

## 2. Data Wrangling / Data Pre-processing:

### 2.1 Data Sources:

I am using the dataset provided by the Seattle Police Department (SPD) and recorded by the Seattle Department of Transportation (SDOT)(made available to me by Coursera). It contains all the collision records and was updated weekly. This dataset contains a total of 1,94,673 records and 37 attributes. To get detailed the detailed information and the metadata of this dataset, you can contact SDOT Traffic Management Division, Traffic Records Group.

### 2.2 Data Cleaning & Features Selection:

There were a lot of problems with the existing dataset. To start with, there were a lot of missing values and they were represented in 2 forms (i.e. NaN values, string values with 'Unknown' Label).

Example: ROADCOND attribute has 5012 missing values (NaN Type), but it also has 15078 values with 'Unknown' label.

```
In [9]: df.ROADCOND.value_counts().to_frame()
```

Out[9]:

ROADCOND	
Dry	124510
Wet	47474
Unknown	15078
Ice	1209
Snow/Slush	1004
Other	132
Standing Water	115
Sand/Mud/Dirt	75
Oil	64

```
In [10]: df.ROADCOND.isnull().sum()
```

Out[10]: 5012

Same was the case for all other attributes. So, the first step I took was to convert all the missing values into NaN type so that it will be easy for me while modelling the dataset. Also, I renamed 'X' and 'Y' attribute to 'LONGITUDE' and 'LATITUDE' respectively to make it more meaningful.

After fixing these problems, now it was time for Dimensionality Reduction. There were a lot of attributes which needs to be removed from data frame like,

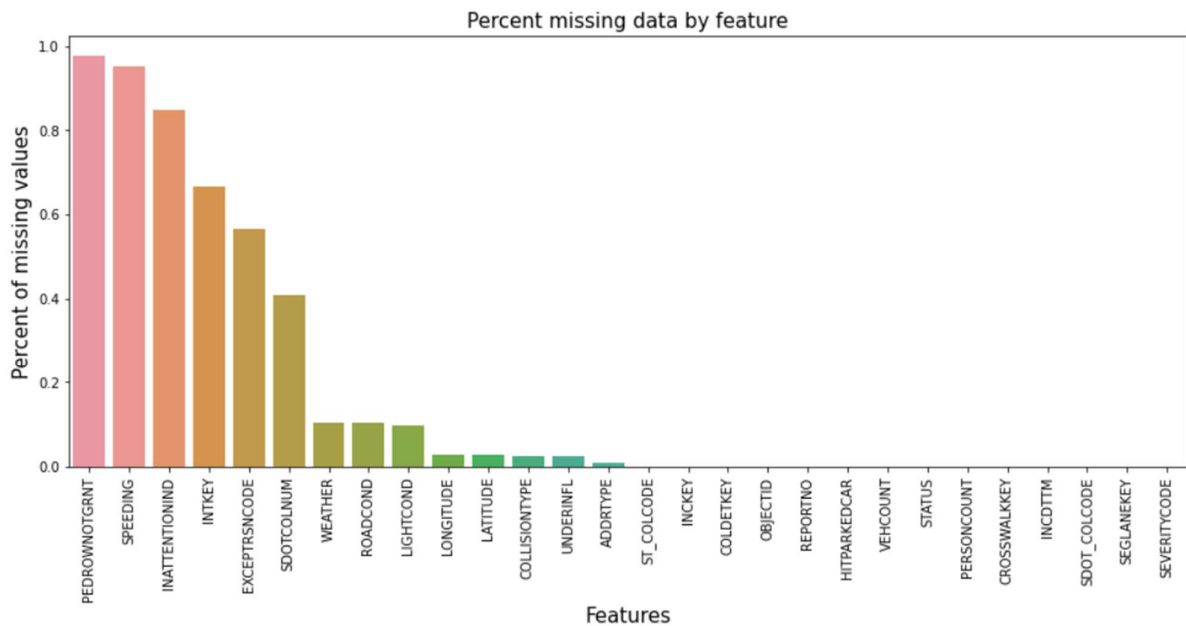
- Handling attributes which were giving the same information repeatedly.
- Handling attributes which have a significant amount of missing data.

Attribute like 'LOCATION' is redundant because this information is already available to us in the form of latitude and longitude. So 'LOCATION' needs to be dropped. Same was the case for the following attributes: SEVERITYCODE.1, SEVERITYDESC, PEDCOUNT, PEDCYLCOUNT, INCDATE, JUNCTIONTYPE, EXCEPTRSNDESC, SDOT\_COLDESC. Hence all these attributes have been dropped.

If an attribute has a lot of missing data, it can create a bias in the model. In this project, I have neglected all those attributes where missing value is more than 40%. Following are the attributes that I dropped:

- PEDROWNOTGRNT, missing value: 97.6%
- SPEEDING, missing value: 95.2%
- INATTENTIONIND, missing value: 84.7%
- INTKEY, missing value: 66.5%
- EXCEPTRSNCODE, missing value: 56.43%
- SDOTCOLNUM, missing value: 43.7%

PEDROWNOTGRNT	190006	0.976026
SPEEDING	185340	0.952058
INATTENTIONIND	164868	0.846897
INTKEY	129603	0.665747
EXCEPTRSNCODE	109862	0.564341



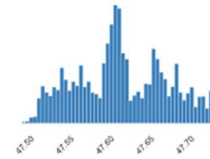
Now LATITUDE & LONGITUDE has 2.7 % missing value. Although I can use imputer to predict the missing values, it is very complicated. These missing values can't be imputed by simple SK Learn Imputer. This is a geographical data and needs a geography-specific library to predict these missing values. But, the missing values % is very small, I decided to drop those missing values.

### LATITUDE

Real number ( $\mathbb{R}_{\geq 0}$ )

MISSING

Distinct	23839	Mean	47.61954252
Distinct (%)	12.6%	Minimum	47.49557292
Missing	5334	Maximum	47.73414158
Missing (%)	2.7%	Zeros	0
Infinite	0	Zeros (%)	0.0%
Infinite (%)	0.0%	Memory size	1.5 MiB

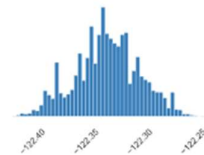


### LONGITUDE

Real number ( $\mathbb{R}$ )

MISSING

Distinct	23563	Mean	-122.3305184
Distinct (%)	12.4%	Minimum	-122.4190911
Missing	5334	Maximum	-122.2389494
Missing (%)	2.7%	Zeros	0
Infinite	0	Zeros (%)	0.0%
Infinite (%)	0.0%	Memory size	1.5 MiB



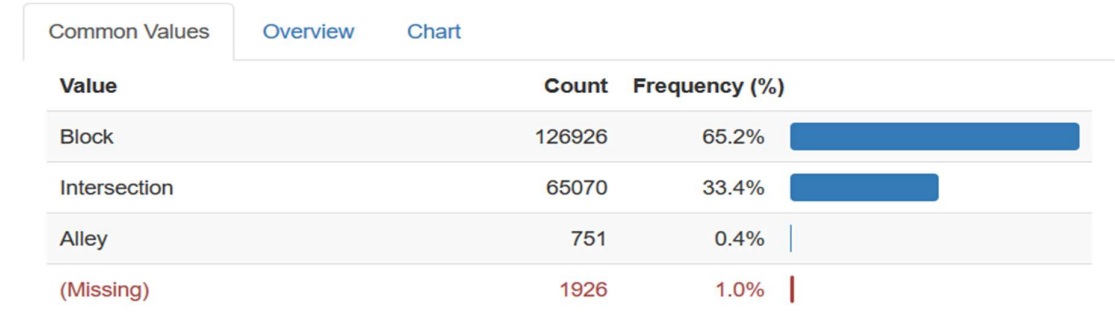
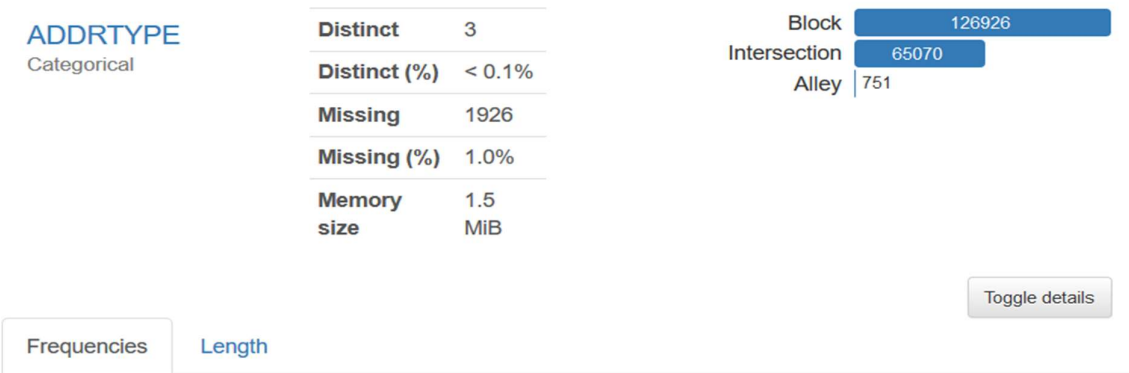
Now there are some attributes which have NO impact on our target variable (i.e. SEVERITYCODE). They are simply post-crash details recorded by the department. These attributes are not responsible to cause any crash and hence are not appropriate in predicting the severity of crashes. Such attributes are:

- OBJECT\_ID
- INCKEY
- COLDETKEY
- REPORTNO
- STATUS
- SDOT\_COLCODE
- ST\_COLCODE
- SEGLANEKEY
- CROSSWALKKEY

I decided to drop all these attributes.

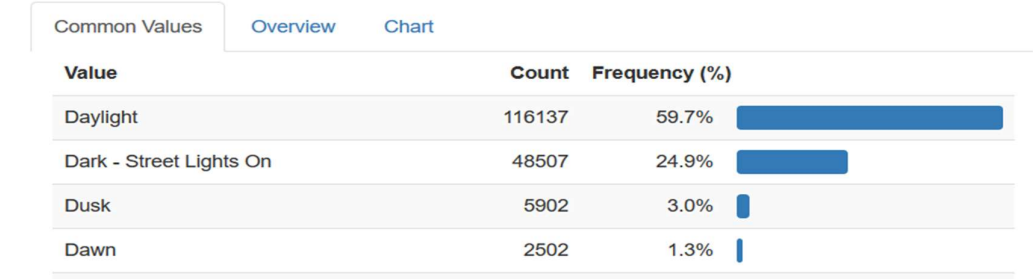
Handling Missing Values in remaining attributes

a. ADDRTYPE



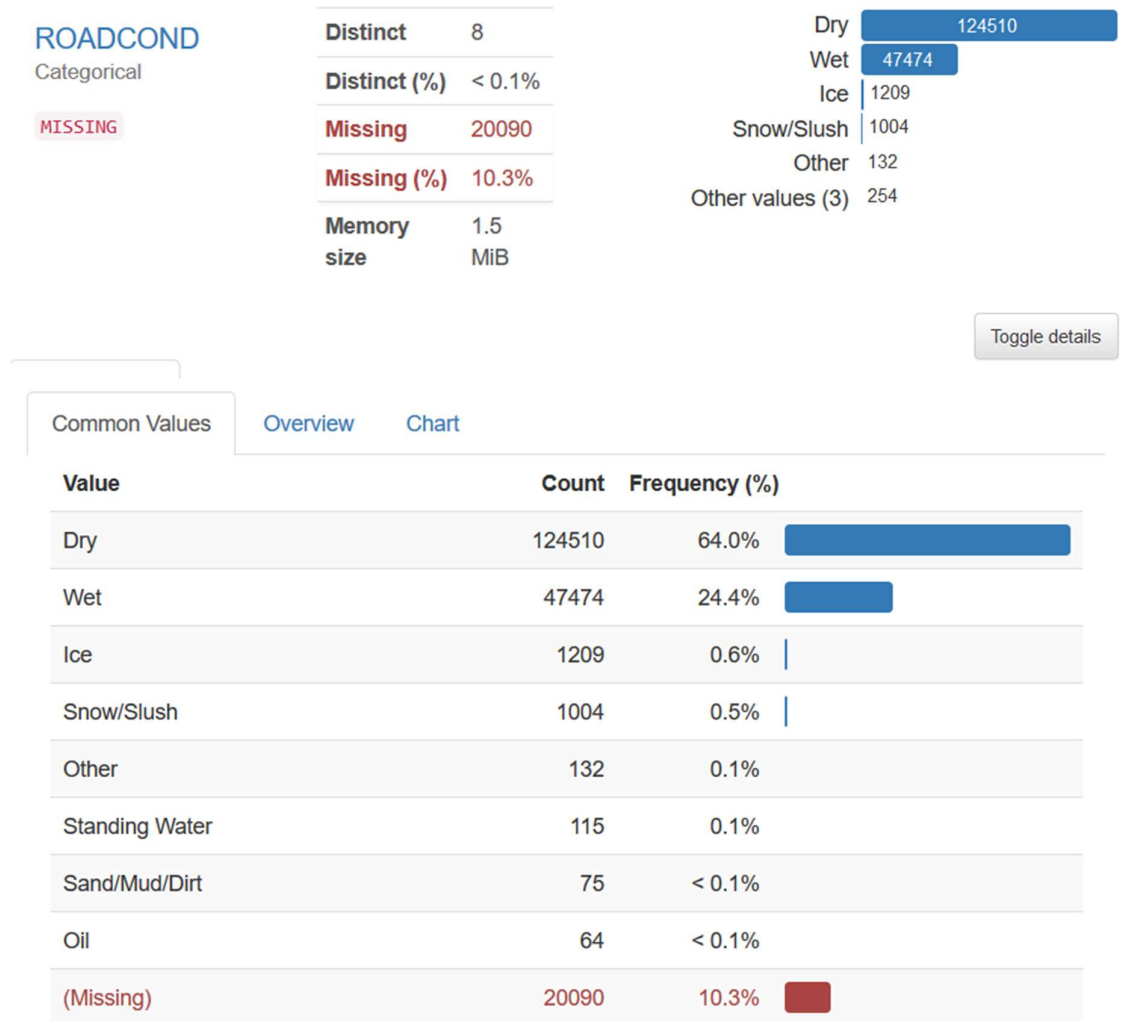
As 65% of crashes happened in Block, there is a high chance that the remaining 1% (missing values) also happened at Block. Therefore, I replaced NaN values by a label value 'Block'.

b. LIGHTCOND



Around 60% of accidents happened during daylight, the probability of those missing values being ‘Daylight’ is high. Hence, I replaced NaN values with the label “Daylight”.

### c. ROADCOND



In 64% of accidents, the condition of the road was ‘Dry’. Hence I decided to replace the NaN values with label “Dry”.

### d. WEATHER



Common Values				Overview	Chart
Value	Count	Frequency (%)			
Clear	111135	57.1%	<div></div>		
Raining	33145	17.0%	<div></div>		
Overcast	27714	14.2%	<div></div>		
Snowing	907	0.5%	<div></div>		
Other	832	0.4%	<div></div>		
Fog/Smog/Smoke	569	0.3%	<div></div>		
Sleet/Hail/Freezing Rain	113	0.1%			
Blowing Sand/Dirt	56	< 0.1%			
Severe Crosswind	25	< 0.1%			
Partly Cloudy	5	< 0.1%			
(Missing)	20172	10.4%	<div></div>		

Around 57% of the cases happened during a clear day. Hence, I replaced NaN values with label 'Clear'.