README.MD Technical Presentation and final cleanup 14 hours ago Technical Presentation.ipynb Technical Presentation and final cleanup 14 hours ago **README.MD** 

## Using Machine Learning for determining Car Accidents

**Phase 3 Project:** 

# causes



#### The objective of the

project:

**OVERVIEW** 

Contributors 2

Auustiino

Jupyter Notebook 99.6%

Languages

Python 0.4%

schahmatist Dmitriy F

Using traffic accident data published on Chicago Data portal to determine if an accident was caused by a driver **DATA** 

#### The row data consists of 3

**Crashes** 

People

Crash data shows information about each traffic crash on city streets within the City of Chicago limits

#### This data contains information about people involved in a crash and if any injuries were sustained.

**Vehicles** 

This dataset contains information about vehicles involved in a traffic crash.

• All three datasets are gziped and saved in data/raw folder

### People and Vehicles datasets were joined using VEHICLE\_ID and CRASH\_RECORD\_ID

**DATA PREPARATION** 

- Crashes dataset was joined with People/Vehicles using CRASH\_RECORD\_ID
- Datasets were filtered to consider only Drivers (passengers, pedestrians, etc were excluded)
- Data was cleaned (Rows with unknown or missing data were removed) • Target variable was bined into two categories: "AT\_FAULT" and "NOT\_AT\_FAULT'
- All the steps for Data preparation and cleaning were put into: "initial\_prep.py" script

The processed data was generated by initial\_prep.py as "crashes.gz" and placed into data/processed folder

#### We considered the following models:

**MODELING** 

• Logistic Regression

- Decision Tree • Extra Tree
- Nearest Neighbor (KNN)
- Random Forest
- Bayesian
- GXBoost
- **Features Selection**

We also tried "bagging" and "voting" Ensembles of the models above

#### • We started with 20 features and filtered them one by one until we end up with only 5. **RESULTS**

GXBoost model was chosen as marginally better model (close second and third were Random Forest and Logistic Regression)

# Uncomment only if you want to reprocess the raw data using initial\_prep.py

## %run src/preprocessing/initial\_prep.py

# The script will re-create crashes.gz, comma separated gzipped file in ../data/processed

# Import required sklearn, pandas, numpy and other libraries: %run src/import\_libraries.py

Loading processed data

%matplotlib inline

(348442, 154)

## full\_df.shape

sample\_df=full\_df.copy()

**Defining Training and testing sets** 

full\_df=pd.read\_csv('data/processed/crashes.gz', compression='gzip', low\_memory=False)

```
y = sample_df['GUILTY']
x = sample_df.drop(['GUILTY'],axis=1)
X_train, X_test, y_train, y_test = train_test_split(x, y, random_state=100, test_size=0.35, strat
nominal_columns = [ 'FIRST_CRASH_TYPE', 'FIRST_CONTACT_POINT', 'MANEUVER', 'PHYSICAL_CONDITION',
                   'TRAFFIC_CONTROL_DEVICE', 'VEHICLE_DEFECT']
```

#sample\_df=full\_df.sample(100000, random\_state=100)

X\_train = X\_train[nominal\_columns] X\_test = X\_test[nominal\_columns]

## from sklearn.base import BaseEstimator, TransformerMixin

cat\_pipeline = Pipeline([

Defining a pipeline

```
('cat_encoder', OneHotEncoder(sparse=False, handle_unknown = 'ignore')),
         ('std_scaler', StandardScaler()),
    ])
XGB_pipeline = Pipeline([ ( "cat_pipeline", cat_pipeline ),
                  ('boost', XGBClassifier(use_label_encoder=False,
                                         eval_metric='auc', n_estimators=150, n_jobs=-1))
                ])
```

**Training and Predicting** 

pipe\_grid={}

gs\_pipe.fit(X\_train, y\_train) y\_pred\_train=gs\_pipe.predict(X\_train) y\_pred\_test=gs\_pipe.predict(X\_test)

gs\_pipe = GridSearchCV(estimator=XGB\_pipeline, param\_grid=pipe\_grid, cv=3, scoring='roc\_auc')

## **EVALUATION**

```
#print(XGB_pipeline.score(X_train, y_train ))
print(gs_pipe.score(X_train, y_train))
print(classification_report(y_train, y_pred_train))
print('----')
print(classification_report(y_test, y_pred_test))
fig, ax = plt.subplots(figsize=(13,6), ncols=2)
ConfusionMatrixDisplay.from_predictions(y_test, y_pred_test, normalize='pred',
                                       display_labels=[1,0],cmap='pink', ax=ax[0])
ConfusionMatrixDisplay.from_predictions(y_test, y_pred_test, normalize='true',
```

ax[0].set\_title("Precision for predictions") ax[1].set\_title("Recall for true cases") ax[0].set\_xlabel('Predicted Cases')

ax[1].set\_ylabel('True Cases'); Our model is much better in determining not\_at\_fault (1) class than "at\_fault" (0).

ax[0].set\_ylabel('True Cases')

ax[1].set\_xlabel('Predicted Cases')

Still it can be helpful for Legal and Insurance Companies and lawyers investigating traffic accidents

© 2022 GitHub, Inc. Terms Privacy Training Security Status Contact GitHub Pricing About

display\_labels=[1,0],cmap='pink', ax=ax[1]);