# IBM Data Science Capstone

## Guidance for pedestrians in Seattle

Jan Stroppel - 10/2020

# Current status

- National Highway Traffic Safety Administraion (NHTSA) sees an increase in the number of accidents

- Accidents with pedestrians often leads to severe injuries

Pedestrian fatalities increased 27% from 2007-2016, while all other traffic deaths decreased by 14%.
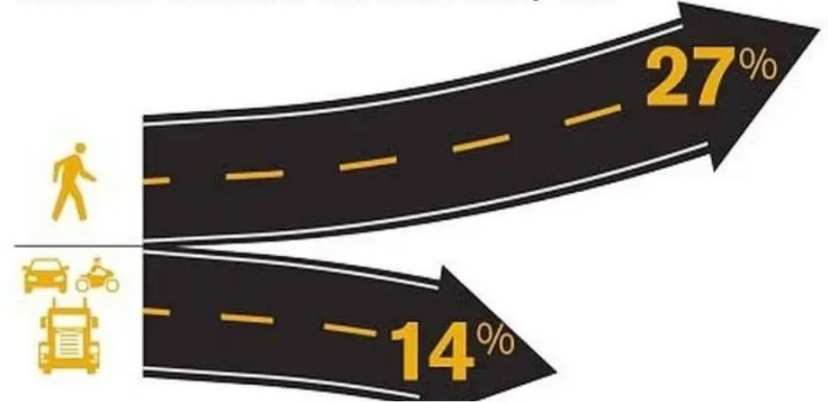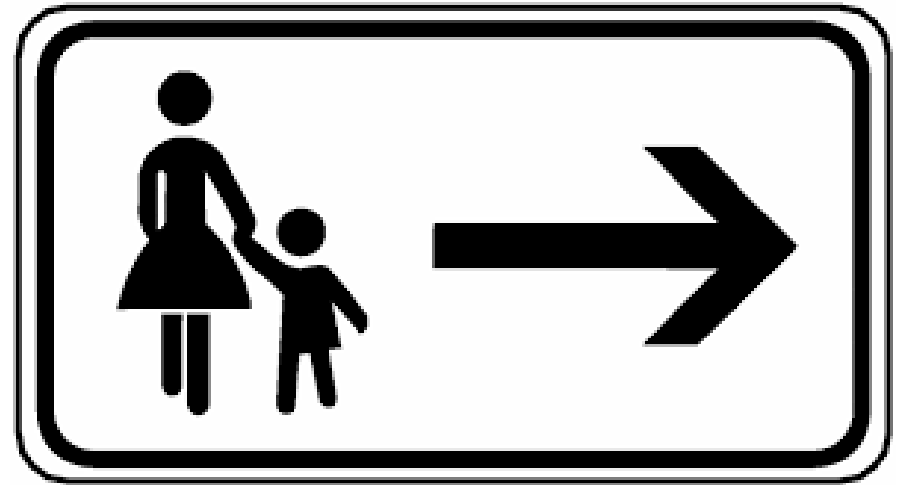
27%

14%

Bild: NHTSA

# Business Problem

- Create a model to predict the probability of pedestrians involed in an accident

- Find parameters which indicates a high risk for pedestrians

- Enable authorities in Seattle to create a guidance for pedestrians

# Data



The dataset „Collisions – All Years" for Seattle was used for creating, training and testing the model

# Selected Columns

- Speeding: Was the vehicle too fast?

- Weather: Was it raining?

- Road Condition: Was the road slippy?

- Light condition: Was the pedestrian visible?

- Date and Time: Are there differences for week days and the time?

- Inattention: Was the driver distracted?

- Drugs: Was the driver under influence of alkohol or drugs?

- Dependent variable was PEDESTRIANINVOLVED, which is a boolean derived from the collision type
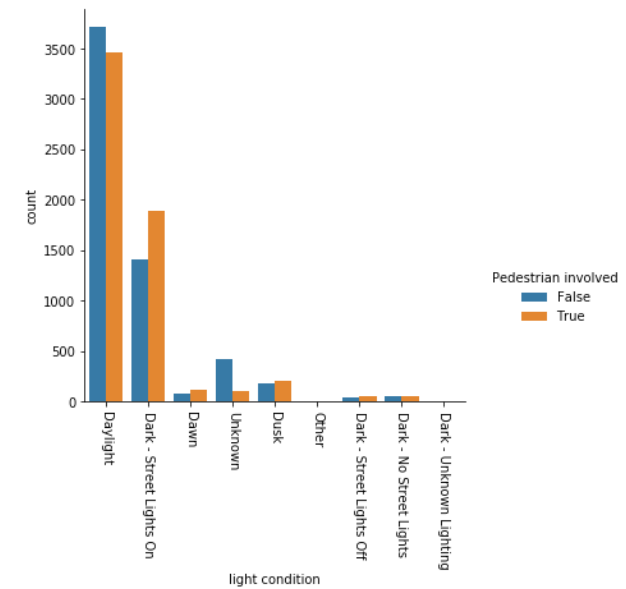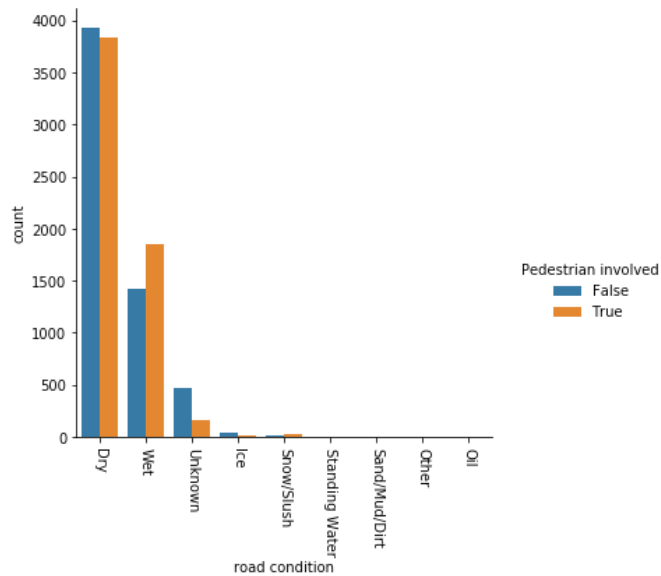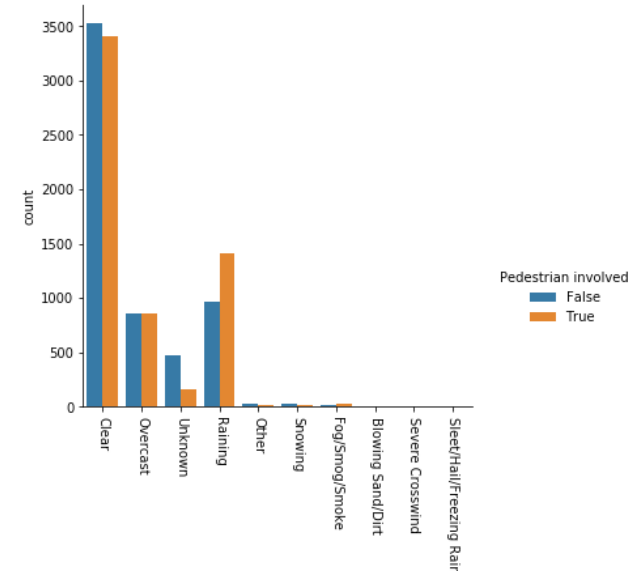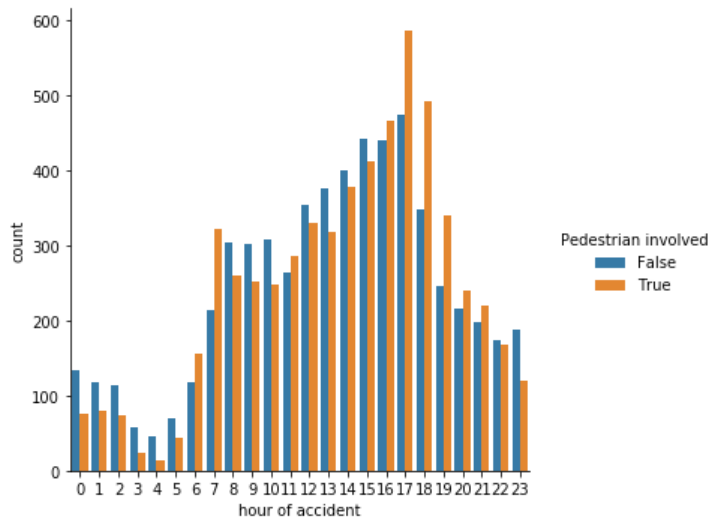
# Data Preparation

- Unify entries for booleans (N,Y instead of 0,1)

- Setting meaningful default values (Unknown for missing data)

- Transforming categorical values into numerical values (0,1 instead of Dry, Wet)

- Balancing data set (same number of rows for involved/not involved pedestrians)

```
Attribute values count:
Clear                       96391
Raining                     28699
Overcast                    23831
Unknown                     13082
Snowing                       776
Other                         678
Fog/Smog/Smoke                521
Sleet/Hail/Freezing Rain       90
Blowing Sand/Dirt              50
Severe Crosswind               24
Partly Cloudy                   5
Name: WEATHER, dtype: int64
```

```python
cat_col = df_balanced.select_dtypes(['object'])
encoding_maps = []

for column in cat_col:
    df_balanced[column] = pd.Categorical(df_balanced[column]).codes
df_balanced.head()
```
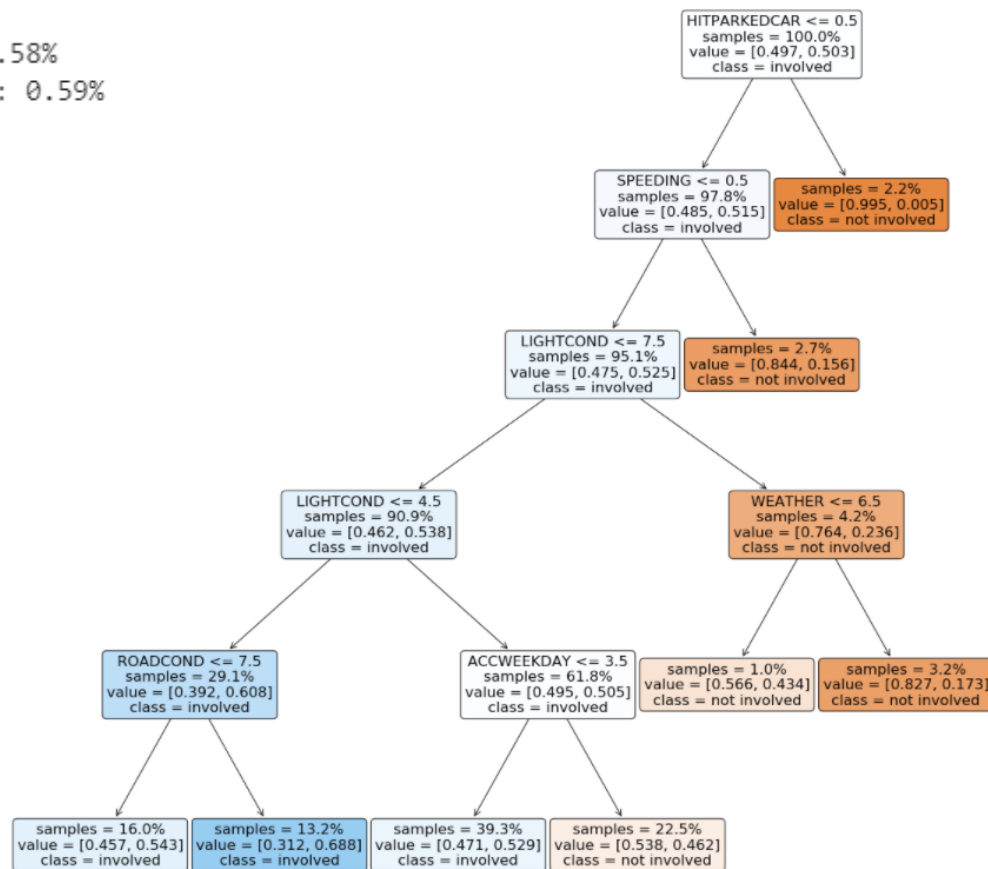
# Data visualization

# Methodology

- As most variables are categorical, a classifier model was chosen

- It was decided to create a decision tree

- The dependent variable is binary, so ideal for a decision tree

- The dataset was split 30/70 in training/test set

# Decision Tree



Accuracy: 0.58%
Recall score: 0.58%
Precision score: 0.59%
F1 score: 0.56%

# Evaluation

- Model performance low (50 to 60%)

- Reason could be not enough data to train and test the model (only 5891 entries for each value)

- Another reason could be that business problem can not be handled by the data set