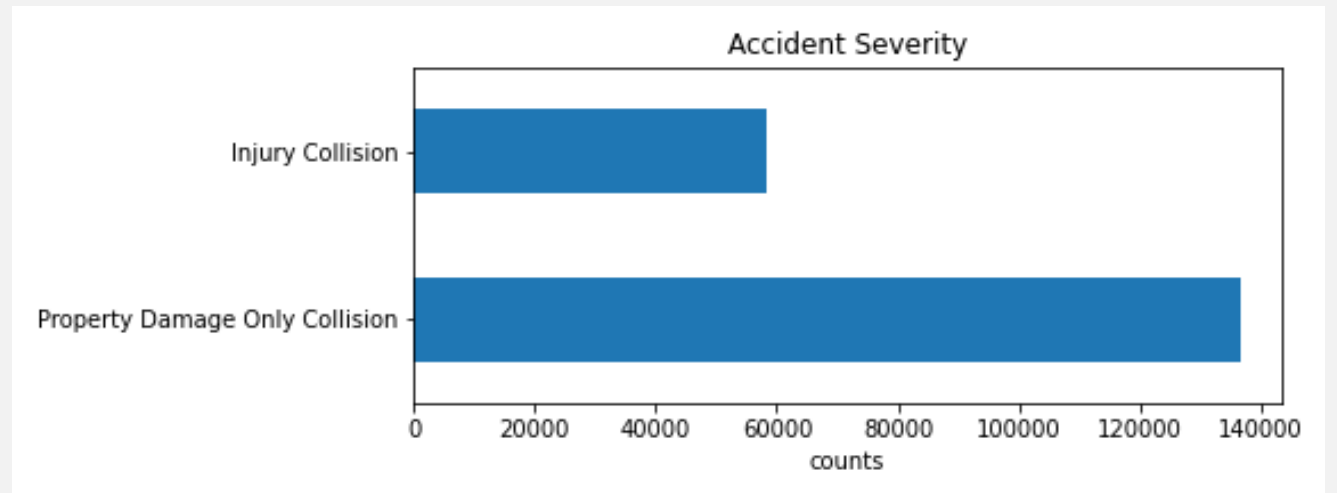# PREDICTING ACCIDENT SEVERITY
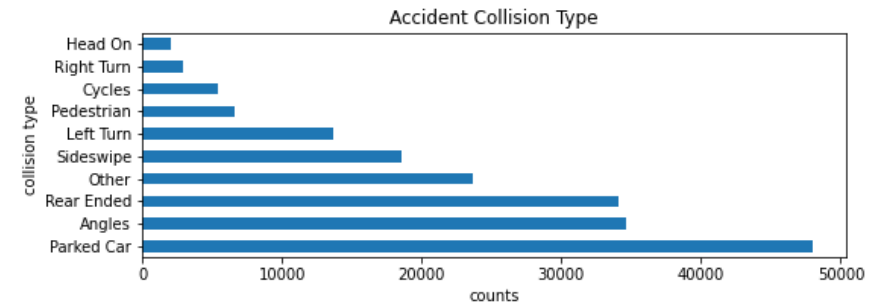
Analisa Hill

September 20th 2020

# ACCIDENT SEVERITY BREAK-DOWN

- Non-injury collision occurs ~70% of the time versus an injury collision ~30% of the time.

- Unbalanced dataset
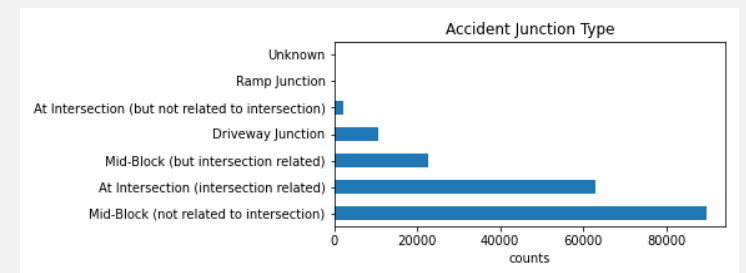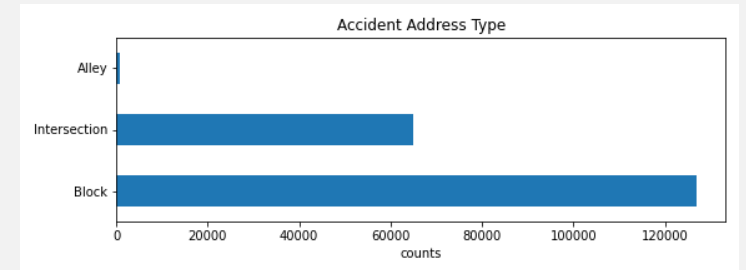
- Future: add data to make it more balanced

# SUMMARY OF ACCIDENT COLLISION TYPE

- Collision types for this dataset are labeled as: head on, right turn, unknown, cycles, pedestrian, left turn, sideswipe, other, rear ended, angles and parked car.

- Accidents involving parked cars occur 25.3%

- Accidents at an angle (18.3%), and then rear-ended (~18.0%).

- The lowest collision types were head-on (1.0%), right-turn (1.6%), and cycles (2.9%).
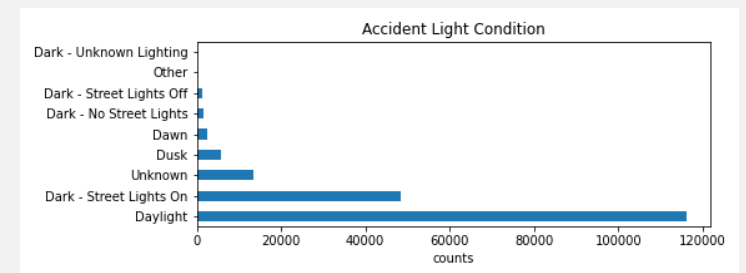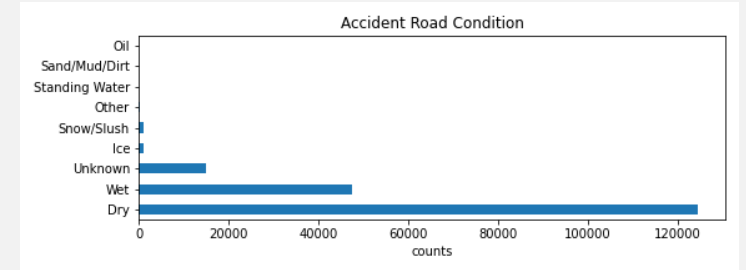
# ROADWAY JUNCTION AND ADDRESS BREAK DOWN



- 66% of the time the accident occurs on a block

- 34% at an intersection

- Less than 1% in an alley

- 48% of these accidents are mid-block and not related to intersections, which coincides with the address type data

- Second prevalent is at an intersection or intersection related, which is also confirmed by the address

# WEATHER, ROAD, AND LIGHT CONDITIONS



- Surprisingly, 66% of the accidents occurred when the road was dry and 25% of the time when the road was wet.

- The weather conditions also confirmed what is seen in the road conditions.

- 59% of the accidents occurred during clear weather and 17% occurred when there was rain.

- The weather is a little more informative than the road condition data because it also gives the option of overcast which contributes almost 15% of the accidents..

# CLASSIFICATION REPORT FOR DIFFERENT MODELS

- Models used: logistic regression, decision tree, random forest

- For each model, the ROC (receiver operator characteristic) curves and AUC (area under the curve) were evaluated.

- The ROC curve calculates the true positive rate versus false positive rate of our classifier. The ROC of each model are comparable.

- The AUC sheds a little more information. The closer the AUC is to one, the better classifier it is.

- Thus, the random forest is a little better than the decision tree and logistic regression with an AUC coming in at 0.7703.



ROC (receiver operator characteristic)

[AUC = 0.7670] | lr
[AUC = 0.7652] | dt
[AUC = 0.7703] | rf



```
classification report for lr
LogisticRegression(C=0.01, solver='liblinear')
              precision    recall  f1-score   support

           1       0.74      0.97      0.84     25317
           2       0.79      0.24      0.37     11262

    accuracy                           0.75     36579
   macro avg       0.76      0.61      0.60     36579
weighted avg       0.76      0.75      0.69     36579


classification report for dt
DecisionTreeClassifier(criterion='entropy', max_depth=10)
              precision    recall  f1-score   support

           1       0.75      0.96      0.84     25317
           2       0.75      0.28      0.40     11262

    accuracy                           0.75     36579
   macro avg       0.75      0.62      0.62     36579
weighted avg       0.75      0.75      0.71     36579


classification report for rf
RandomForestClassifier(max_depth=10)
              precision    recall  f1-score   support

           1       0.74      0.98      0.84     25317
           2       0.81      0.23      0.36     11262

    accuracy                           0.75     36579
   macro avg       0.78      0.60      0.60     36579
weighted avg       0.76      0.75      0.69     36579
```
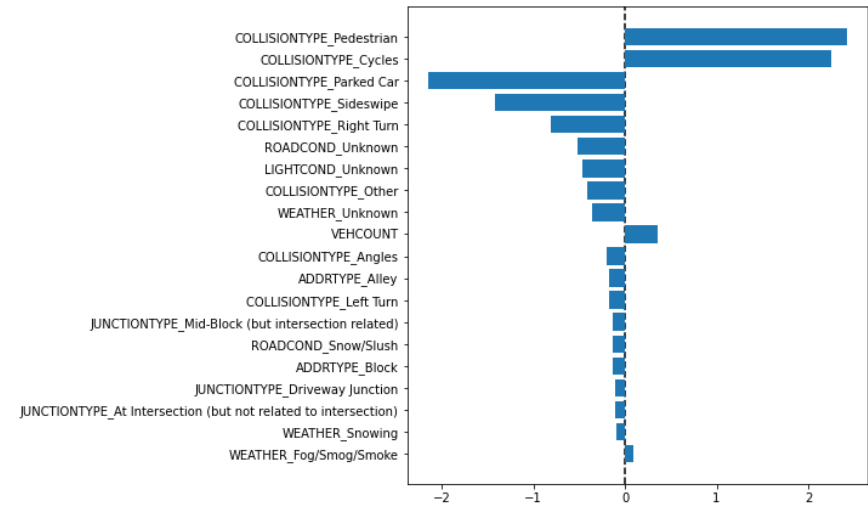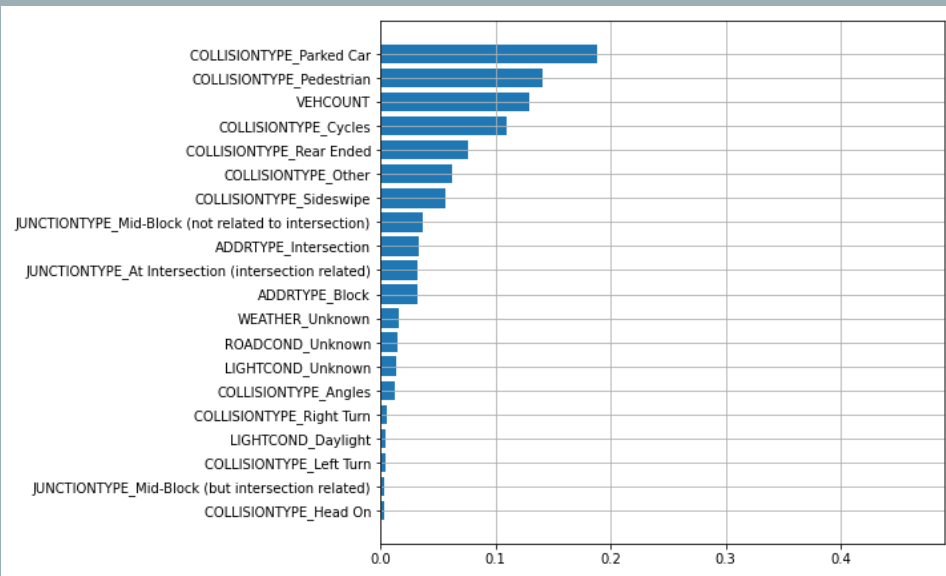
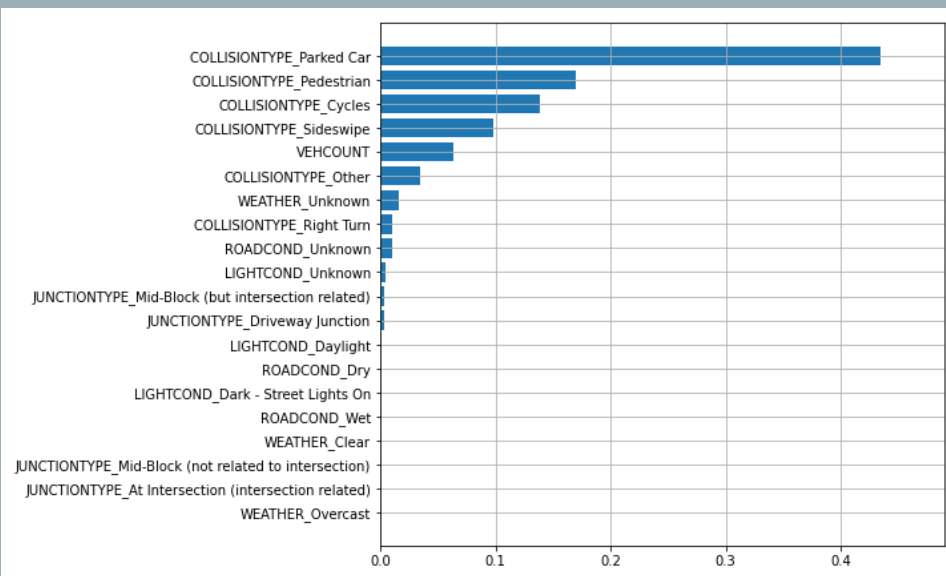# IMPORTANT FEATURES ACQUIRED FROM LOGISTIC REGRESSION MODEL

• Injury accidents are caused when the collision type involves a pedestrian or cyclist.

• Third largest contributing factor is the vehicle count. These are shown by the positive values on the bar plot.

• It makes logical sense that the more vehicles involved in an accident means that there is a higher probability of getting injured because there are more people involved.

• It also makes sense that a non-injury incident happens most with parked cars.

• This non-injury accident is determined by the negative value on the bar plot.
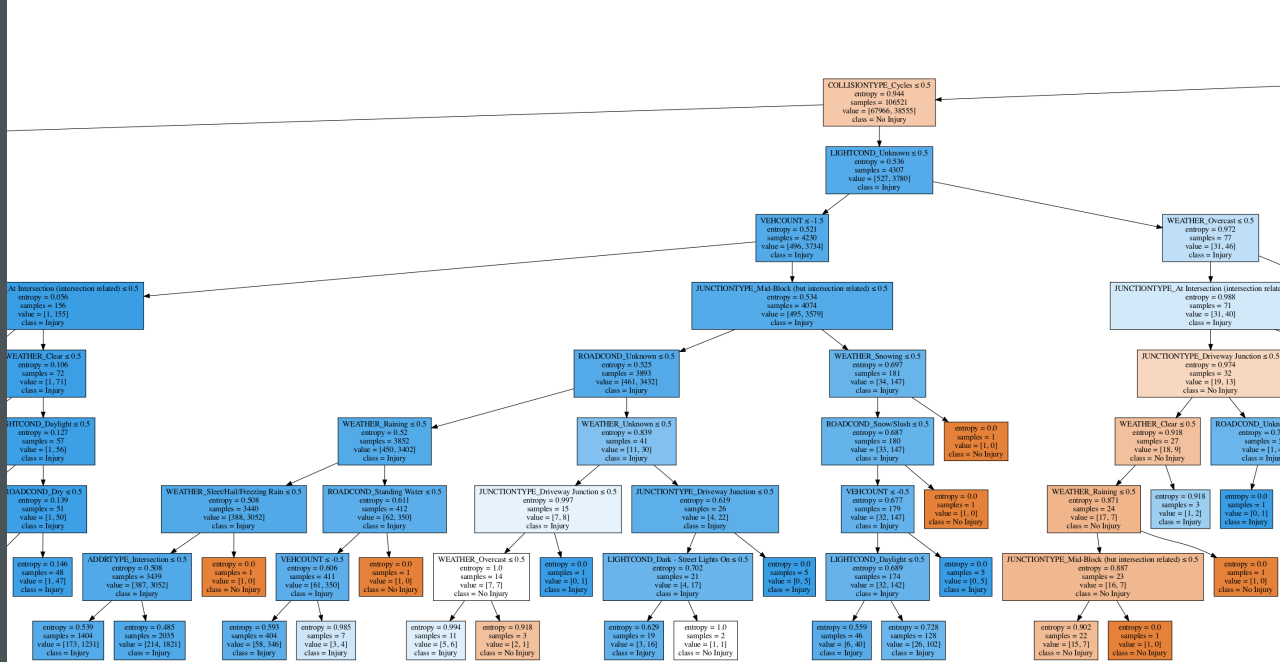
# IMPORTANT FEATURES ACQUIRED FROM DECISION TREE AND RANDOM FOREST MODELS

- These bar plots tell us that parked car and pedestrian collision type plays the largest role in deciding whether it was an injury or no injury accident.

- This is similar to the decision tree and almost nearly the same, but are slightly different because weighted of the ensemble nature of the model, but gives us the same trend as with the decision tree.

- This confirms what we found in the logistic regression model, although the logistic regression model tells us more about which features contribute directly from it's negative and positive values in the bar plot.

DECISION TREE AND WHAT LEADS TO INJURY

DECISION TREE AND WHAT LEADS TO NO INJURY

# CONCLUSIONS

- In this study, I looked at accident severity to predict if someone who was in an accident what is the chance that they were injured.
- The important features that I identified injury accidents are caused when the collision type that involves a pedestrian or cyclist.
- The next contributing factor is the vehicle count.
- I built a logistic regression and decision tree models which turned out to be not as robust as the random forest.
- The random forest model can help first responders determine whether they should dispatch paramedics, police, or none.
- Future directions could be finding a method to determine how many responders should be sent depending on the accident severity.