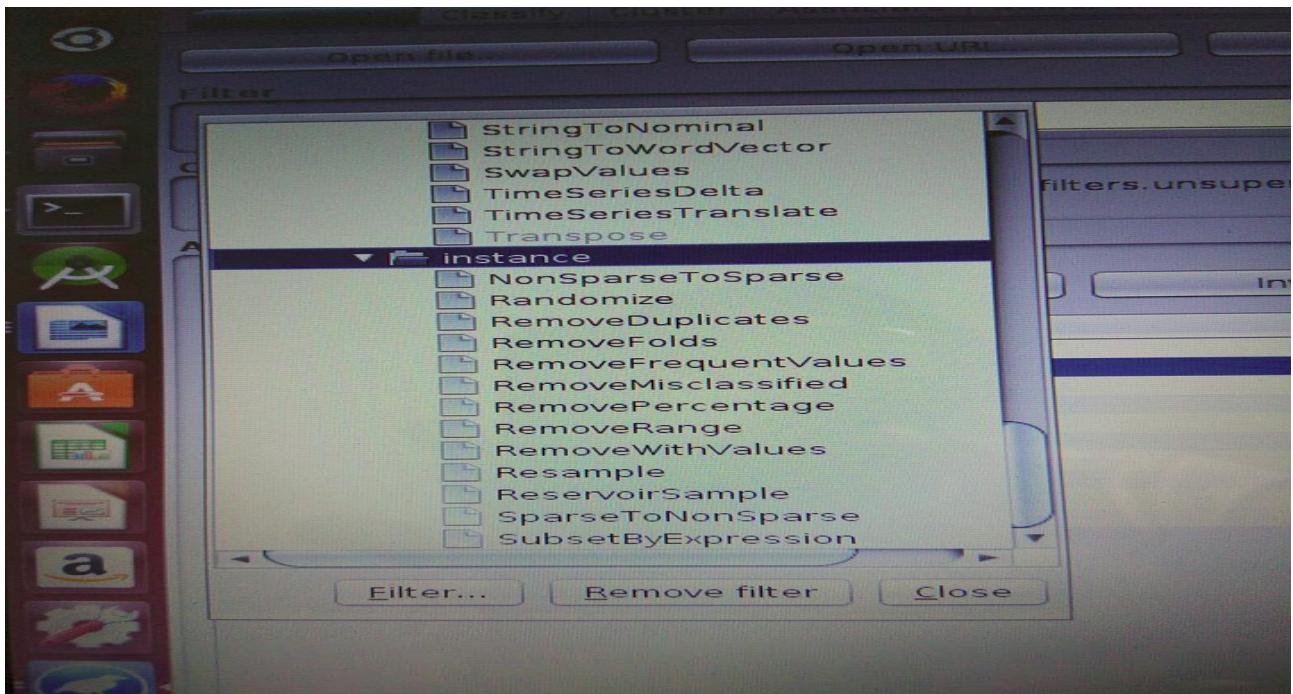
Click the “OK” button to use the configuration for the filter. Click the “Apply” button to apply the filter.

Now the unwanted rows got deleted.



2) Removing Duplicate Data:

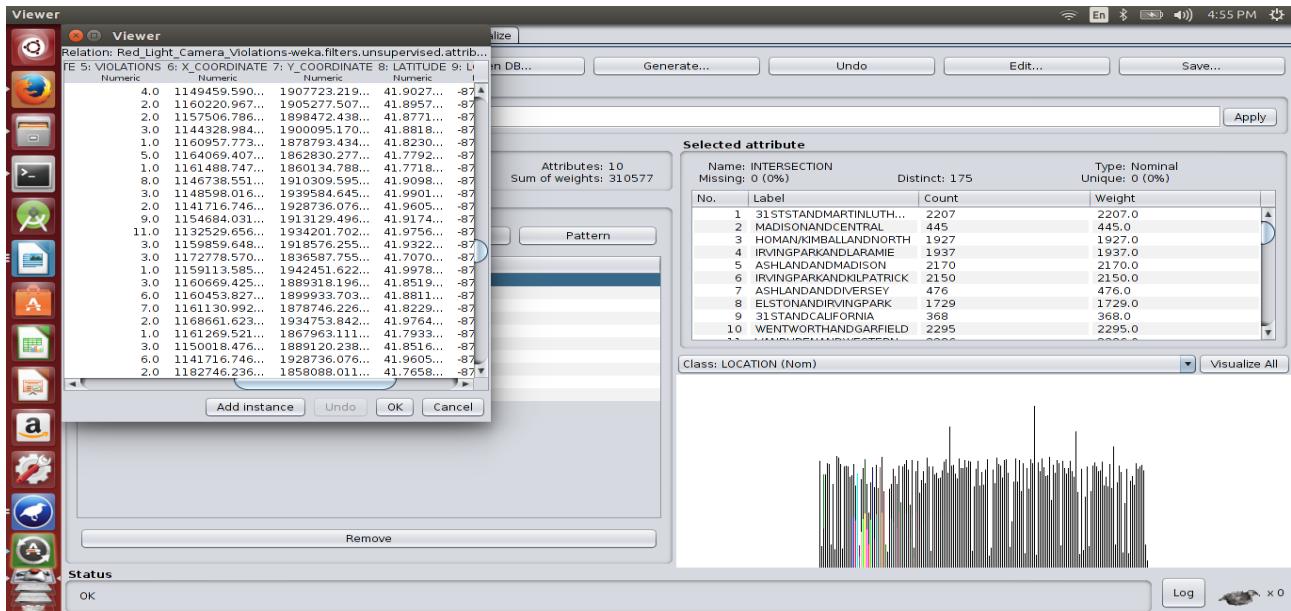
Click the “Choose” button for the Filter and select Remove Duplicate Data, it is under unsupervised.instance.RemoveDuplicate Data.



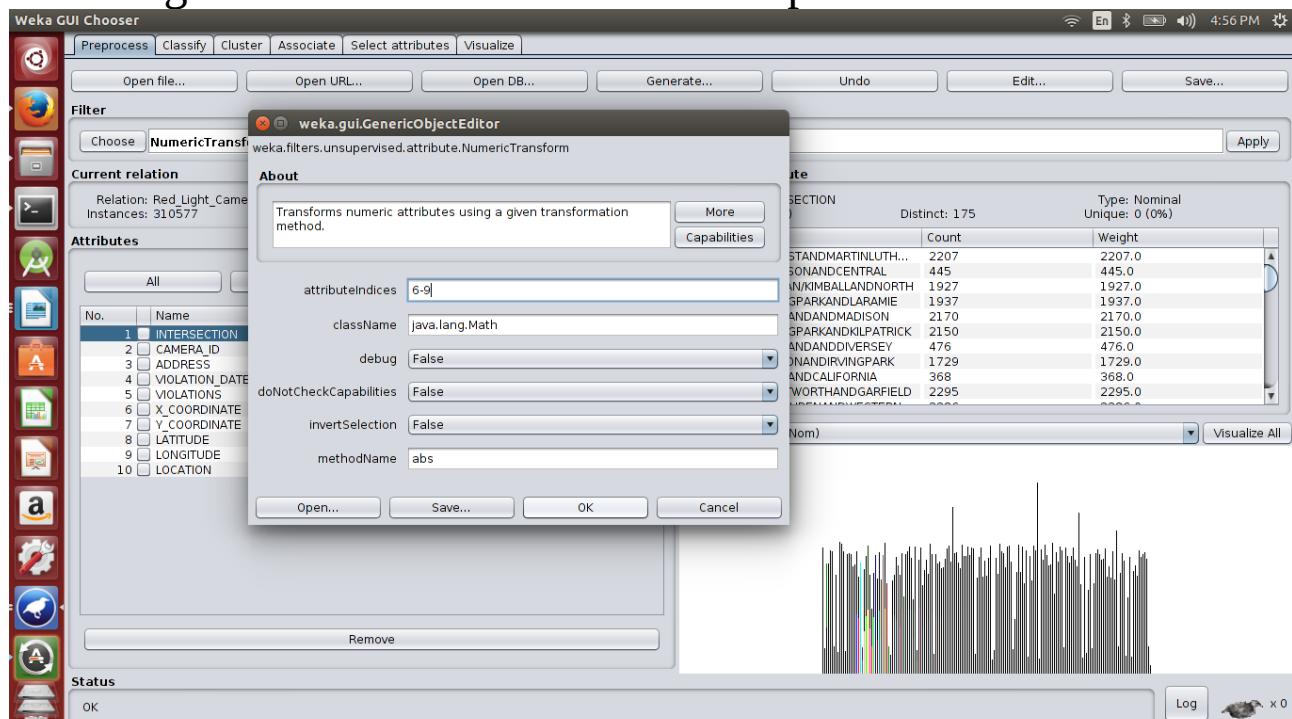
3) Numeric Transform of Attributes:

As the given data consists all numeric attributes and as we filled the missing values with mean of all other values ,so as the average will be in decimal in case if we need it to convert to absolute value or floor of the value we use numeric transform. We can select this from unsupervised.attribute.numeric transform

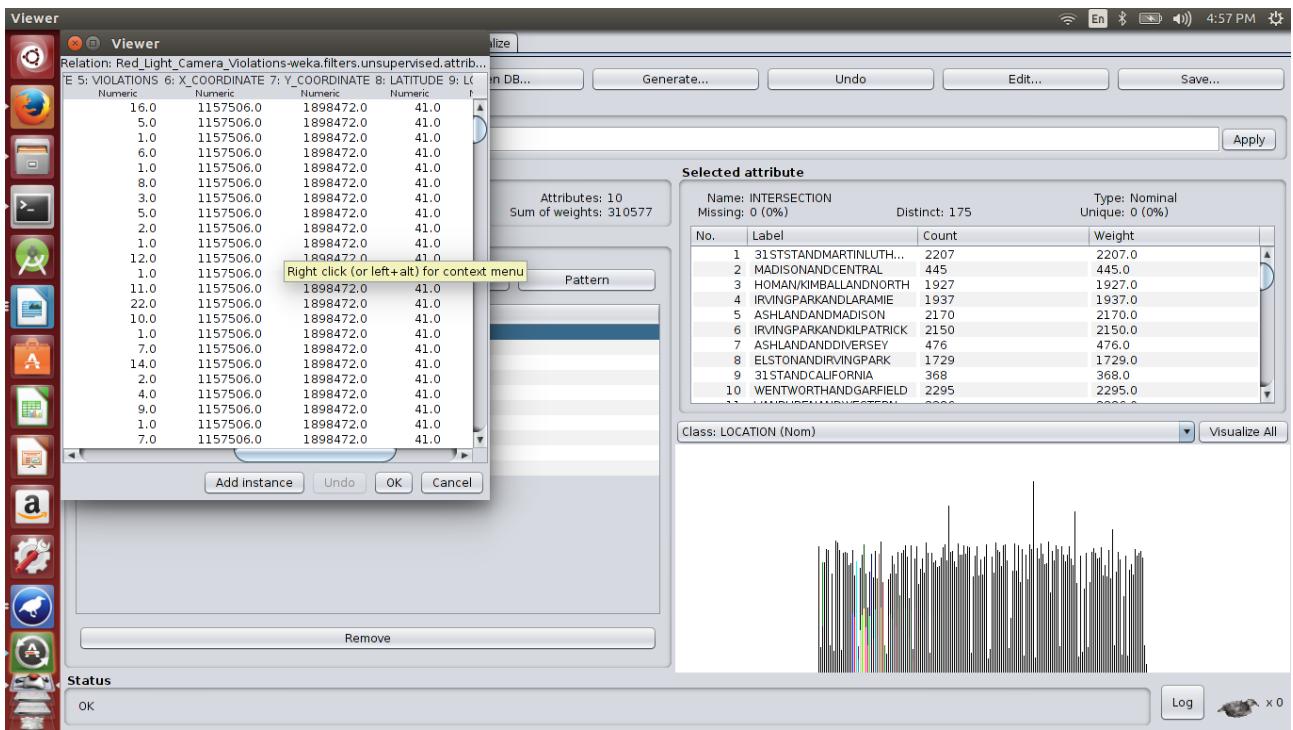
Before numeric transform:



Selecting Numeric Transform from unsupervised Attribute:



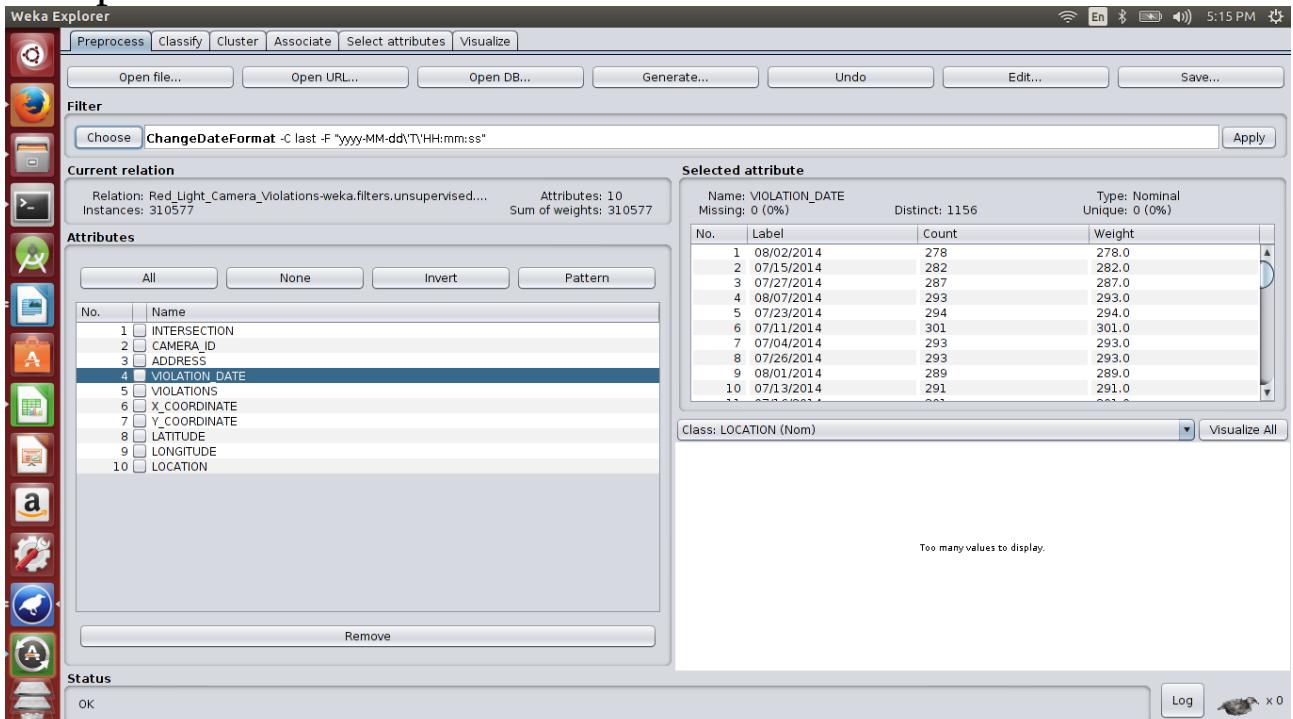
Now after selecting method name as floor the data looks like:



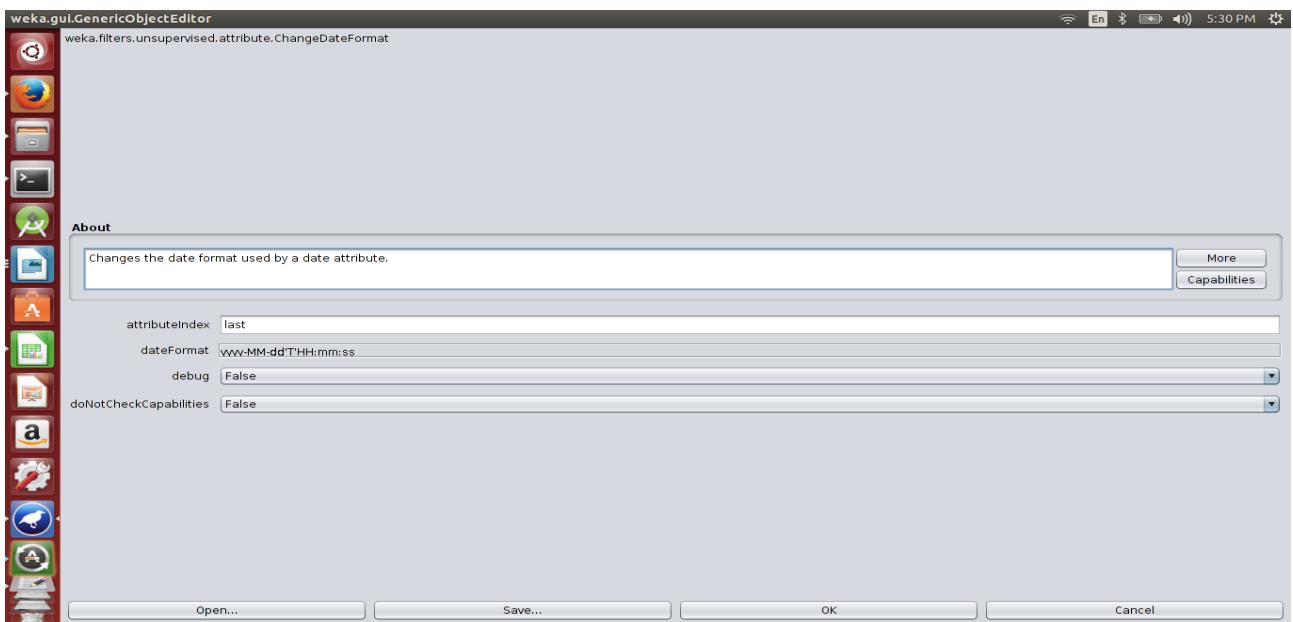
4) Changing Date format:

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 unsupevised folder.



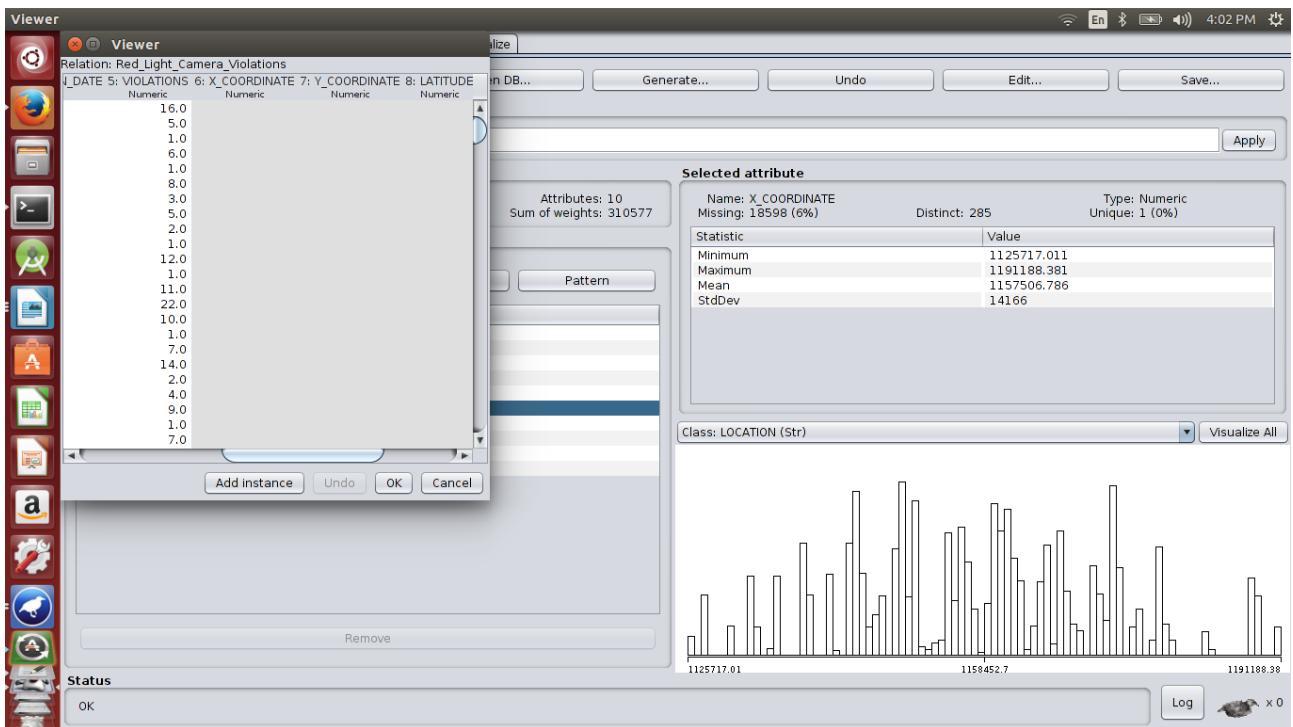
Changing date format:



5)FILLING MISSING VALUES:

Data Set initially with missing values:

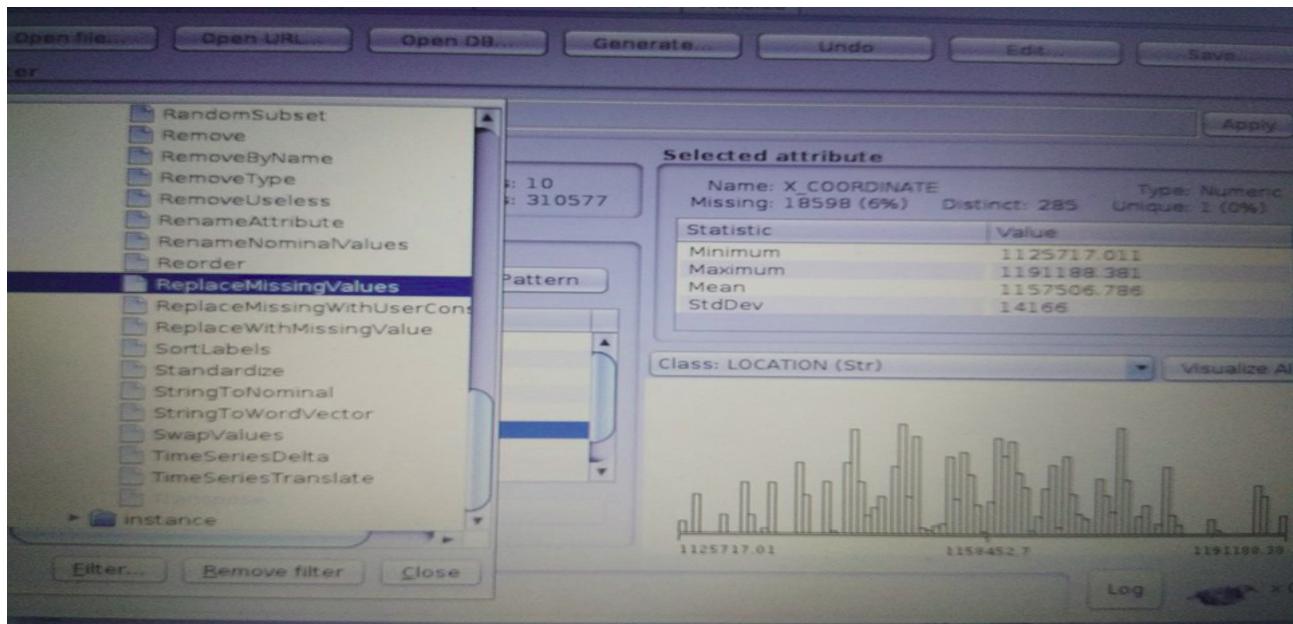
Let us fill missing values for X-coordinate .So,We select X-coordinate from preprocess tab and proceed.



ReplaceMissingValues filter. Click the “Choose” button for the Filter and select ReplaceMissingValues, it is under unsupervised.attribute.ReplaceMissingValues.

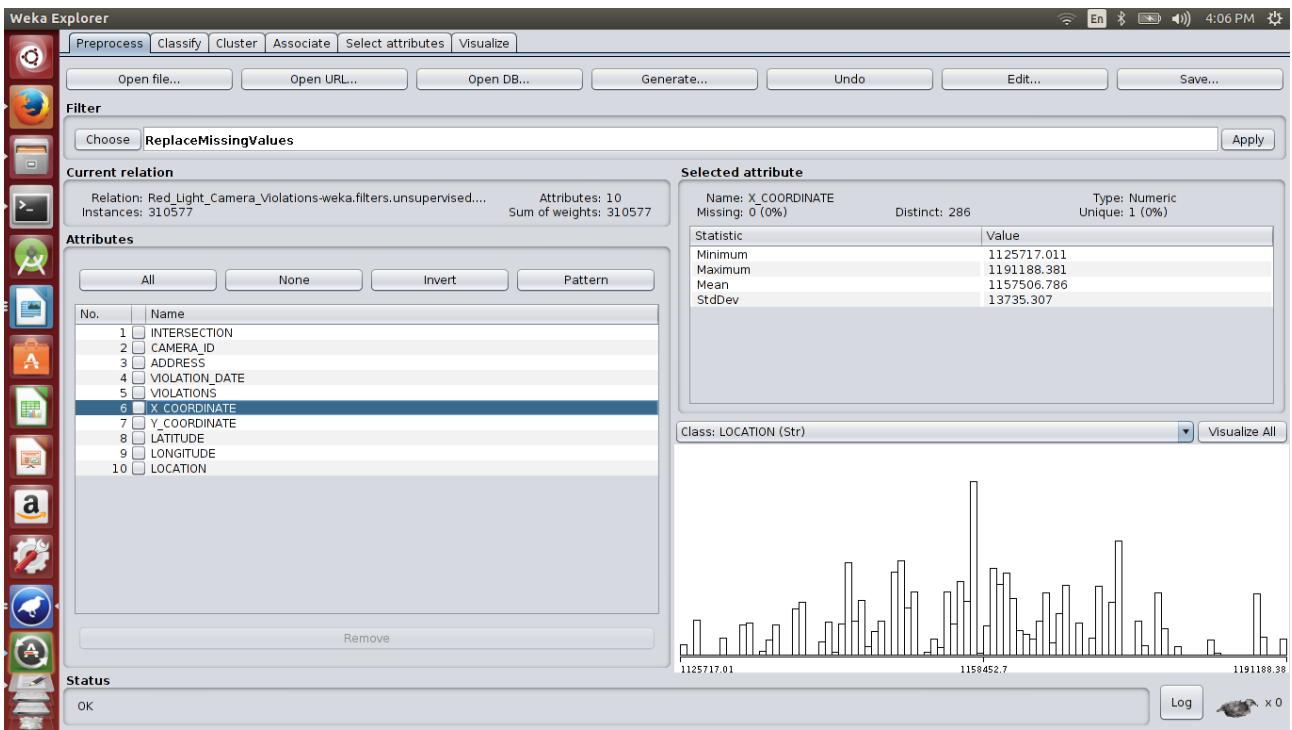
So, the missing values got replaced with mean.

Also we can choose the indices of the attributes of whose missing values we want to fill. Default it fills all the missing values of all attributes.



After clicking on Apply button:

We observe there are no more missing values.



Now the dataset looks like this:

All the missing values got replaced with mean of values.

Viewer

Relation: Red_Light_Camera_Violations-weka.filters.unsupervised.attribute.ReplaceMissingValues-weka.filters.unsupervised.attribute.StringToNominal.Rfirst-last-weka.filters.unsupervised.attribute.ReplaceMis...

No. 1: INTERSECTION 2: CAMERA_ID 3: ADDRESS 4: VIOLATION_DATE 5: VIOLATIONS 6: X_COORDINATE 7: Y_COORDINATE 8: LATITUDE 9: LONGITUDE 10: LOCATION

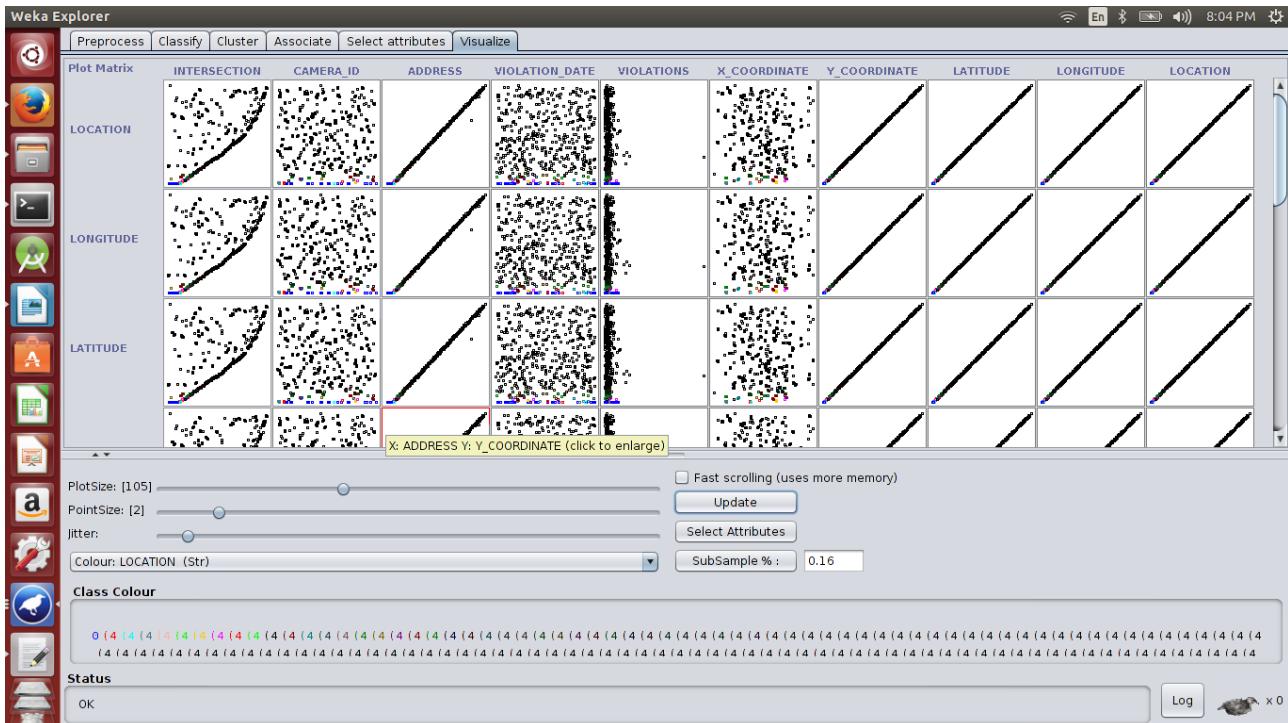
	INTERSECTION	CAMERA_ID	ADDRESS	VIOLATION_DATE	VIOLATIONS	X_COORDINATE	Y_COORDINATE	LATITUDE	LONGITUDE	LOCATION
1	WESTERNAND...	2672.0	67005W...	07/15/2016	9.0	1161488.747...	1860134.788...	41.7718...	-87.68358...	(41.771884...
2	PULASKIAND...	1542.0	2800NP...	07/15/2016	2.0	1149200.724...	1918332.222...	41.9318...	-87.72712...	(41.931831...
3	WESTERNAND...	1212.0	6400NW...	07/15/2016	4.0	1159113.585...	1942451.622...	41.9978...	-87.69003...	(41.997817...
4	111THANDHA...	2422.0	11100SH...	07/10/2016	14.0	1172923.657...	1831209.044...	41.6922...	-87.64252...	(41.692263...
5	STONYISLAND...	2731.0	67005CO...	07/07/2016	21.0	1188316.454...	1860851.137...	41.7732...	-87.5852246	(41.773251...
6	BELMONTAND...	1372.0	3200NKE...	07/09/2016	28.0	1154426.977...	1921110.466...	41.9393...	-87.70784...	(41.939351...
7	FOSTERANDN...	1031.0	5200NN...	07/13/2016	13.0	1132529.656...	1934201.702...	41.9756...	-87.78802...	(41.975685...
8	87THANDVINC...	2412.0	87005VL...	07/15/2016	11.0	1172026.328...	1847142.496...	41.7360...	-87.64534...	(41.736006...
9	DAMENAND63...	2701.0	63005DA...	07/13/2016	4.0	1164069.407...	1862830.277...	41.7792...	-87.67405...	(41.779226...
10	111THANDH...	1761.0	16000NH...	07/15/2016	2.0	1170678.283...	1910880.991...	41.9109...	-87.64841...	(41.910940...
11	HALSTEDAND...	1134.0	1000WF...	06/08/2017	2.0	1168689.247...	1934758.578...	41.9765...	-87.65503...	(41.976505...
12	SHERIDANAN...	1734.0	20000WDI...	06/18/2016	13.0	1162815.710...	1908034.672...	41.9032...	-87.67738...	(41.903298...
13	DIVISIONAND...	1814.0	3000WC...	07/13/2016	21.0	1156075.266...	1905216.604...	41.8957...	-87.70221...	(41.895704...
14	SACRAMENTO...	2212.0	2400WS...	06/08/2017	3.0	1161277.012...	1868185.749...	41.7939...	-87.68414...	(41.793981...
15	55THANDWES...	1514.0	2800WDI...	07/13/2016	45.0	1157212.542...	1918526.674...	41.9322...	-87.69767...	(41.932205...
16	CALIFORNIAA...	2422.0	11100SH...	06/09/2017	7.0	1172923.657...	1831209.044...	41.6922...	-87.64252...	(41.692263...
17	111THANDH...	2464.0	1600079...	07/14/2016	3.0	1188284.704...	1852964.829...	41.7516...	-87.58559...	(41.751611...
18	STONEYISLAN...	1134.0	1000WF...	07/13/2016	3.0	1168689.247...	1934758.578...	41.9765...	-87.65503...	(41.976505...
19	STATEAND79TH...	2651.0	79005ST...	07/26/2016	9.0	1177618.166...	1952621.870...	41.7509...	-87.62469...	(41.750918...
20	FOSTERANDB...	1111.0	5200NBR...	07/13/2016	5.0	1167334.508...	1934714.008...	41.9764...	-87.66001...	(41.976412...
21	111THANDH...	2572.0	71005CO...	07/15/2016	7.0	1182746.236...	1858088.011...	41.7658...	-87.60572...	(41.765800...
22	MILWAUKEE...	1364.0	5075WM...	06/30/2016	5.0	1141716.746...	1928736.076...	41.9605...	-87.75437...	(41.960522...
23	79THANDHAL...	2543.0	800W79...	07/15/2016	1.0	1172320.105...	1862517.981...	41.7507...	-87.64410...	(41.750751...
24	FULLERTONAN...	1553.0	6400WF...	07/14/2016	2.0	1133317.594...	1915293.661...	41.9237...	-87.78556...	(41.923786...
25	AUSTINANDA...	1491.0	3600NA...	07/15/2016	2.0	1135691.599...	1923335.967...	41.9458...	-87.77665...	(41.945813...
26	HALSTEDAND...	1714.0	800WDV...	07/14/2016	19.0	1170767.827...	1908259.644...	41.9037...	-87.64816...	(41.903745...
27	75THANDSTATE...	2621.0	75000ST...	07/13/2016	14.0	1177538.763...	1855235.551...	41.7580...	-87.62490...	(41.758092...
28	MONTROSEAN...	1221.0	4400NW...	07/15/2016	19.0	1159539.733...	1929183.918...	41.9614...	-87.68883...	(41.961401...
29	MILWAUKEEA...	1364.0	5075WM...	07/16/2016	1.0	1141716.746...	1928736.076...	41.9605...	-87.75437...	(41.960522...
30	CICEROANDAR...	2001.0	2000NCL...	07/01/2016	4.0	1144031.525...	1912881.830...	41.9169...	-87.74626...	(41.916973...
31	DAMENAND63...	2702.0	63005DA...	07/13/2016	4.0	1164069.407...	1862830.277...	41.7792...	-87.67405...	(41.779226...
32	WESTERNAND...	1562.0	3600NW...	07/13/2016	10.0	1159695.573...	1923871.086...	41.9468...	-87.68840...	(41.946819...
33	4700WESTERN...	2142.0	4700SW...	06/06/2017	3.0	1161120.438...	1873431.066...	41.8083...	-87.68457...	(41.808378...
34	CENTRALAND...	1604.0	5600WB...	07/14/2016	4.0	1138453.617...	1920755.972...	41.9386...	-87.76656...	(41.938683...
35	IRVINGPARKA...	1024.0	6400WIR...	07/13/2016	5.0	1132945.939...	1925946.304...	41.9530...	-87.78668...	(41.953024...
36	HAMILANDILA...	1881.0	3000NA...	07/05/2016	7.0	1150970.026...	1901381.688...	41.8852...	-87.72106...	(41.885282...
37	PULASKIAND6...	2342.0	6300SPU...	07/15/2016	13.0	11505746.003...	1862516.747...	41.7786...	-87.72290...	(41.778636...
38	KEDZIEAND26...	2382.0	26005KE...	07/15/2016	1.0	1155404.235...	1886524.460...	41.8444...	-87.70518...	(41.844424...
39	BELMONTAND...	1371.0	3200NKE...	07/10/2016	2.0	1154426.977...	1921110.466...	41.9393...	-87.70784...	(41.939351...
40	CORTLANDAN...	1721.0	1900NAS...	06/28/2016	38.0	1165348.084...	1912736.454...	41.9161...	-87.66794...	(41.916147...
41	JEFFERYAND9...	2522.0	9500SE...	06/07/2017	3.0	1191132.420...	1842335.735...	41.7223...	-87.57549...	(41.722375...

Add instance Undo OK Cancel

6)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

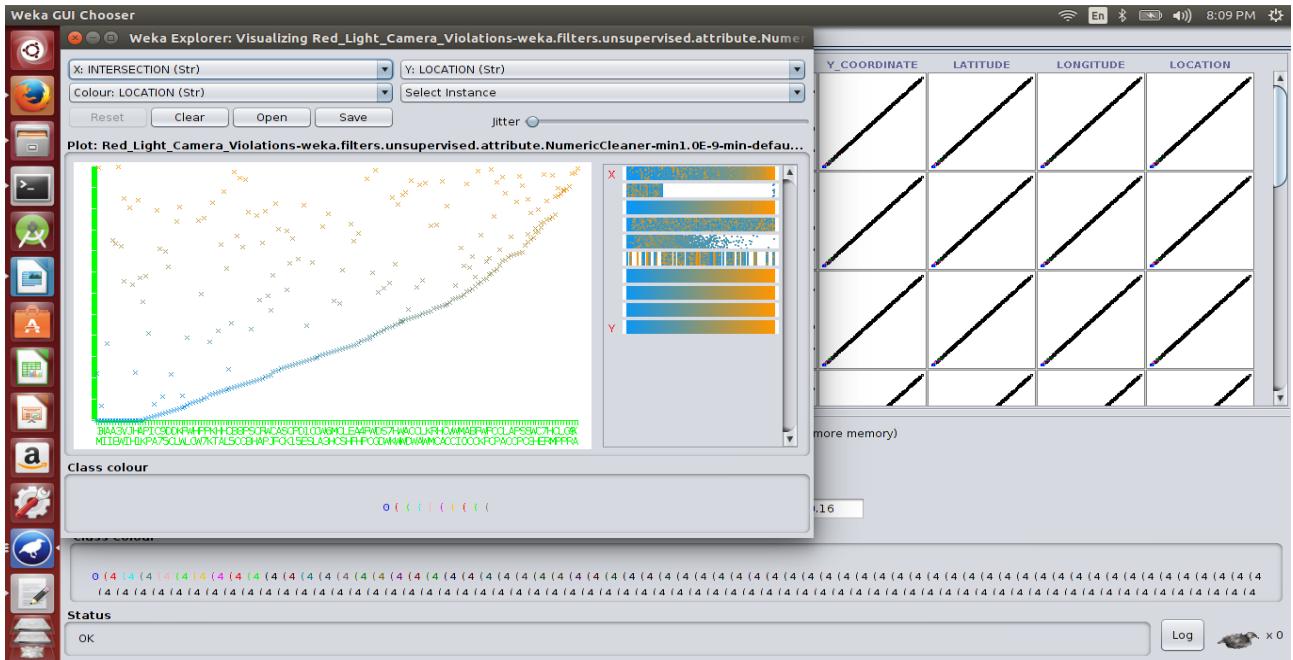
or negative correlations. So, Weka provides us Visualize tab for this purpose.



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.

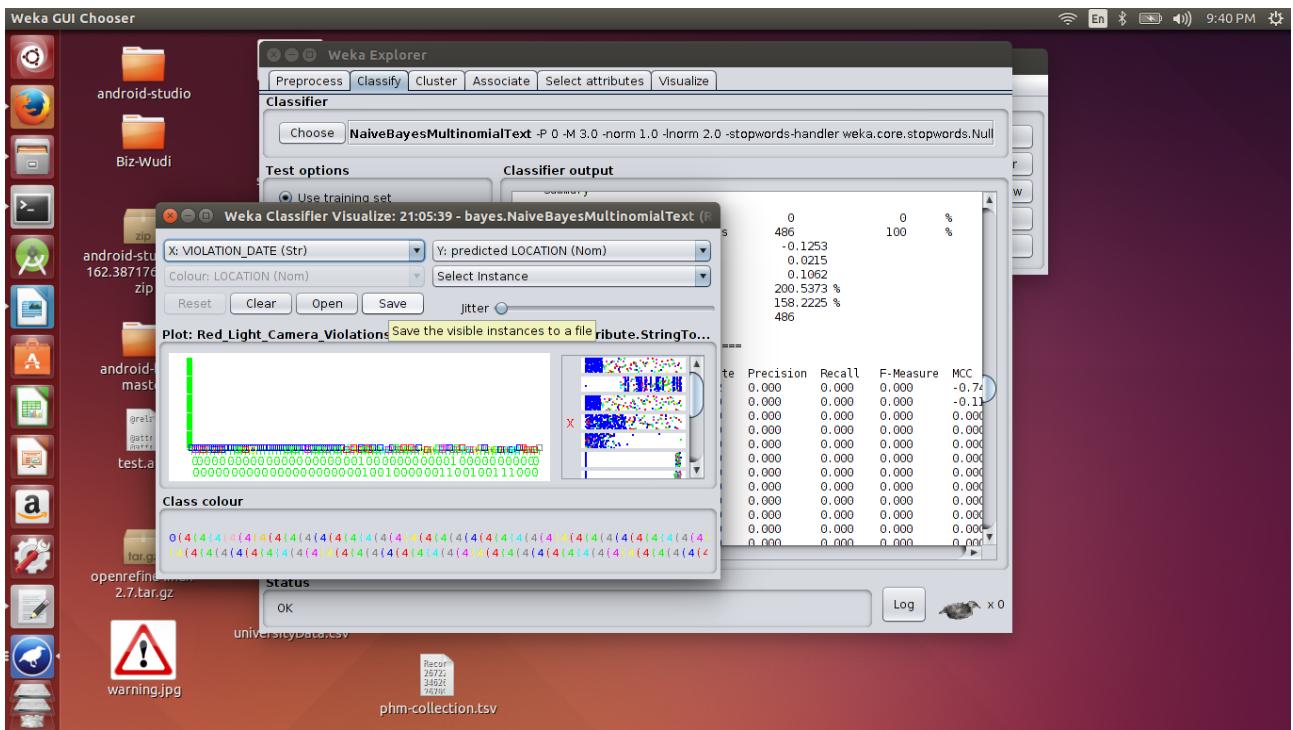
We can also see that each plot appears twice, first in the top left triangle and again in the bottom right triangle with the axes flipped. We can also see a series of plots starting in the bottom left and continuing to the top right where each attribute is plotted against itself.

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.

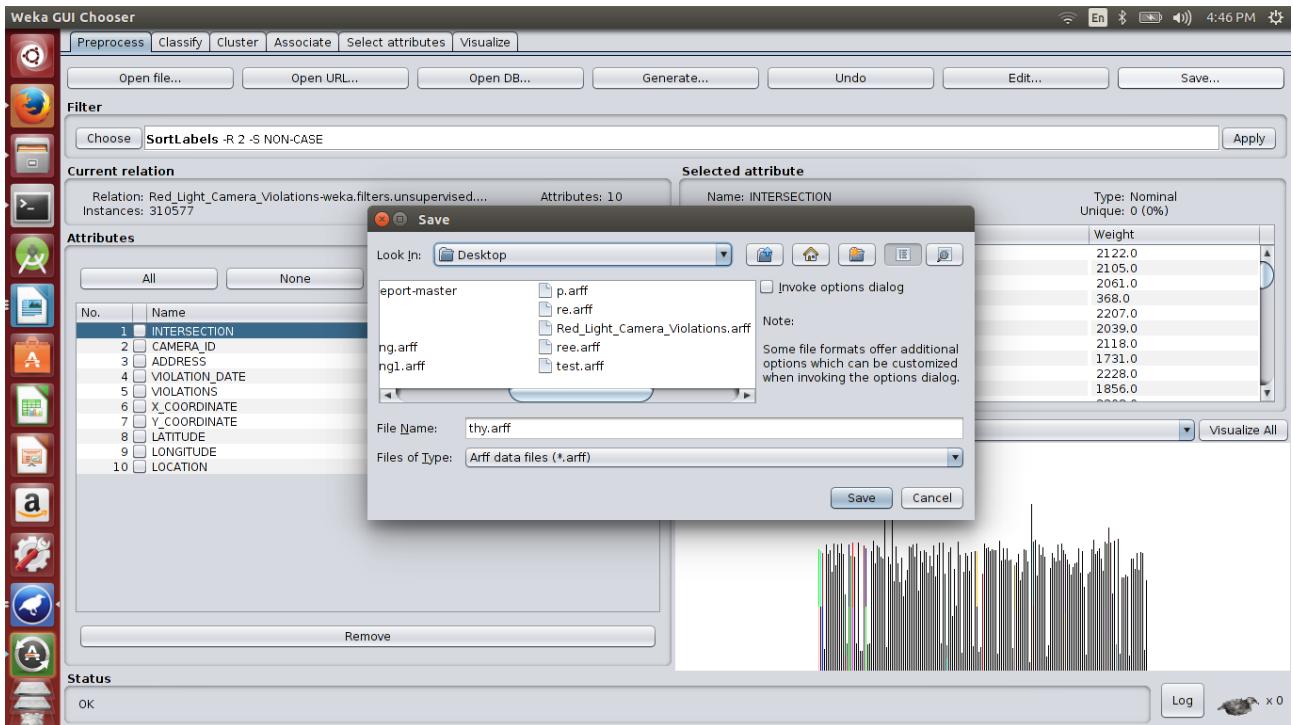


7)Explore Data:

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

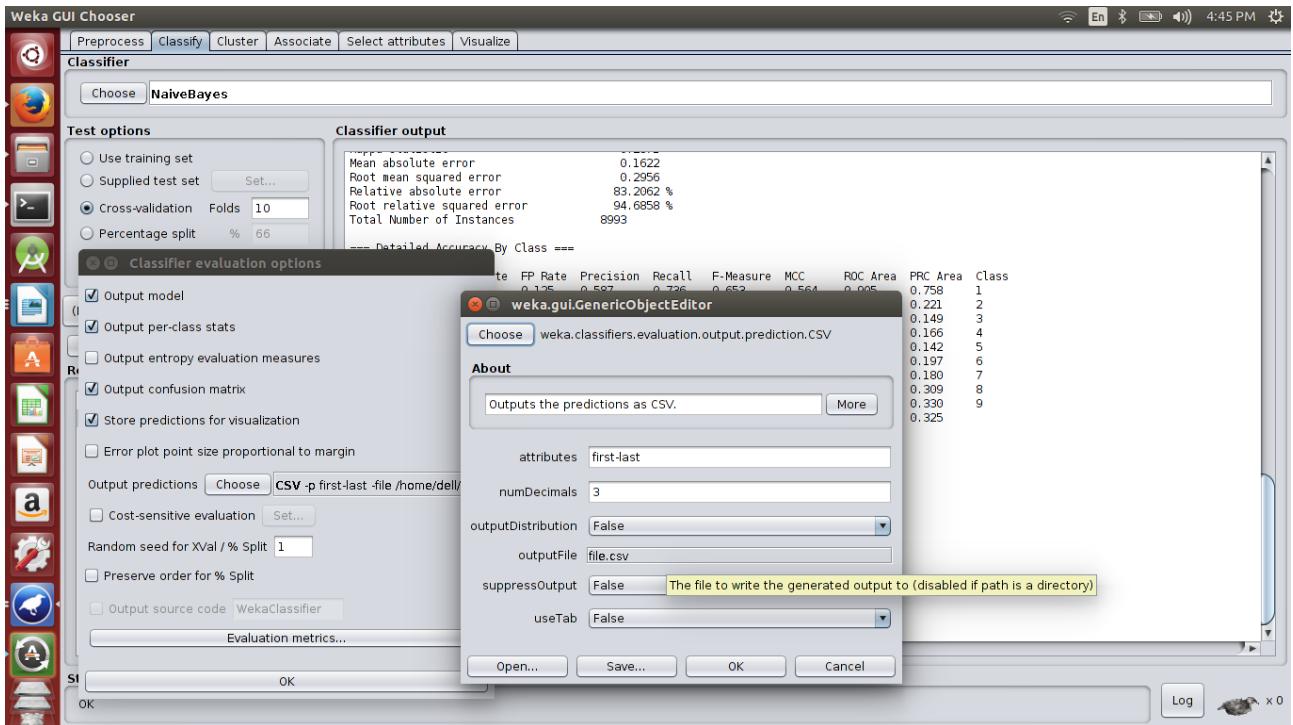


Or in the preprocess tab we have a save option where we can save the data.



And if we want data in csv format back We can do that by

selecting More options button from classify tab and then give the file format that we need and then the output file name.



Observed differences between weka and openrefine:

- 1) Weka provides many classifiers like naive bayes, Serialized classifier and many unlike open refine
- 2) Also for filling missing values weka has a direct option unlike openrefine.
- 3) Also we can change attributes from one type to other(nominal to numeric, String to nominal and many other) .
- 4) Scatterplot is obtained in both.
- 5) Both have options for clustering and merging.
- 6) So, whatever open refine have weka also has, but there are some which are not in openrefine but in weka.