

## DATA

We have the data of road accident of Seattle city, which consists of 1,94,673 cases and have 37 attributes or columns for all these accidents. In our data, the target label is “SEVERITYCODE” where the code for severity of accident is 1 for property damage collision and 2 for Injury collision. We have 37 columns in our data set where many of the columns in our data set are not fit for analysis like how many pedestrians, bicycles, vehicles or individuals involved in accident. There are various columns which are purely for administrative records. We exclude these unnecessary data and then make a bar plot of possible attributes in our data set to see which attributes gives us insight about the road accident severity. Data, in its original form, has many attributes which are unnecessary. First we remove the unnecessary columns.

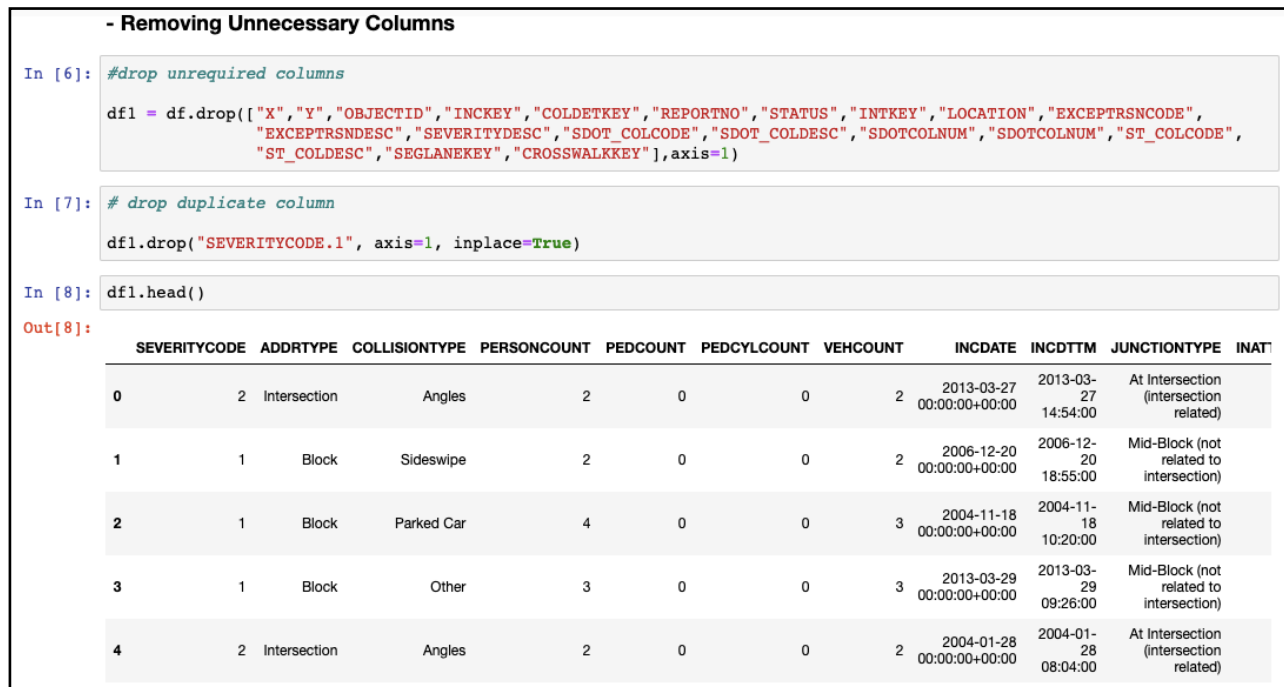


Figure 1

Original data had various missing values which must be replaced or removed. We either replaced them with mostly occurring variables of that column, or completed removed the corresponding rows from the data set.

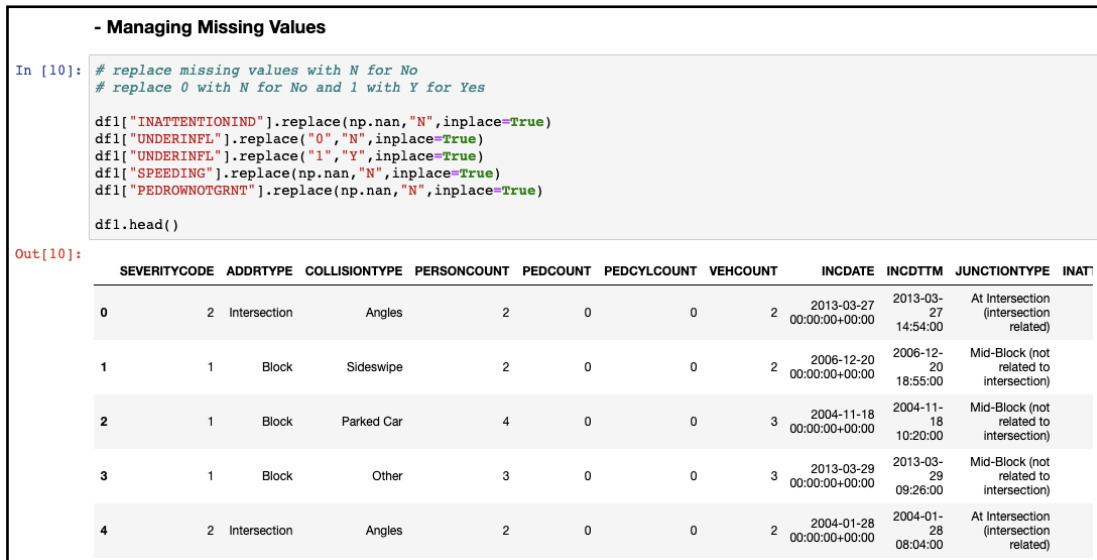


Figure 2

After plotting bar graphs we saw that data is highly unbalanced, where severity code in class 1 is almost three times the size of class 2. We balanced the data by down sampling it.

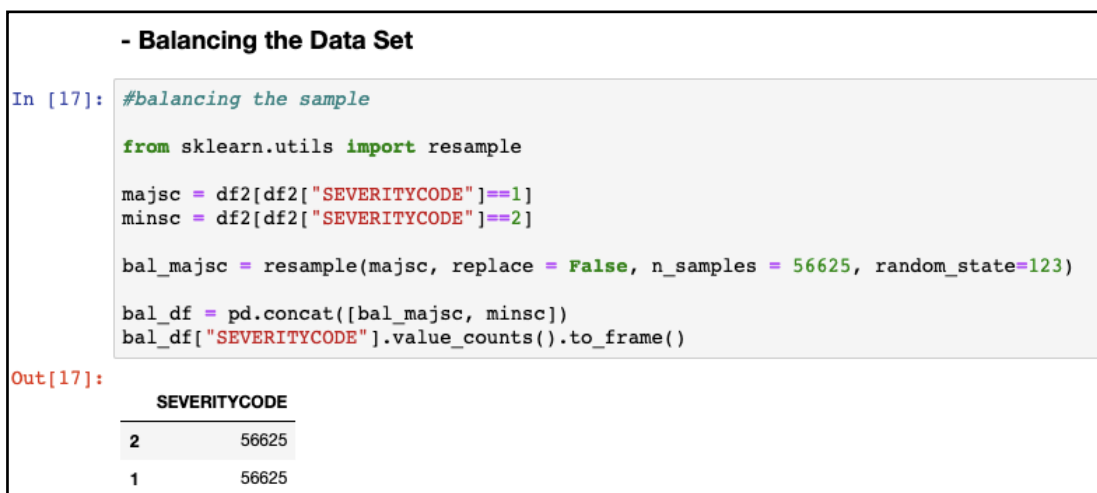


Figure 3

We made a bar plot of attributes like speeding, under influence of alcohol, and inattention of driver to study the accidents due to errors on part of individual.

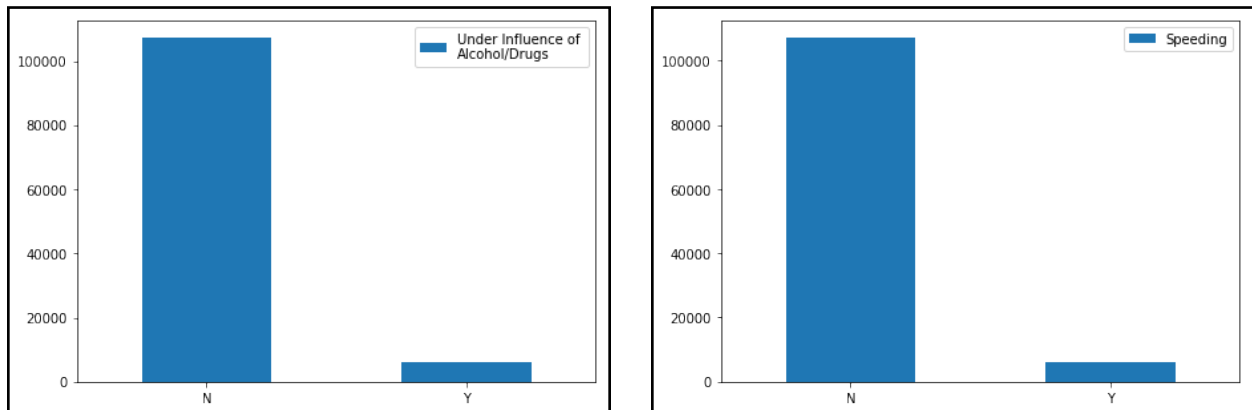


Figure 4

We also checked if the day of the week resembles with chances of accident and its severity. We noticed that working days have a higher frequency of collisions.

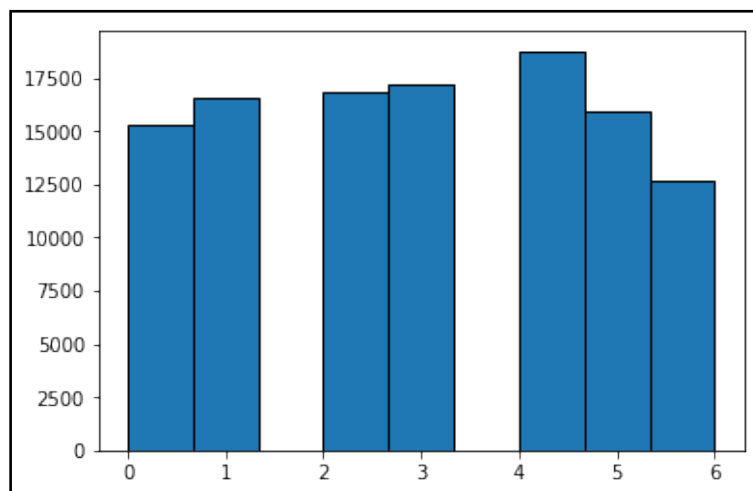
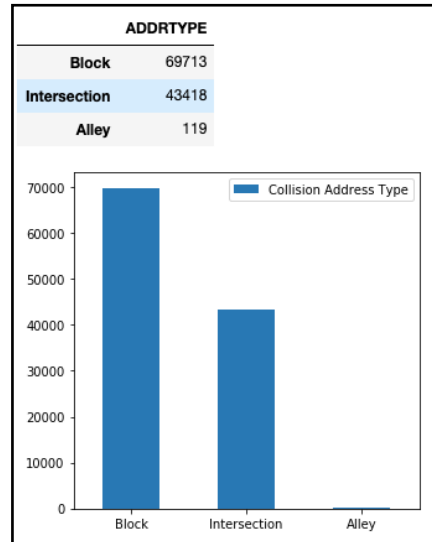


Figure 5

Attributes in our data set are mostly of string types so we used label encoding to make our data suitable for fitting with model. A sample of attribute “ADDRTYPE” is shown here.



```
In [29]: bal_df['ADDRTYPE'].replace(to_replace=['Block', 'Intersection', 'Alley'], value=[0,1,2],inplace=True)
bal_df.head()
```

Out[29]:

	SEVERITYCODE	ADDRTYPE	COLLISIONTYPE	PERSONCOUNT	PEDCOUNT	PEDCYLCOUNT	VEHCOUNT	INCDATE	INCDTTM
129091	1	0	Sideswipe	2	0	0	2	2014-05-17 00:00:00+00:00	2014-05-17 15:35:00
175353	1	1	Left Turn	5	0	0	2	2018-03-11 00:00:00+00:00	2018-03-11 15:04:00
110094	1	0	Parked Car	2	0	0	2	2012-08-26 00:00:00+00:00	2012-08-26 03:55:00
46167	1	1	Left Turn	5	0	0	2	2007-03-14 00:00:00+00:00	2007-03-14 18:00:00
38310	1	0	Parked Car	2	0	0	2	2006-10-19 00:00:00+00:00	2006-10-19 08:38:00

Figure 6

After all this preprocessing we finally select our attributes for model fitting which have significant impact on severity of accident and which can be used for model fitting. The attributes are:

1. **WEATHER**: It reveals about weather conditions of the day
2. **ROADCOND**: It gives information about road conditions.
3. **LIGHTCOND**: It gives information about the kind of light available during accident.
4. **ADDRTYPE**: It reveals the location or address of accident out of Alley, Block and Intersection.

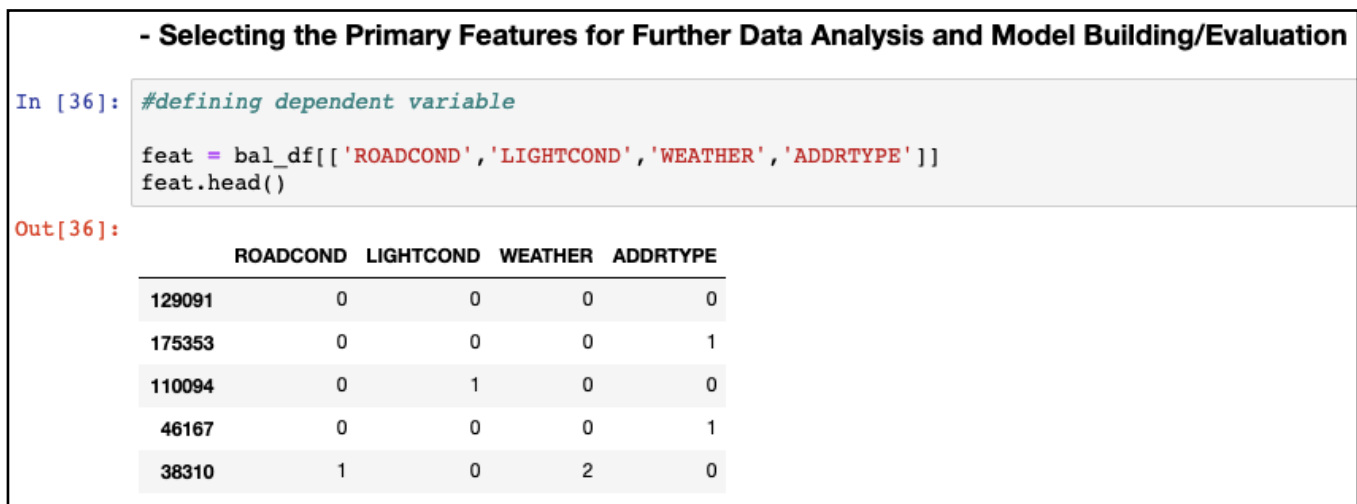


Figure 7

We got a perfectly balanced dataset with our chosen attributes to determine the target.