

## Introduction

In this IBM Data Science Professional Certificate at Coursera, the final Capstone Project consist of predicting the severity of a car crash using a real world dataset from the city of Seattle. There are plenty of business applications for these project from government using that insights collected to know the most dangerous places in the city to insurance companies that could use this to offer more customizable packages to their clients.

## Data Understanding

At Coursera we were given a dataset from Seattle Open Data and here are the attributes that were included in the dataset, also their data type, length and description. Attributes in **bold** were the ones that got selected as features for the project.

Attribute	Data type and length	Description
<b>SEVERITYCODE</b>	Text, 100	A code that corresponds to the severity of the collision: • 1—property damage • 2—injury
X	Not specified	Not specified
Y	Not specified	Not specified
OBJECTID	Not specified	Not specified
INCKEY	Long	Not specified
COLDEKEY	Long	Not specified
REPORTNO	Not specified	Not specified
STATUS	ObjectID	ESRI unique identifier
<b>ADDRTYPE</b>	Text, 12	Collision address type: • Alley • Block • Intersection
INTKEY	Double	Key that corresponds to the intersection associated with a collision
LOCATION	Text, 255	Description of the general location of the collision
EXCEPTRSNCODE	Text, 10	Not specified
EXCEPTRSNDESC	Text, 300	Not specified
SEVERITYDESC	Text	A detailed description of the severity of the collision
<b>COLLISIONTYPE</b>	Text, 300	Collision type
<b>PERSONCOUNT</b>	Double	The total number of people involved in the collision
PEDCOUNT	Double	The number of pedestrians involved in the collision. This is entered by the state.
<b>PEDCYLCOUNT</b>	Double	The number of bicycles involved in the collision. This is entered by the state.
<b>VEHCOUNT</b>	Double	The number of vehicles involved in the collision. This is entered by the state.
INCDATE	Date	The date of the incident.

<b>INCDTTM</b>	Text, 30	The date and time of the incident.
<b>JUNCTIONTYPE</b>	Text, 300	Category of junction at which collision took place
<b>SDOT_COLCODE</b>	Text, 10	A code given to the collision by SDOT.
<b>SDOT_COLDESC</b>	Text, 300	A description of the collision corresponding to the collision code.
<b>INATTENTIONIND</b>	Text, 1	Whether or not collision was due to inattention. (Y/N)
<b>UNDERINFL</b>	Text, 10	Whether or not a driver involved was under the influence of drugs or alcohol.
<b>WEATHER</b>	Text, 300	A description of the weather conditions during the time of the collision.
<b>ROADCOND</b>	Text, 300	The condition of the road during the collision.
<b>LIGHTCOND</b>	Text, 300	The light conditions during the collision.
<b>PEDROWNOTGRNT</b>	Text, 1	Whether or not the pedestrian right of way was not granted. (Y/N)
<b>SDOTCOLNUM</b>	Text, 10	A code provided by the state that describes the collision.
<b>SPEEDING</b>	Text, 1	Whether or not speeding was a factor in the collision. (Y/N)
<b>ST_COLCODE</b>	Text, 10	A code provided by the state that describes the collision.
<b>ST_COLDESC</b>	Text, 300	A description that corresponds to the state's coding designation.
<b>SEGLANEKEY</b>	Long	A key for the lane segment in which the collision occurred.
<b>CROSSWALKKEY</b>	Long	A key for the crosswalk at which the collision occurred.
<b>HITPARKEDCAR</b>	Text, 1	Whether or not the collision involved hitting a parked car. (Y/N)

Also, we can see in the following table the number of non-null values for every attribute.

<b>SEVERITYCODE</b>	194673 non-null int64
<b>X</b>	189339 non-null float64
<b>Y</b>	189339 non-null float64
<b>OBJECTID</b>	194673 non-null int64
<b>INCKEY</b>	194673 non-null int64
<b>COLDETKEY</b>	194673 non-null int64
<b>REPORTNO</b>	194673 non-null object
<b>STATUS</b>	194673 non-null object
<b>ADDRTYPE</b>	192747 non-null object
<b>INTKEY</b>	65070 non-null float64
<b>LOCATION</b>	191996 non-null object
<b>EXCEPTRSNCODE</b>	84811 non-null object
<b>EXCEPTRSNDESC</b>	5638 non-null object
<b>SEVERITYCODE.1</b>	194673 non-null int64
<b>SEVERITYDESC</b>	194673 non-null object
<b>COLLISIONTYPE</b>	189769 non-null object
<b>PERSONCOUNT</b>	194673 non-null int64
<b>PEDCOUNT</b>	194673 non-null int64
<b>PEDCYLCOUNT</b>	194673 non-null int64
<b>VEHCOUNT</b>	194673 non-null int64
<b>INCDATE</b>	194673 non-null object
<b>INCDTTM</b>	194673 non-null object
<b>JUNCTIONTYPE</b>	188344 non-null object
<b>SDOT_COLCODE</b>	194673 non-null int64
<b>SDOT_COLDESC</b>	194673 non-null object
<b>INATTENTIONIND</b>	29805 non-null object
<b>UNDERINFL</b>	189789 non-null object
<b>WEATHER</b>	189592 non-null object
<b>ROADCOND</b>	189661 non-null object
<b>LIGHTCOND</b>	189503 non-null object
<b>PEDROWNOTGRNT</b>	4667 non-null object
<b>SDOTCOLNUM</b>	114936 non-null float64
<b>SPEEDING</b>	9333 non-null object
<b>ST_COLCODE</b>	194655 non-null object
<b>ST_COLDESC</b>	189769 non-null object
<b>SEGLANEKEY</b>	194673 non-null int64
<b>CROSSWALKKEY</b>	194673 non-null int64
<b>HITPARKEDCAR</b>	194673 non-null object

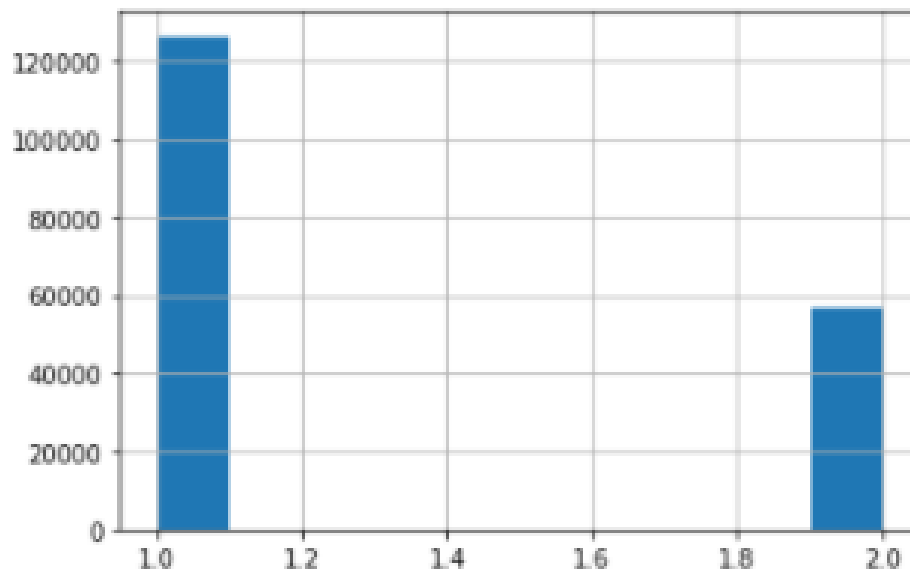
## Data Preparation

After features were selected we had to handle missing values. Here is our approach.

Feature	Non-null values	Null values approach
SEVERITYCODE	194673	No null values.
ADDRTYPE	192747	Since there were only 1926 null values, we decided to drop them.
COLLISIONTYPE	189769	Since there were only 4797 null values, we decided to drop them.
PEDCYLCOUNT	194673	No null values
VEHCOUNT	194673	No null values.
INATTENTIONIND	29805	Since we believe inattention would be evident, we made the assumption that null values were negative.
UNDERINFL	189789	We made an assumption that null values weren't under influence of substances.
WEATHER	189592	We made the assumption that null values were on clear days.
ROADCOND	189661	We made the assumption that null values were on clear roads.
LIGHTCOND	189503	We made the assumption that null values were on daylight.
PEDROWNOTGRNT	4667	We made the assumption that null values were cases in which pedestrian permission wasn't granted.
SPEEDING	9333	We made the assumption that null values weren't cases of speeding.

Assumptions made on "WEATHER", "ROADCOND" and "LIGHTCOND" were made because our reasoning was that people preferably decides to drive on the better conditions.

We made a plot in order to know how balanced was our target ("SEVERITYCODE") and we discovered that it was skewed and that we had an imbalanced dataset.



In order to balance our target outcomes, we decided to use the SMOTE (Synthetic Minority Oversampling Technique) algorithm in order to have the same number and that our models could work with balanced outcomes and give us better performance.

## Methodology

Since we want to predict if a car cash results in only a property damage or if it is more severe (injury) then we can say we have a classification problem. Classification models used for this project were the following:

- K Nearest Neighbors.
- Decision Tree.
- Support Vector Machine.
- Logistic Regression.

### K Nearest Neighbors

Is a simple algorithm that stores all available data points and classifies new cases bases on similarity measures, such as distance.

### Decision Tree

Is a model that separates possible outcomes into nodes and then it goes downstream as it breaks down into more nodes.

## Support Vector Machine

Is a model that is able to generalize between 2 different classes if the set of labeled data is provided in the training set to the algorithm,

## Logistic Regression

Is a model that uses a sigmoid function to obtain probabilities and therefore it can divide data points to differentiate classes.

## Results

Models were trained on 70% of the original dataset. With sklearn library it was possible to get accuracy scores that are determined on the test set to evaluate models. We will make our decision to select the best model with the highest F1 score. In the following table the results are presented.

<b>Model</b>	<b>Jaccard Index</b>	<b>F1 score</b>	<b>Log loss</b>
<i>KNN</i>	0.72	0.71	NA
<i>Decision Tree</i>	0.67	0.68	NA
<i>SVM</i>	0.66	0.68	NA
<i>Logistic Regression</i>	0.66	1	0.55

## Conclusion

In this project, we built classification models in order to get accurate results trying to predict whether a car crash could only end in a property damage or an injury. We have shown that we the features selected, severity code of a car crash could be predictable. Of all the four models selected for this project, Logistic Regression was the one with the better F1 score, alongside with that advantage that Logistic Regression has the fastest runtime.

## Future directions

Even though Logistic Regression was the model with better score, the F1 score of 1 could possibly be flawed and because of that suggestions are pretty welcomed.

Alongside the prediction of severity of a car crash another prediction that could be worth of working on could be an algorithm that returns the probability of the occurrence of a car crash according to determined features. This algorithm could be sold to car makers for them to implement it in their own cars and alert users when it predicts risk of a car crash.