1 Data collection and understanding

1.1Data source

We used shared data for Seattle city as basis to deal with the accidents data. At first glance at the CSV file, we could see what type of data we have with us. The target label for the data set is Severity, which describes fatality of an accident. The remaining columns have different types of attributes. Also noticed that the data had some unbalanced attributes which need to be attended in next steps.

We also used the collisions meta data available for the years to understand the nature of all attributes. Having about 1.95L data observations, we could plan for split of the observations that could be used to train and test the prospective model.

1.2 Data understanding

The dataset basics are provided as follows;

Title: Collisions—All Years

Abstract: All collisions provided by Traffic Records.

Description: This includes all types of collisions. Collisions are displayed at the

intersection or mid-block of a segment in the Annexure.

Timeframe: 2004 to Present.

Keyword(s): SDOT, Seattle, Transportation, Accidents, Bicycle, Car, Collisions, Pedestrian,

Traffic, Vehicle

Types: The data is a mix of numerical and categorical types.

The data set provides labelled data for severity of accident. It shows a dual class categorical type of variable. The attributes (36 columns) convey information mainly about;

- the incident: such as identification no., location coordinates, date, time etc.
- the collision: such as code, type, description, injuries, fatalities etc.
- the impact: such as count of pedestrians, cyclists, vehicles involved etc.
- preconditions: such as inattention, influence of drugs, road condition, weather, speeding etc.

Attributes are almost complete with the information such as name, data type, length and description as shown in next section. State Collision Code Dictionary comprising about 85 codes with descriptions is also provided in supplement.

With the given dataset, severity code is identified as the target variable (labelled or dependent) while rest of the fields are noted as independent variables or the attributes. The case objective along with given data does qualify it as a classification problem of the supervised machine learning. All columns that could influence the cause and impact of an accident need to be selected for training and testing the model.

2 Data preparation

2.1Basic insight of dataset

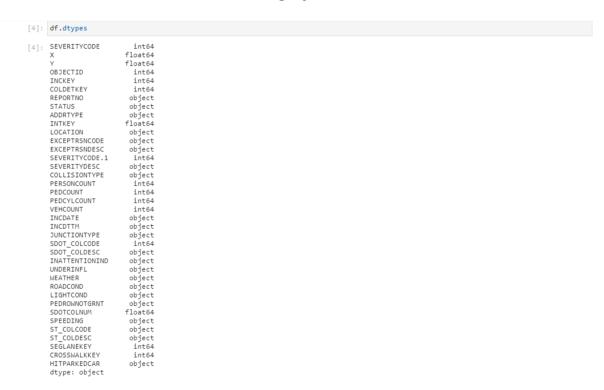
After reading data into Pandas data frame, it became a good start to explore the dataset. Following ways are followed to obtain essential insights of the data to help better understand the dataset:

Columns:

It provides list of columns that exist in the dataset.

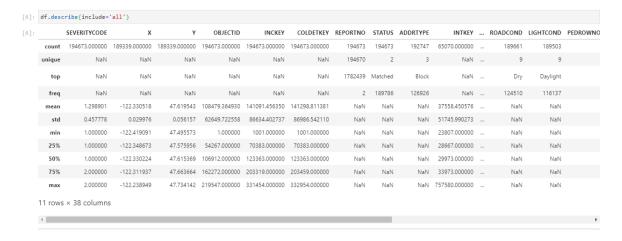
Data types:

This step is to know the variety of types viz. object, float, int, bool and datetime64. In order to better learn about each attribute, it is necessary to know the data type of each column, which were identified using Python as in screen shot below;



Data description:

We could get statistical summary, such as count, unique value, column mean value, column standard deviation, etc of each column. It provides various summary statistics, excluding NaN (Not a Number) values.



Info:

Similar to above, but it provides a concise summary of the DataFrame as shown below;

2.2Feature selection

In the first screening, it was noticed that some of the attributes were not significant to the cause of or to assess the impact of the severity. So these could be dropped for removing bias during design of the model. Rest of the columns were retained for further analysis.

2.3Cleansing data

In the second step, columns were analysed for missing values. Columns were treated for substituting missing values as shown in the table below;

Sr. No.	Attribute	Data type, length	Description	Wrangling Method	gRationale
1	OBJECTID	OBJECTID	ESRI unique identifier	Dropped	Insignificance
2	Χ	Longitude	ESRI geometry field	Dropped	Insignificance
3	Υ	Latitude	ESRI geometry field	Dropped	Insignificance
4	ADDRTYPE	Text, 12	Collision address type:	Retained	1% missing data,
			• Alley		replaced by max.
			• Block		frequency
			 Intersection 		
5	INTKEY	Double	Key that corresponds to		Insignificance
			intersection associated	with	
			a collision		
6	LOCATION	Text, 255	Description of the gene	raDropped	Insignificance
			location of the collision		
7	EXCEPTRSNCODE	Text, 10		Dropped	Insignificance
8	EXCEPTRSNDESC	Text, 300		Dropped	Insignificance
9	SEVERITYCODE	Text, 100	A code that corresponds		Target variable
			the severity of the		
			collision:		
			• 3—fatality		
			 2b—serious injury 		
			• 2—injury		
			• 1—prop damage		
			• 0—unknown		
10	SEVERITYDESC	Text	A detailed description o		Target variable
			the severity of the collis	ion	
11	COLLISIONTYPE	Text, 300	Collision type	Retained	2.5% missingdata,
		,	31		replaced by max.
					frequency
12	PERSONCOUNT	Double	The total number of peo	pRe tained	
			involved in the collision		
13	PEDCOUNT	Double	The number of pedestri		
			involved in the collision		
	DED 0) # 201 ··· :=		This is entered by the s		
14	PEDCYLCOUNT	Double	The number of bicycles		
			involved in the collision		
1 -	VELICOLINIT	Davidal -	This is entered by the s		
15	VEHCOUNT	Double	The number of vehicles		
			involved in the collision		
			This is entered by the s	ıace.	

Sr. No.	Attribute	Data type, length	Description	Wranglir Method	R ationale
16	INJURIES	Double	The number of total injustration involved in the collision. This is entered by the st		
17	SERIOUSINJURIES	Double	The number of serious injuries involved in the collision. This is entered the state.	Retained by	
18	FATALITIES	Double	The number of fatalities involved in the collision. This is entered by the st		
19	INCDATE	Date	The date of the incident	. Dropped	Insignificance
20	INCDTTM	Text, 30	The date and time of the incident.	eDropped	Insignificance
21	JUNCTIONTYPE	Text, 300	Category of junction at which collision took place	e	3.3% missingdata, replaced by max. frequency
22	SDOT_COLCODE	Text, 10	A code given to the coll by SDOT.	si lono pped	Insignificance
23	SDOT_COLDESC	Text, 300	A description of the collision corresponding the collision code.	to	
24	INATTENTIONIND	Text, 1	Whether or not collision was due to inattention (Y/N).	Dropped	85% data is missing
25	UNDERINFL	Text, 10	Whether or not a driver involved was under the influence of drugs or alcohol.	Dropped	only 3% observation are influencing
26	WEATHER	Text, 300	A description of the weather conditions duri the time of the collision	_	2.6% missingdata, replaced by max. frequency
27	ROADCOND	Text, 300	The condition of the roa during the collision.	dRetained	2.6% missingdata, replaced by max. frequency
28	LIGHTCOND	Text, 300	The light conditions dur the collision.	rRgetained	2.7% missingdata, replaced by max. frequency
29	PEDROWNOTGRNT	Text, 1	Whether or not the pedestrian right of way not granted. (Y/N)	Dropped was	97.6% data missing
30	SDOTCOLNUM	Text, 10	A number given to the collision by SDOT.	Dropped	Insignificance

Sr. No.	Attribute		Description	Wranglin Method	gRationale
31	SPEEDING	Text, 1	Whether or not speeding was a factor in the collision. (Y/N)		only 3% observations are influencing, rest data unavailable
32	ST_COLCODE	Text, 10	A code provided by the state that describes the collision. For more information about these codes, please see the State Collision Code Dictionar	tate	Insignificance
33	ST_COLDESC	Text, 300	A description that corresponds to the state coding designation.	Retained s's	2.5% missingdata, replaced by max. frequency
34	SEGLANEKEY	Long	A key for the lane segm in which the collision occurred.	e D itopped	Insignificance
35	CROSSWALKKEY	Long	A key for the crosswalk which the collision occurred.	alDropped	Insignificance
36	HITPARKEDCAR	Text, 1	Whether or not the collination involved hitting a parke car. (Y/N)		only 3.7% observations are influencing
37	STATUS	Text, 10	Matched, Unmatched	Dropped	Insignificance
38	REPORTNO	Long	Sr. No. of report for inte purposes	rDarbpped	Insignificance

2.4Transforming data

The last step in data cleansing would be to check and make sure that all data is in the correct format (int, float, text or other). To use categorical variables for regression analysis, indicator variables (or dummy variable) were used for transforming categorical variables into binary values (0s and 1s). This would make the data ready for next tests of correlation and determining significance. The results are as shown in screen shots below;

```
[15]: df_clean.info(max_cols=157)
                                                                                                <class 'pandas.core.frame.DataFrame'>
RangeIndex: 194673 entries, 0 to 194672
Data columns (total 156 columns):
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            Non-Null Count Dtype
                                                                                                           # Column
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    194673 non-null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 int64
                                                                                                                                                                   SEVERITYCODE
                                                                                                                                                                       SEVERITYDESC
                                                                                                                                                                       PERSONCOUNT
                                                                                                                                                                       PEDCOUNT
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 int64
                                                                                                                                                 PEDCYLCOUNT
VEHCOUNT
VEHCOUNT
VEHCOUNT
VEHCOUNT
ADDRTYPE_Alley
ADDRTYPE_Block
ADDRTYPE_Intersection
COLLISIONTYPE_Magles
COLLISIONTYPE_Cycles
COLLISIONTYPE_Cycles
COLLISIONTYPE_Head On
COLLISIONTYPE_Pedestrian
COLLISIONTYPE_Pedestrian
COLLISIONTYPE_Pedestrian
COLLISIONTYPE_Rape Ended
COLLISIONTYPE_Rape Ended
COLLISIONTYPE_Rape Linder
COLLISIONTYPE_Rape Linder
COLLISIONTYPE_Sideswipe
JUNCTIONTYPE_AI Intersection (intersection related)
JUNCTIONTYPE_AI Intersection (intersection related)
JUNCTIONTYPE_Mid-Block (but intersection related)
JUNCTIONTYPE_Mid-Block (but intersection related)
JUNCTIONTYPE_Name Junction
JUNCTIONTYPE_Name Junction
JUNCTIONTYPE_Name Junction
JUNCTIONTYPE_Name Junction
JUNCTIONTYPE_Unknown
SDOT_COLDESC_DRIVERLESS VEHICLE RAN OFF ROAD - HIT FIXED OBJECT
SOOT_COLDESC_DRIVERLESS VEHICLE RAN OFF ROAD - NO COLLISION
SDOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE FRONT END AT ANGLE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE LEFT SIDE AT ANGLE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE LEFT SIDE AT ANGLE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE LEFT SIDE SIDESWIPE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE LEFT SIDE SIDESWIPE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE LEFT SIDE SIDESWIPE
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE REAR END
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE REAR
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE REAR
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK MOTOR VEHICLE REAR
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK BODESTRIAN
SOOT_COLDESC_DRIVERLESS VEHICLE STRUCK BODESTR
                                                                                                                                                                   PEDCYLCOUNT
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 int64
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 int64
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     uint8
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     uint8
                                                                                                                                                    SDOT_COLDESC_MOTOR_VEHICLE_OVERTURNED IN ROAD

SDOT_COLDESC_MOTOR VEHICLE STRUCK PEDESTRIAN

SDOT_COLDESC_MOTOR VEHICLE STRUCK PEDESTRIAN

SDOT_COLDESC_MOTOR VEHICLE ANN OFF ROAD - HIT FIXED OBJECT

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, FRONT END AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, FRONT END AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, LEFT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, LEFT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, LEFT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, RIGHT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, RIGHT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, RIGHT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK MOTOR VEHICLE, RIGHT SIDE AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK PEDALCYCLIST, FRONT END AT ANGLE

SDOT_COLDESC_MOTOR VEHICLE STRUCK PEDALCYCLIST, LEFT SIDE SIDESWIPE

SDOT_COLDESC_MOTOR VEHICLE STRUCK PEDALCYCLIST, RIGHT SIDE SIDESWIPE

SDOT_COLDESC_PEDALCYCLIST OVERTURNED IN ROAD

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE FRONT END AT ANGLE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE FRONT END AT ANGLE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE REAR END

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE RIGHT SIDE SIDESWIPE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE RIGHT SIDE AT ANGLE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE RIGHT SIDE AT ANGLE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE RIGHT SIDE AT ANGLE

SDOT_COLDESC_PEDALCYCLIST STRUCK MOTOR VEHICLE RIGHT S
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              194673 non-null uint8
```

uint8

```
80 ROADCOMD_Stand/Mud/Dirt
10 ROADCOMD_Standing Water
11 ROADCOMD_Standing Water
12 ROADCOMD_Standing Water
13 ROADCOMD_Standing Water
13 ROADCOMD_Standing Water
13 ROADCOMD_Standing Water
14 ROADCOMD_Standing Water
15 ROADCOMD_Water
16 ROADCOMD_Water
17 ROADCOMD_Water
18 ROADCOMD_
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         uint8
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              uint8
114 ST_COLDESC_From same direction - one right turn - one straight
115 ST_COLDESC_From same direction - one right turn - one straight
115 ST_COLDESC_Not stated
116 ST_COLDESC_Not stated
117 ST_COLDESC_One car entering driveway access
118 ST_COLDESC_One car entering parked position
119 ST_COLDESC_One car leaving driveway access
120 ST_COLDESC_One car leaving parked position
121 ST_COLDESC_One car leaving parked position
121 ST_COLDESC_One car leaving parked position
122 ST_COLDESC_One car leaving parked position
123 ST_COLDESC_One car leaving parked position
124 ST_COLDESC_Dead_Cyclist Strikes Moving Vehicle
125 ST_COLDESC_Pedalcyclist Strikes Moving Vehicle
125 ST_COLDESC_Pedalcyclist Strikes Pedalcyclist or Pedestrian
126 ST_COLDESC_Pedalcyclist Strikes Pedalcyclist or Pedestrian
126 ST_COLDESC_Pedalcyclist Strikes Pedalcyclist
127 ST_COLDESC_Railway Vehicle Strikes Pedalcyclist
128 ST_COLDESC_Railway Vehicle Strikes Pedalcyclist
129 ST_COLDESC_Railway Vehicle Strikes Pedalcyclist
120 ST_COLDESC_Same direction -- both turning left -- one stopped -- rear end
120 ST_COLDESC_Same direction -- both turning left -- one stopped -- rear end
120 ST_COLDESC_Same direction -- both turning left -- one stopped -- rear end
120 ST_COLDESC_Same direction -- both turning left -- one stopped -- rear end
120 ST_COLDESC_Same direction -- both turning right -- both moving -- rear end
120 ST_COLDESC_Same direction -- both turning right -- both moving -- rear end
120 ST_COLDESC_Same direction -- both turning right -- both moving -- rear end
120 ST_COLDESC_Same direction -- both turning right -- one stopped -- rear end
120 ST_COLDESC_Same direction -- both turning right -- both moving -- rear end
120 ST_COLDESC_Same direction -- both turning right -- one stopped -- sideswipe
120 ST_COLDESC_Same direction -- both turning right -- one stopped -- sideswipe
120 ST_COLDESC_Same direction -- both turning right -- one stopped -- sideswipe
120 ST_COLDESC_Same direction -- both turning right -- one stopped -- rear end
120 ST_COLDESC_Same dir
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      194673 non-null uint8
```

2.5Test of correlation and significance

To get a better measure of the important characteristics, we would look at the correlation of the variables vis-a-vis target variable i.e. Accident Severity. 2 measures are followed for the analysis.

A. Pearson Correlation:

The Pearson Correlation measures the linear dependence between two variables X and Y of 'int64' or 'float64' types.

The resulting coefficient is a value between -1 and 1 inclusive, where:

- 1: Total positive linear correlation.
- 0: No linear correlation, the two variables most likely do not affect each other.
- -1: Total negative linear correlation.

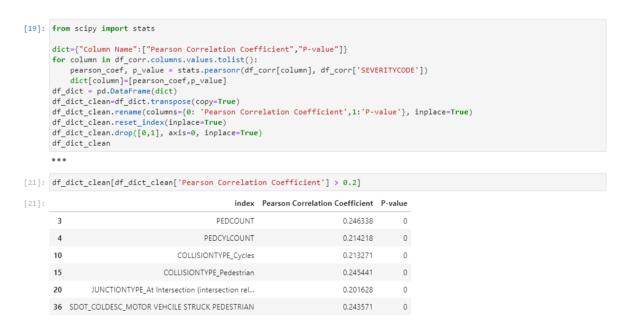
The closeness to terminal values (-1 and 1) would decide strength of the correlation.

B. P-value:

The P-value is the probability value that the correlation between these two variables is statistically significant. Normally, we choose a significance level of 0.05, which means that we are 95% confident that the correlation between the variables is significant. We would use "stats" module in the "Scipy" library to get the P-value.

By convention, when the

p-value is < 0.001: we say there is strong evidence that the correlation is significant. the p-value is < 0.05: there is moderate evidence that the correlation is significant. the p-value is < 0.1: there is weak evidence that the correlation is significant. the p-value is > 0.1: there is no evidence that the correlation is significant.



2.6Conclusion: Important Variables

By now we have a better idea of what our data looks like and which variables are important to take into account when predicting the 'Severity' class.

As we move into building machine learning modelling steps to automate our analysis, feeding the model with variables that meaningfully affect our target variable would help improve the model's prediction performance.