# PREDICTING CAR ACCIDENT SEVERITY IN SEATTLE, WASHINGTON

IBM APPLIED DATA SCIENCE CAPSTONE PROJECT

SARAH KRISTOLAITIS

# THE PROBLEM

- Car accidents are a common occurrence in the United States, and can cause significant distress, disability, decreased quality of life, financial hardship, and even death

- According to the National Safety Council, 38,800 people died in car crashed in the US in 2019, and 4.4 million people had to seek medical attention for car accident related injuries

# INTEREST IN THE PROBLEM

## Government:

- Reducing car accidents would reduce healthcare costs and negative impacts on worker productivity, GDP, etc.

- The data from this project could drive more targeted improvement to infrastructure, policing strategies, or public education efforts

## Private Businesses:

- May be able to use the data to create or improve products, such as GPS systems, that would assist their customers with reducing their accident risk

# DATA ACQUISITION AND CLEANING

- Data was obtained from the Applied Data Science Capstone project site, with additional details on the data taken from the Seattle GeoData website

- The selected attributes were those deemed to be most likely to provide early warnings of accident risk (i.e. would be able to be implemented into a GPS type of system in order to provide useful warnings to drivers)

- The target variable was the accident severity, rated at 1 (property damage only) or 2 (injury to a person occurred)
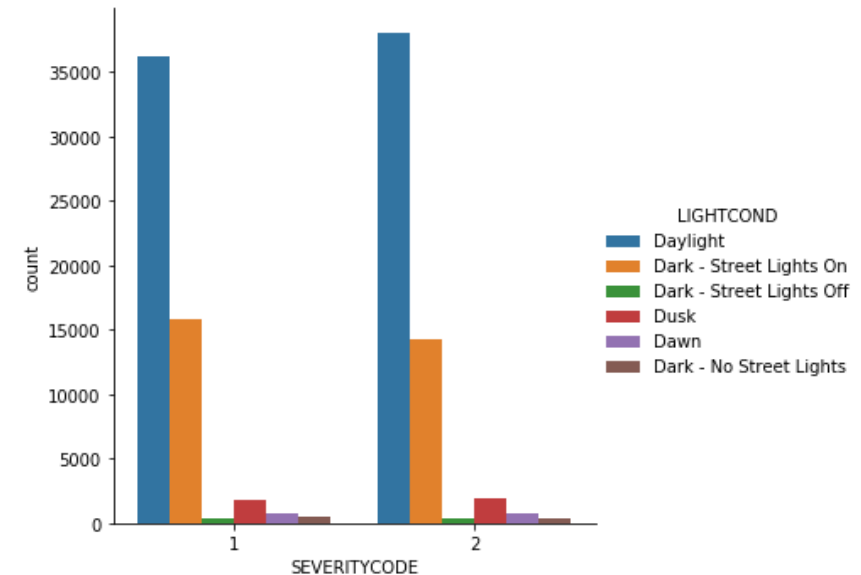
# DATA ACQUISITION AND CLEANING

- The 7 chosen attributes for analysis were Road Conditions, Light Conditions, Inattentiveness, Being Under an Influence, Speeding, Weather, and Address Type

- Data was balanced so that type 1 accidents did not overpower type 2

- Null values for Inattentiveness, Speeding, and Being Under an Influence were determined to be "N"

- Null values for all other fields resulted in the observation being eliminated

- There were 111,308 observations remaining for analysis and model development

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND LIGHT CONDITIONS

Accidents were most likely to result in personal injury:

- At dawn
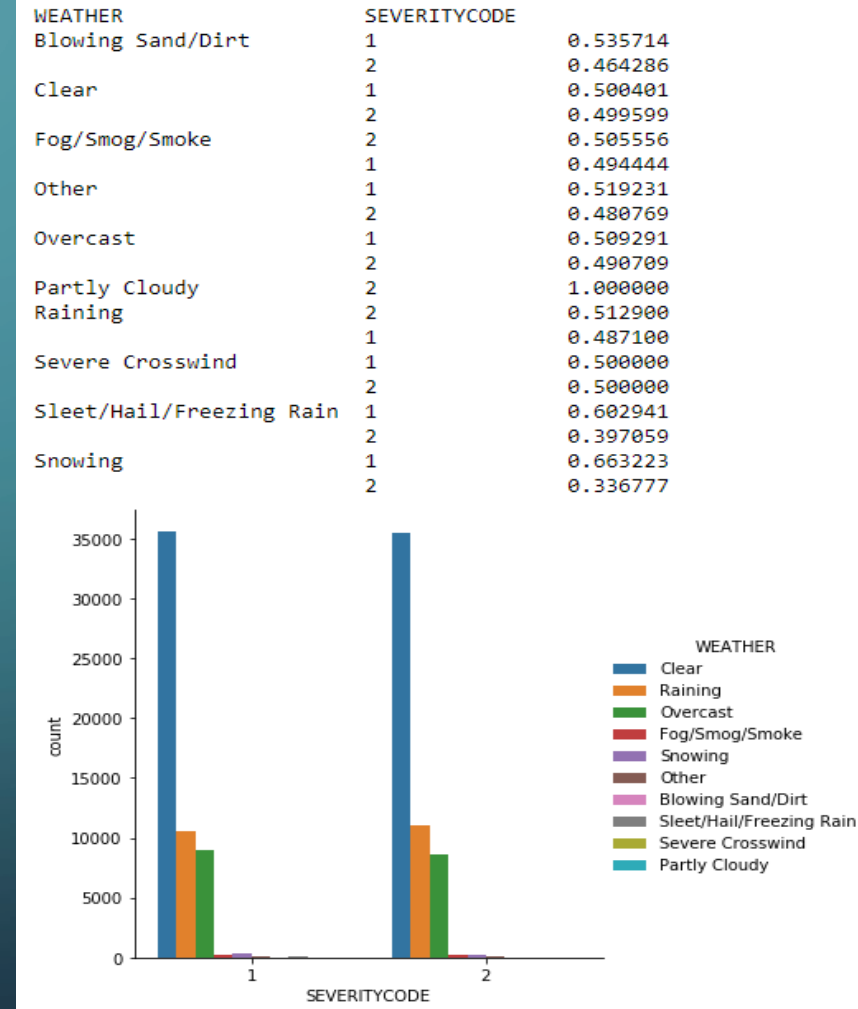- At dusk
- During daylight hours

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND WEATHER

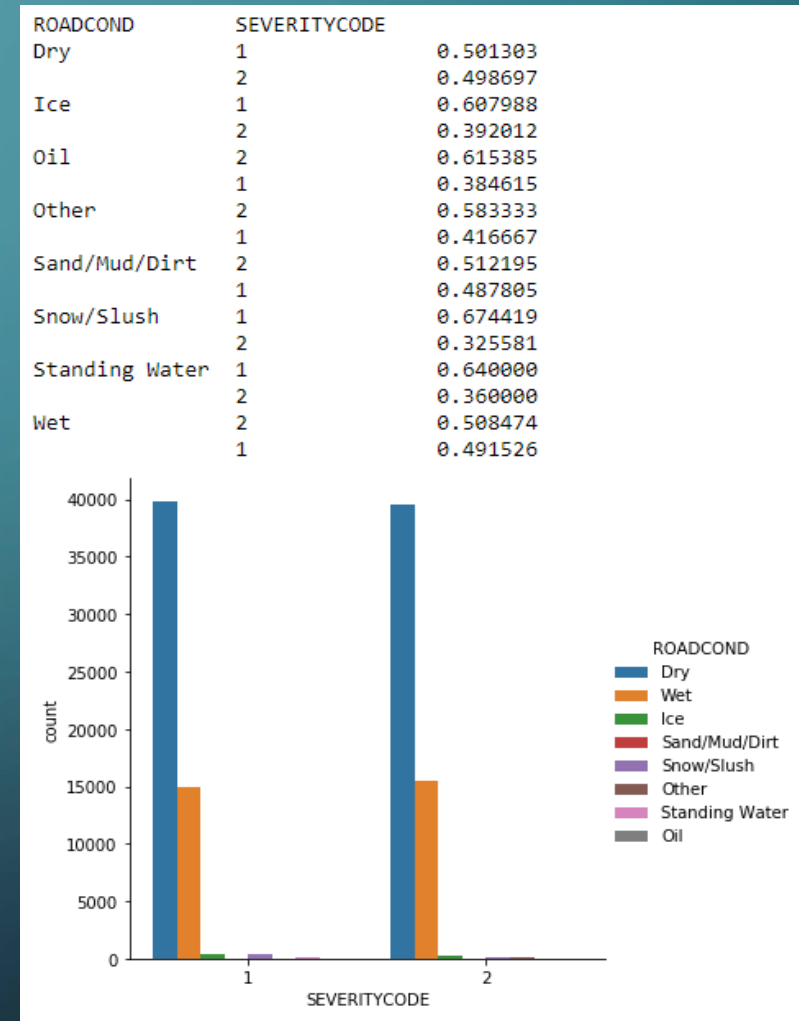Accidents were most likely to result in personal injury:

- With Fog/Smog/Smoke
- On Partly Cloudy days
- When Raining

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND ROAD CONDITIONS
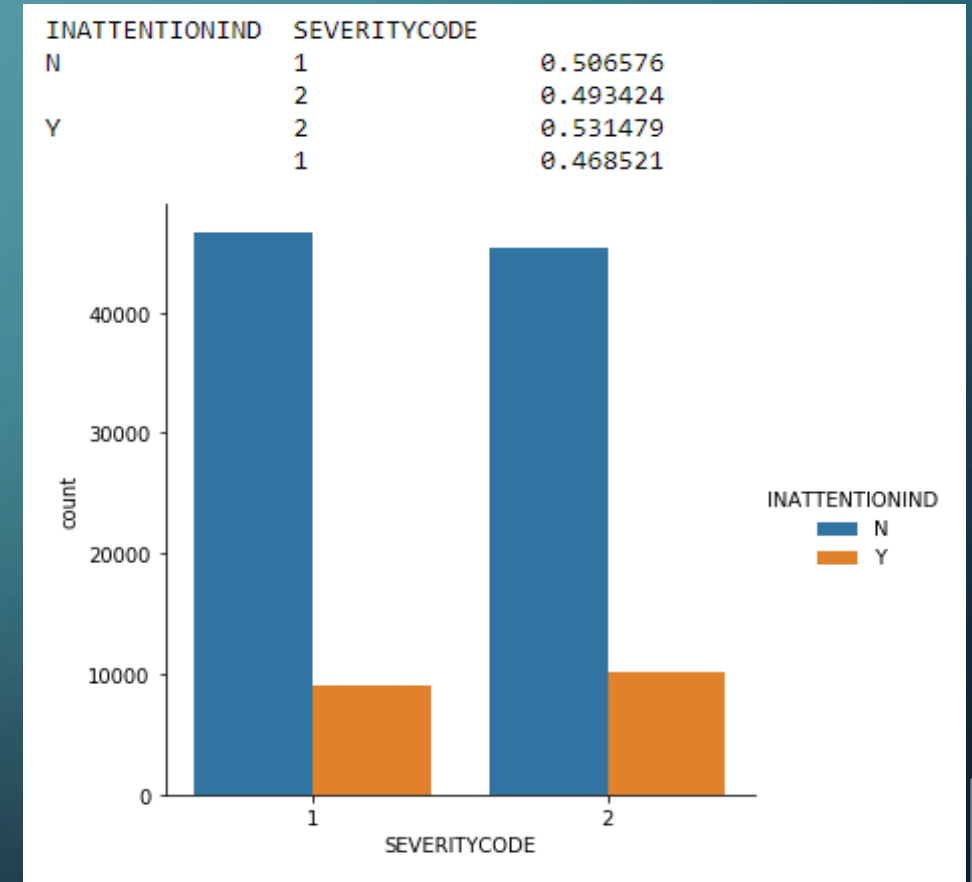
Accidents were most likely to result in personal injury:

- In Oil
- In Sand/Mud/Dirt
- On Wet Roads
- In "other" conditions

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND DRIVER INATTENTION

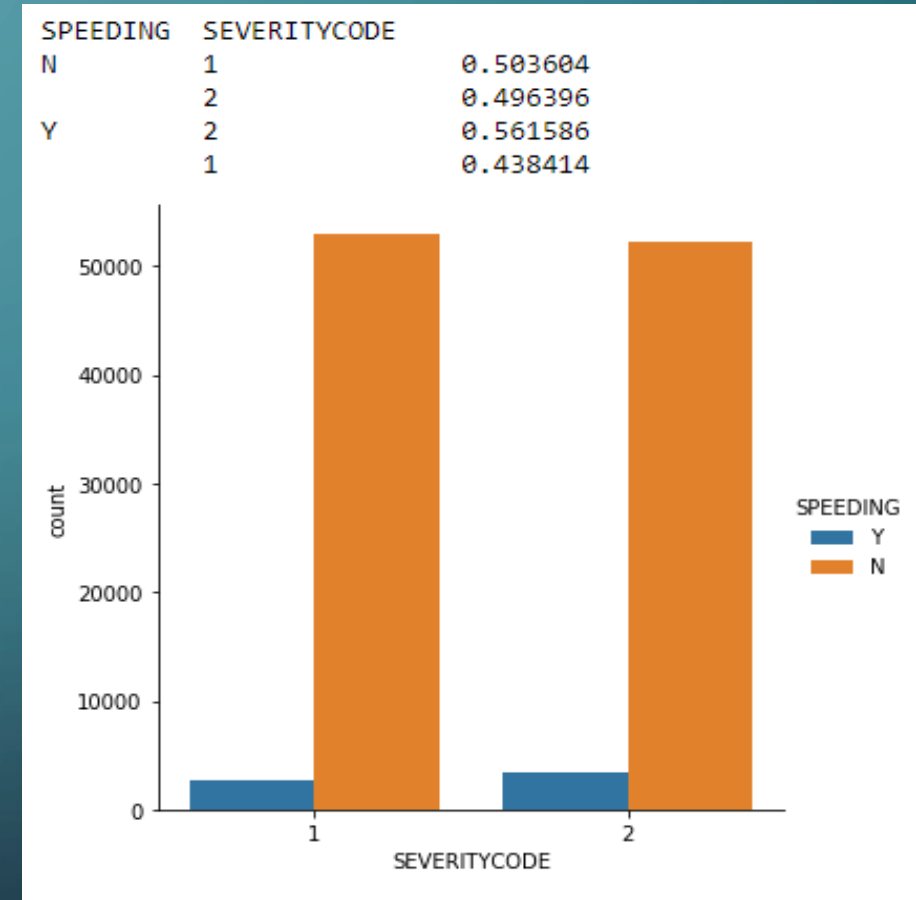Accidents were most likely to result in personal injury:

- When a driver is deemed to have been distracted

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND SPEEDING

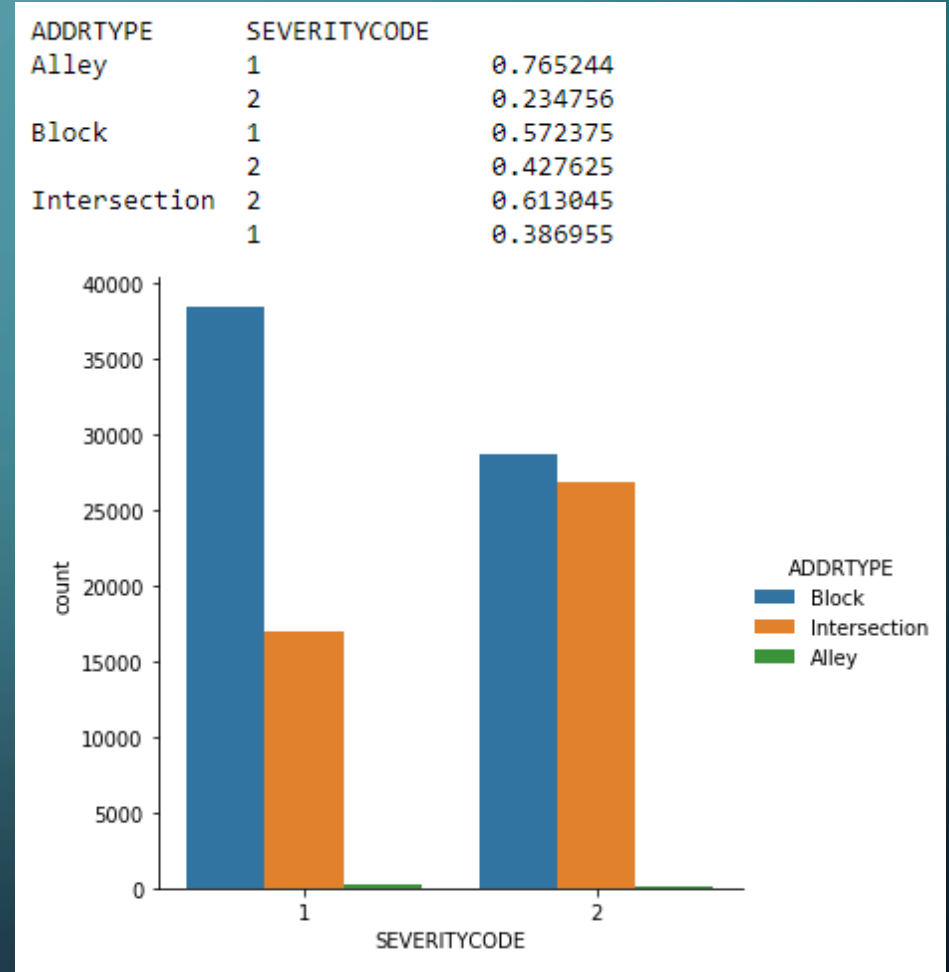Accidents were most likely to result in personal injury:

- When a driver is speeding

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND ADDRESS TYPE

Accidents were most likely to result in personal injury:

- In an intersection

# RELATIONSHIP BETWEEN ACCIDENT SEVERITY AND BEING UNDER AN INFLUENCE

Accidents were most likely to result in personal injury:

- When a driver is under the influence of drugs or alcohol

# PREDICTIVE MODELING

Four different machine learning models were developed and tested:

- K Nearest Neighbour

- Decision Tree

- Support Vector Machine

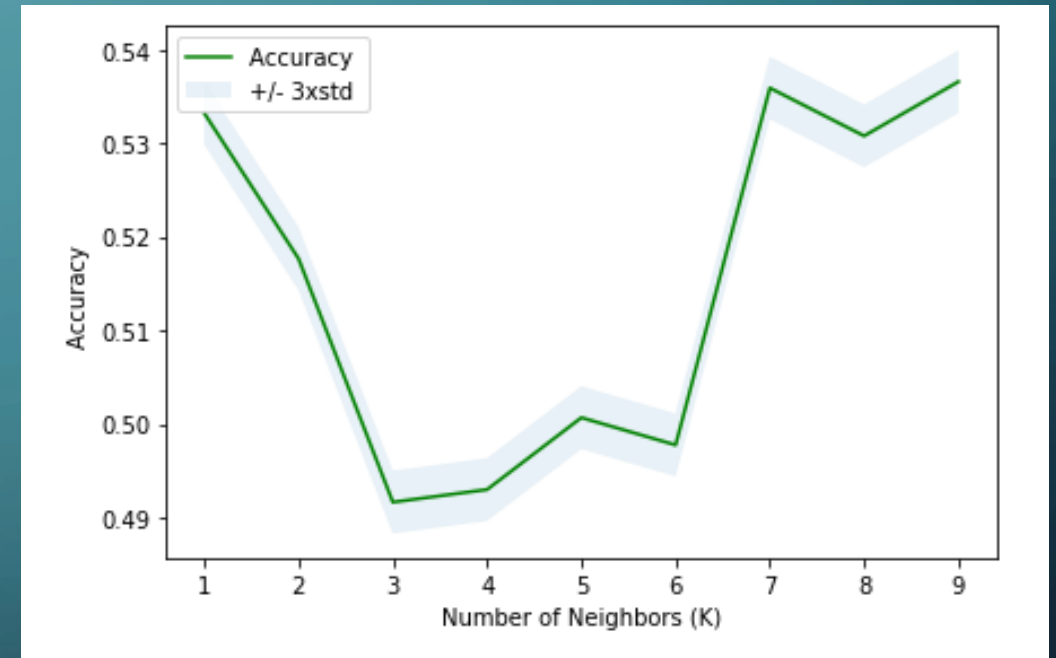- Logistic Regression

# DATA SPLITTING

The 111, 308 observations were split into a training and a test set

- 80% of the data became the training set
- 20% of the data became the test set

# K NEAREST NEIGHBOUR

- An initial value of k = 4 was used

- Further assessment revealed an optimum k value of 9

- The model was retrained with k = 9

- Train set accuracy = 0.542

- Test set accuracy = 0.536

# DECISION TREE

- A decision tree model was fit using the Scikit-learn library

- Prediction accuracy was determined to be 0.596

# SUPPORT VECTOR MACHINE

- An SVM modeled was constructed and fit using the Scikit-learn library

- The model accuracy was determined to be 0.597

# LOGISTIC REGRESSION

- A logistic regression model was developed and fit using the Scikit-learn library

- Model accuracy was determined to be 0.596

```
#predict probabilities
yhat_prob = LR.predict_proba(X_test)
yhat_prob[0:5]

array([[0.58662205, 0.41337795],
       [0.62561619, 0.37438381],
       [0.38966538, 0.61033462],
       [0.58662205, 0.41337795],
       [0.58662205, 0.41337795]])
```

# MODEL METRICS

The metrics performed rather poorly, with Jaccard Index and F1-scores being between 0.54 and 0.6 for all models

Log loss for the Logistic Regression model was higher than desired, at 0.671

| | Algorithm | Jaccard | F1-score | Log Loss |
|---|---|---|---|---|
| 0 | KNN | 0.541309 | 0.528938 | NA |
| 1 | Decision Tree | 0.592859 | 0.592369 | NA |
| 2 | SVM | 0.597046 | 0.596881 | NA |
| 3 | Logistic Regression | 0.593354 | 0.591432 | 0.671243 |

# CONCLUSIONS

- The models definitely need improvement prior to being deployed for any practical purpose

- The models may benefit from greater or fewer attributes, a reduction in the number of potential values associated with each attribute, or from a different set of attributes altogether

- Visualization of purely categorical data is surprisingly challenging!

# FUTURE DIRECTIONS

- The project may be revisited once the inexperienced author further develops her data science and Python skills, as the project has highlighted important areas where further knowledge is needed

- Repeating the project with a dataset for a region with different weather and other conditions may provide some interesting information regarding potential regional differences in accident patterns

- New questions emerged from the data analysis section, such as: "Why are accidents more likely to be severe when they occur in intersections?" This question could potentially be answered with the same data set