# Study of Severity of US Car Accidents Using Machine Learning

**Shambhu KC**
**September 21, 2020**

# Existing Problem

❖ Road accidents kills 1.35 million/ year and injures 50 million/year

❖ Causes impact on the traffic

❖ Numbers can be minimized in many cases

# Project Goal

❖ Build a model by Machine Learning

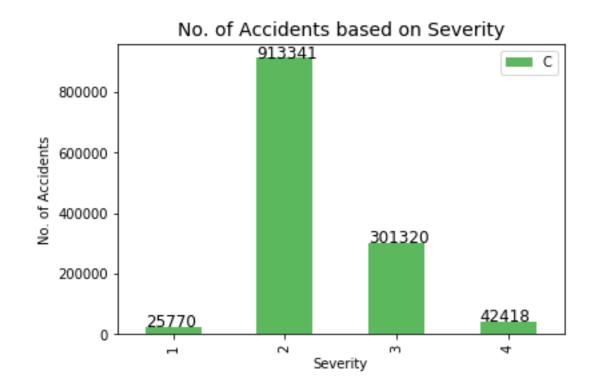❖ Suggest any recommendation to minimize the impact of accident

# Data Collection

❖ Historical data of accidents from Kaggle

❖ Record from Feb 2016 to June 2020 covering 49 states of US

❖ 3.5 million accident records

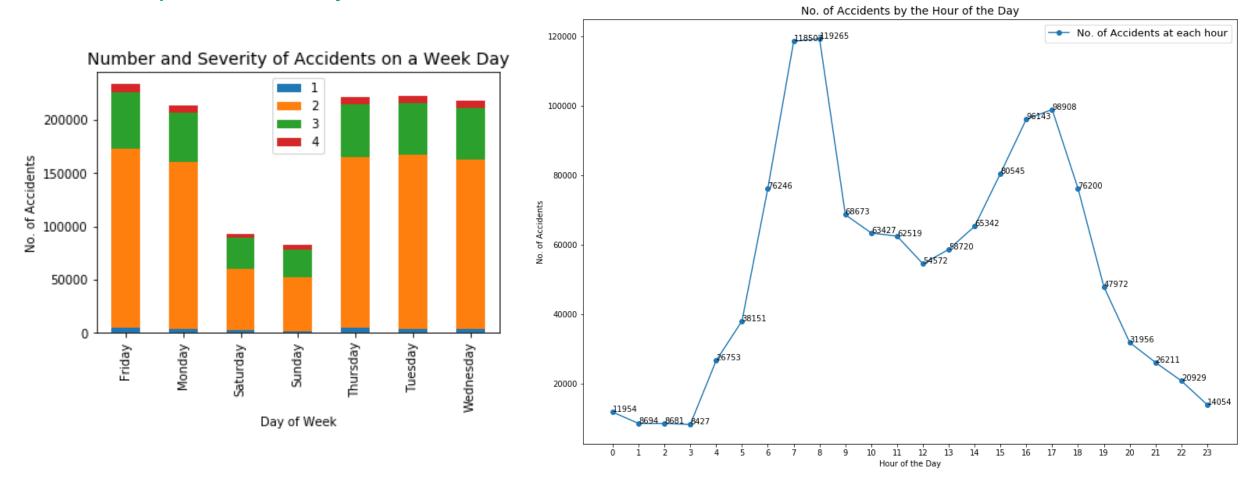❖ 46 columns covering many features; location, weather, traffic signals

# Data Cleaning

❖ Target variable is Severity in terms of impact on the traffic

❖ Features related to time, weather and traffic objects

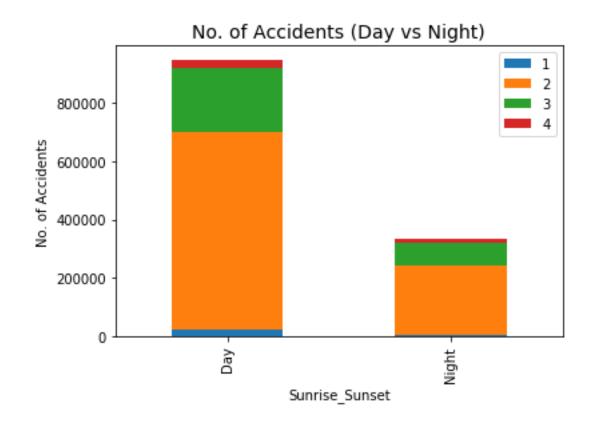❖ Removing columns with NaN values

❖ Removing unnecessary columns

# Data Visualization and Analysis

How many accidents of each category?



Majority of accidents were of severity level 2.

# Data Visualization and Analysis

How a particular day and time affects?



More accidents occurred during rush hours.

# Data Visualization and Analysis

What is the affect of day and Night?



More accidents occurred during day

# Data Visualization and Analysis

Correlation between different features and Severity

There is some degree of correlation.

# Data Visualization and Analysis

How does weather condition impact?



| | Precipitation(in) |
|---|---|
| **Severity** | |
| 1 | 0.005408 |
| 2 | 0.006966 |
| 3 | 0.011394 |
| 4 | 0.007558 |

# Data Visualization and Analysis

How does weather condition impact?



| Visibility(mi) | |
| --- | --- |
| **Severity** | |
| 1 | 9.500240 |
| 2 | 8.880002 |
| 3 | 8.728982 |
| 4 | 8.843497 |

| Wind_Chill(F) | |
| --- | --- |
| **Severity** | |
| 1 | 70.144137 |
| 2 | 59.569060 |
| 3 | 59.150609 |
| 4 | 58.066307 |

# Building a Model

K Nearest algorithm is used.

```python
# importing the necessary Library
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import f1_score, accuracy_score, f1_score, precision_score, recall_score, classification_report, confusion_r
from sklearn import metrics

# create an instance
knn = KNeighborsClassifier()

# we are going use GridSearchCV to find the best set of parameters
# define the hyperparameters

params = {'n_neighbors': [5,10],
          'weights': ['uniform','distance'],
          'algorithm':['auto', 'brute']}

grid_knn = GridSearchCV(estimator = knn, param_grid = params,
                        scoring = 'accuracy', cv =2,  # cv is number of cross-validation to try for each selected set of hyperpo
                        verbose =1, n_jobs = -1)

grid_knn.fit(X_train, y_train)

# extract best estimator
print("The best set of parameters ", grid_knn.best_estimator_)

# let's make prediction
yhat = grid_knn.predict(X_test)

# Let's print out evaluation metrices
print("Accuracy: ", accuracy_score(yhat, y_test))
print("F1_score:", f1_score(yhat, y_test))
print("Precision:", precision_score(yhat, y_test))
print("Recall:", recall_score(yhat, y_test))
```

Tried to optimize the parameters but the
computer couldn't perform the task.

# Result

❖ More accidents occurred during the day, typically at rush hour

❖ Outside temperature, precipitation, visibility, and wind chill impacts the seriousness of the accident.

❖ More accidents occurred in the proximity of traffic objects.

# Recommendation

❖ More stringent rule should be implemented at junction or at the traffic lights.

❖ The drivers should be extra cautious when the weather is not ideal.

❖ If possible, the rush should be avoided.

# Future Work

❖ Extend the model to cover other factors such as fatalities, and features like driver's sex, speed, vehicle type.