

DATA

Data understanding:

Data which I have used in this project I have collected it from example dataset which can be downloaded from [here](#). This dataset is in comma separated file format (csv). It has 38 columns and 194673 entries. Many of attributes are null values or not defined which must be removed or replaced with values like mean, standard deviation or median according to suitability of data cleaning. Among 38 columns all columns cannot be used for modelling. We need to exploratory analysis and some correlation analysis to select features. I have done some analysis which are shown in form of figure in next I will do data cleaning which will include null value handling, feature selection on basis of different analysis.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 194673 entries, 0 to 194672
Data columns (total 38 columns):
#   Column                Non-Null Count  Dtype
---  -
0   SEVERITYCODE           194673 non-null int64
1   X                      189339 non-null float64
2   Y                      189339 non-null float64
3   OBJECTID              194673 non-null int64
4   INCKEY                194673 non-null int64
5   COLDETKEY            194673 non-null int64
6   REPORTNO              194673 non-null object
7   STATUS                194673 non-null object
8   ADDRTYPE              192747 non-null object
9   INTKEY                65070 non-null float64
10  LOCATION              191996 non-null object
11  EXCEPTRSNCODE         84811 non-null object
12  EXCEPTRSNDESC         5638 non-null object
13  SEVERITYCODE.1        194673 non-null int64
14  SEVERITYDESC          194673 non-null object
```

Figure 1: Description of Dataset

```
Index(['SEVERITYCODE', 'X', 'Y', 'OBJECTID', 'INCKEY', 'COLDETKEY', 'REPORTNO',
      'STATUS', 'ADDRTYPE', 'INTKEY', 'LOCATION', 'EXCEPTRSNCODE',
      'EXCEPTRSNDESC', 'SEVERITYCODE.1', 'SEVERITYDESC', 'COLLISIONTYPE',
      'PERSONCOUNT', 'PEDCOUNT', 'PEDCYLCOUNT', 'VEHCOUNT', 'INCDATE',
      'INCDTTM', 'JUNCTIONTYPE', 'SDOT_COLCODE', 'SDOT_COLDDESC',
      'INATTENTIONIND', 'UNDERINFL', 'WEATHER', 'ROADCOND', 'LIGHTCOND',
      'PEDROWNOTGRNT', 'SDOTCOLNUM', 'SPEEDING', 'ST_COLCODE', 'ST_COLDDESC',
      'SEGLANEKEY', 'CROSSWALKKEY', 'HITPARKEDCAR'],
      dtype='object')
```

Figure 2: All columns of dataset

```
dtypes: float64(4), int64(12), object(22)
memory usage: 56.4+ MB
```

Figure 3: Data types of different columns