#### Contents 2 ♥

```
0.1 Global Func
▼ 1 Complete jupyte
  ▼ 1.1 Logistic Reg
    ▼ 1.1.1 Model 2
        1.1.1.1 Mod
    ▼ 1.1.2 Model 3
        1.1.2.1 Mod
    ▼ 1.1.3 Model 4
        1.1.3.1 Mod
      1.1.4 Best LR
  ▼ 1.2 Random Fo
     ▼ 1.2.1 Model 1
        1.2.1.1 Mod
    ▼ 1.2.2 Model 2
        1.2.2.1 Moc
  2 Data Interpretat
```

3 Recommendation

#### In [4]:

```
#data wrangling and visualization packages
   import pandas as pd
   import numpy as np
 3
   import matplotlib.pyplot as plt
   import seaborn as sns
 6
   plt.style.use('fivethirtyeight')
   import statsmodels.api as sm
7
   import scipy.stats as stats
9
   #feature engineering packages
10
   from sklearn.impute import SimpleImputer
11
12
   from sklearn.preprocessing import OneHotEncoder, StandardS
13
14
   #feature selection packages
   from feature engine.selection import DropDuplicateFeatures
15
16
   from feature_engine.selection import DropConstantFeatures
17
18
   #modeling packages
   from sklearn.model_selection import train_test_split
   from sklearn.dummy import DummyClassifier
20
21
   from sklearn.linear_model import LogisticRegression, Logis
   from sklearn.pipeline import Pipeline
   from sklearn.feature_selection import SelectFromModel
24
   from sklearn.ensemble import RandomForestClassifier
25
26 #modeling evaluation packages
   from sklearn.metrics import precision_score
27
   from sklearn.metrics import recall score
29 from sklearn.metrics import accuracy_score
30 from sklearn.metrics import f1_score
   from sklearn.model_selection import cross_val_score
   from sklearn.metrics import plot_confusion_matrix
   from sklearn.metrics import confusion_matrix
   from sklearn.metrics import classification_report
   from sklearn.metrics import plot_roc_curve, roc_curve, auc
   from sklearn.metrics import get_scorer
36
37
38 #optimization packages
39 from sklearn.model selection import GridSearchCV
```

## 0.1 Global Function



1 #function to look at plots and stats of column with or wit

#### In [5]:

	<b>▼</b> 2 3	<pre>def model_eval(model, X_train, y_train, X_test, y_test,</pre>
	4	prev_X_test=None, prev_y_test=None):
	5	prev_x_test=None; prev_y_test=None).
Contents <i>⊋</i> <b>♦</b>	6	This function takes in a fit model and provides classi
o.1 Global Func	7	metrics on that model. Optionally, a previous model ca
▼ 1 Complete jupyte	8	improvement metrics between the current model and the
▼ 1.1 Logistic Reg	9	Keyword Arguments:
•	10	- model: A fit model
▼ 1.1.1 Model 2	11	- X_train, y_train, X_test, y_test: Training and testi
1.1.1.1 Mod	12	the "model" stated above was fit on.
▼ 1.1.2 Model 3	13	- prev model: Another fit model
1.1.2.1 Mod	14	<ul><li>prev_X_train, prev_y_train, prev_X_test, prev_y_test</li></ul>
▼ 1.1.3 Model 4	15	testing dataframes which the previous model was fit on
1.1.3.1 Mod	16	iii
1.1.4 Best LR	17	#current model predictions on testing dataframe
▼ 1.2 Random Fo	18	<pre>y_hat_test = model.predict(X_test)</pre>
▼ 1.2.1 Model 1	19	#get scores of current model on testing dataframe
1.2.1.1 Mod	20	recall_model = get_scorer('recall')(model, X_test, y_t
▼ 1.2.2 Model 2	21	f1_model = get_scorer('f1')(model, X_test, y_test).rou
1.2.2.1 Moc	22	accuracy_model = get_scorer('accuracy')(model, X_test,
2 Data Interpretat	23	<pre>auc_model = get_scorer('roc_auc')(model, X_test, y_test)</pre>
3 Recommendation	24	dac_model gee_seerer ( rec_ade /(model) //_eese, y_ees
3 Recommendant	25	<pre>recall_model_train = get_scorer('recall')(model, X_tra</pre>
	26	#if statement to check for availability of a previous
	27	#and testing dataframes
	<b>▼</b> 28	<pre>if prev_model != None:</pre>
	<b>▼</b> 29	if prev_X train is None:
	30	#if prevoius model has a different training an
	31	may provided model had a day for one or a smany and
	32	#get previous model predictions and scores
	33	y hat test prev = prev model.predict(X test)
	34	)a.c_cosc_p. c.
	35	<pre>recall_prev = get_scorer('recall')(prev_model,</pre>
	36	f1_prev = get_scorer('f1')(prev_model, X_test,
	37	accuracy_prev = get_scorer('accuracy')(prev_mo
	38	auc_prev = get_scorer('roc_auc')(prev_model, X
	39	add_prev gee_seerer( red_ade /(prev_modes) //
	40	<pre>print('MODEL EVAL VS PREVIOUS (TEST)')</pre>
	41	print('====================================
	42	F•(
	43	#create dataframe comparing current and previo
	<b>▼</b> 44	<pre>df = pd.DataFrame(index=['Recall','F1','Accura</pre>
	45	columns=['Previous Model','C
	46	, ,
	47	<pre>df.loc['Recall','Current Model'] = recall mode</pre>
	48	<pre>df.loc['Recall','Previous Model'] = recall_pre</pre>
	49	df.loc['F1','Current Model'] = f1_model
	50	df.loc['F1','Previous Model'] = f1_prev
	51	<pre>df.loc['Accuracy','Current Model'] = accuracy_</pre>
	52	<pre>df.loc['Accuracy','Previous Model'] = accuracy</pre>
	53	<pre>df.loc['AUC','Current Model'] = auc_model</pre>
	54	<pre>df.loc['AUC','Previous Model'] = auc_prev</pre>
	55	df.loc['Recall','Delta'] = (recall_model - rec
	56	df.loc['F1','Delta'] = (f1_model - f1_prev).ro
	57	df.loc['Accuracy','Delta'] = (accuracy_model -
4	58	df.loc['AUC','Delta'] = (auc_model - auc_prev)

59

Ontents   O.1 Global Func  1 Complete jupyte  1.1 Logistic Reg  1.1.1 Model 2  1.1.1.1 Model 3  1.1.2.1 Mod  1.1.3 Model 4  1.1.3.1 Mod  1.1.4 Best LR  1.2 Random Fo  1.2.1 Model 1  1.2.1.1 Mod  1.2.2 Model 2  1.2.2.1 Moc  2 Data Interpretat  3 Recommendatic
<b>←</b>

```
60
                display(df)
 61
                print('\n')
 62
                #display current model scores, reports and gra
 63
 64
                print(f"CURRENT MODEL: {'Overfit (Recall)' if
                65
                print('\n')
 66
 67
                print(f'Recall on Training: {recall model trai
 68
                print(f'Recall on Test: {recall_model}')
 69
 70
                print('\n')
 71
 72
                print('Classification Reports-----
 73
                print(classification_report(y_test, y_hat_test
 74
 75
                print('Test Graphs------
 76
               fig, ax = plt.subplots(ncols=2, figsize=(15,8)
 77
 78
                plot_confusion_matrix(model, X_test, y_test, c
 79
                                           normalize='true',
                                    display_labels=["NO INJU
 80
 81
                                    ax=ax[0]);
 82
 83
                plot_roc_curve(model, X_test, y_test, ax=ax[1]
 84
               plt.show()
 85
 86
                #display previous model graphs
                print("PREVIOUS MODEL")
 87
 88
                89
                fig, ax = plt.subplots(ncols=2, figsize=(15,8)
 90
 91
                plot_confusion_matrix(prev_model, X_test, y_te
                                           normalize='true',
 92
                                    display_labels=['NO INJU
 93
 94
                                    ax=ax[0]);
 95
 96
                plot_roc_curve(prev_model, X_test, y_test, ax=
 97
                plt.show()
 98
 99
                plt.tight_layout()
            else:
100
                #if previous model has the same dataframes as
101
               y_hat_test_prev = prev_model.predict(prev_X_te
102
103
                recall prev = get scorer('recall')(prev model,
104
105
                                                 prev_y_test
106
                f1_prev = get_scorer('f1')(prev_model, prev_X_
107
                                         prev_y_test).round(
108
                accuracy_prev = get_scorer('accuracy')(prev_mo
109
                                                    prev_y_
110
                auc_prev = get_scorer('roc_auc')(prev_model, p
111
                                               prev y test).
112
                print('MODEL EVAL VS PREVIOUS (TEST)')
113
114
                115
116
                df = pd.DataFrame(index=['Recall','F1','Accura
                                 columns=['Previous Model','C
117
118
                df.loc['Recall','Current Model'] = recall mode
119
                df.loc['Recall','Previous Model'] = recall_pre
120
```

# Contents 2 \* 0.1 Global Func ▼ 1 Complete jupyte ▼ 1.1 Logistic Reg ▼ 1.1.1 Model 2 1.1.1.1 Mod ▼ 1.1.2 Model 3 1.1.2.1 Mod ▼ 1.1.3 Model 4 1.1.3.1 Mod 1.1.4 Best LR ▼ 1.2 Random Fo ▼ 1.2.1 Model 1 1.2.1.1 Mod ▼ 1.2.2 Model 2 1.2.2.1 Mod 2 Data Interpretat 3 Recommendation

```
121
               df.loc['F1','Current Model'] = f1_model
122
               df.loc['F1','Previous Model'] = f1_prev
               df.loc['Accuracy','Current Model'] = accuracy_
123
124
               df.loc['Accuracy','Previous Model'] = accuracy
               df.loc['AUC','Current Model'] = auc_model
125
               df.loc['AUC','Previous Model'] = auc_prev
126
               df.loc['Recall','Delta'] = (recall_model - rec
127
               df.loc['F1','Delta'] = (f1_model - f1_prev).ro
128
               df.loc['Accuracy','Delta'] = (accuracy_model -
129
               df.loc['AUC','Delta'] = (auc_model - auc_prev)
130
131
132
               display(df)
133
               print('\n')
134
135
136
               print(f"CURRENT MODEL: {'Overfit (Recall)' if
               137
138
               print('\n')
139
               print(f'Recall on Training: {recall_model_trai
140
               print(f'Recall on Test: {recall_model}')
141
142
               print('\n')
143
144
               print('Classification Reports-----
145
               print(classification_report(y_test, y_hat_test
146
147
               print('Test Graphs-----
148
               fig, ax = plt.subplots(ncols=2, figsize=(15,8)
149
150
               plot_confusion_matrix(model, X_test, y_test, c
151
                                          normalize='true',
                                   display_labels=['NO INJU
152
153
                                    ax=ax[0]);
154
155
               plot_roc_curve(model, X_test, y_test, ax=ax[1]
156
               plt.show()
157
               print("PREVIOUS MODEL")
158
159
               fig, ax = plt.subplots(ncols=2, figsize=(15,8)
160
161
162
               plot_confusion_matrix(prev_model, prev_X_test,
163
                                          normalize='true',
164
                                    display_labels=['NO INJU
165
                                    ax=ax[0]);
166
               plot_roc_curve(prev_model, prev_X_test, prev_y
167
168
                             ax=ax[1]).ax_.plot([0,1],[0,1])
169
               plt.show()
170
               plt.tight_layout()
171
        else:
172
173
            #if there is no previous model, get current model
174
            print(f"CURRENT MODEL: {'Overfit (Recall)' if reca
175
            print('\n')
176
177
            print('Classification Reports------
178
            print(classification_report(y_test, y_hat_test))
179
180
            print('Test Graphs-----
181
           fig, ax = plt.subplots(ncols=2, figsize=(15,8))
```

```
Contents 2 ♥
```

```
o.1 Global Func

▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2
```

- ▼ 1.1.1 Model 2 1.1.1.1 Model
- ▼ 1.1.2 Model 3 1.1.2.1 Mod
- ▼ 1.1.3 Model 4 1.1.3.1 Mod 1.1.4 Best LR
- ▼ 1.2 Random Fo
  - ▼ 1.2.1 Model 1
    - 1.2.1.1 Mod
  - ▼ 1.2.2 Model 2 1.2.2.1 Mod
- 2 Data Interpretat
- 3 Recommendation

```
182
183
             plot_confusion_matrix(model, X_test, y_test, cmap=
184
                                            normalize='true',
185
                                    display_labels=['NO INJURY',
186
                                    ax=ax[0]);
187
188
             plot_roc_curve(model, X_test, y_test, ax=ax[1]).ax
189
             plt.show()
190
191
192
         return
```

#### In [50]:

#### In [6]:

```
1 X_train_dt = pd.read_csv("crash_train.csv")
2 X_test_dt = pd.read_csv("crash_test.csv")
```

#### In [7]:

```
1 X_train_lr_sel = X_train_dt.copy()
2 X_test_lr_sel = X_test_dt.copy()
```

#### In [8]:

```
1  y_train = pd.read_csv("crash_y_train.csv")
2  y_test = pd.read_csv("crash_y_test.csv")
3
```

#### In [9]:

```
y_train.drop("Unnamed: 0", axis=1, inplace=True)
y_test.drop("Unnamed: 0", axis=1, inplace=True)
```

#### In [10]:

```
1 X_test_dt.drop("Unnamed: 0", axis=1, inplace=True)
2 X_train_dt.drop("Unnamed: 0", axis=1, inplace=True)
```

#### In [11]:

1 X\_train\_dt

#### Out[11]:

Contents <i>⊋</i> ❖
o.1 Global Func
▼ 1 Complete jupyte
▼ 1.1 Logistic Reg
▼ 1.1.1 Model 2
1.1.1.1 Mod
▼ 1.1.2 Model 3
1.1.2.1 Mod
▼ 1.1.3 Model 4
1.1.3.1 Mod
1.1.4 Best LR
▼ 1.2 Random Fo
▼ 1.2.1 Model 1
1.2.1.1 Mod
▼ 1.2.2 Model 2
1.2.2.1 Moc
2 Data Interpretat

3 Recommendation

	LIGHTING_CONDITION_DARKNESS, LIGHTED ROAD	LIGHTING_CONDITION_DAWN	LIGHT
0	0.0	0.0	
1	1.0	0.0	
2	0.0	0.0	
3	0.0	0.0	
4	0.0	0.0	
82794	0.0	0.0	
82795	1.0	0.0	
82796	0.0	0.0	
82797	0.0	0.0	
82798	1.0	0.0	

# 1 Complete jupyter nb

## 1.1 Logistic Regression

#### 1.1.1 Model 2

82799 rows × 360 columns

#### In [12]:

```
1 lr_2 = LogisticRegression()
2 lr_2.fit(X_train_lr_sel, y_train)
```

c:\python\python37\lib\site-packages\sklearn\utils\validation.py
versionWarning: A column-vector y was passed when a 1d array was
lease change the shape of y to (n\_samples, ), for example using
 return f(\*\*kwargs)

#### Out[12]:

LogisticRegression()

#### 1.1.1.1 Model Evaluation

#### 1.1.2 Model 3

```
Contents 2 5
    0.1 Global Func
▼ 1 Complete jupyte
  ▼ 1.1 Logistic Reg
    ▼ 1.1.1 Model 2
        1.1.1.1 Mod
    ▼ 1.1.2 Model 3
         1.1.2.1 Mod
    ▼ 1.1.3 Model 4
         1.1.3.1 Mod
      1.1.4 Best LR
  ▼ 1.2 Random Fo
     ▼ 1.2.1 Model 1
         1.2.1.1 Mod
    ▼ 1.2.2 Model 2
        1.2.2.1 Mod
  2 Data Interpretat
  3 Recommendation
```

```
In [13]:
   1 lr 3 = LogisticRegression()
In [14]:
      params_3 = {"C": [0.001, 0.01, 0.1, 1, 1e7, 1e8, 1e9, 1e10]}
                 "penalty": ["11", "12"],
   2
                 "solver":["liblinear"]}
   3
In [15]:
      lr_3_gridsearch = GridSearchCV(lr_3, params_3,
                                      scoring="recall", n_jobs=-1
   3 | lr_3_gridsearch.fit(X_train_lr_sel, y_train)
c:\python\python37\lib\site-packages\sklearn\utils\validation.py
versionWarning: A column-vector y was passed when a 1d array was
lease change the shape of y to (n_samples, ), for example using
  return f(**kwargs)
Out[15]:
GridSearchCV(cv=5, estimator=LogisticRegression(), n_jobs=-1,
             param_grid={'C': [0.001, 0.01, 0.1, 1, 10000000.0, :
                               1000000000.0, 10000000000.0],
                         'penalty': ['l1', 'l2'], 'solver': ['lil
```

#### 1.1.2.1 Model Evaluation

scoring='recall')

```
In [16]:
    1 lr_3_gridsearch.best_params_
Out[16]:
{'C': 1, 'penalty': 'l1', 'solver': 'liblinear'}
```

#### In [17]:

1 model\_eval(lr\_3\_gridsearch.best\_estimator\_, X\_train\_lr\_sel

#### MODEL EVAL VS PREVIOUS (TEST)

\_\_\_\_\_\_

	Previous Model	Current Model	Delta
Recall	0	0.8	0.8
F1	0	0.83	0.83
Accuracy	0.67	0.89	0.22
AUC	0.51	0.96	0.45

CURRENT MODEL: Not Overfit (Recall)

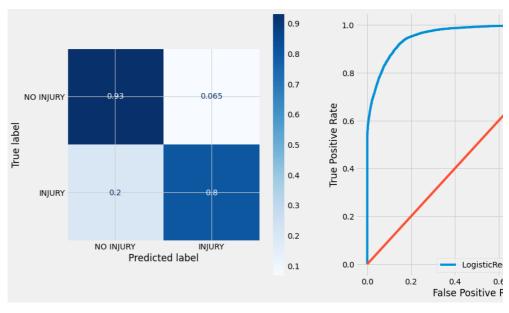
-----

Recall on Training: 0.81

Recall on Test: 0.8

Classifica	ation	Reports				
CIUSSITIC		precision	recall	f1-score	support	
	0	0.91	0.93	0.92	23838	
	1	0.86	0.80	0.83	11648	
accura	асу			0.89	35486	
macro a	avg	0.88	0.87	0.87	35486	
weighted a	avg	0.89	0.89	0.89	35486	

Test Graphs------



PREVIOUS MODEL

#### Contents 2 &

0.1 Global Func▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2 1.1.1.1 Model

▼ 1.1.2 Model 3

1.1.2.1 Mod

▼ 1.1.3 Model 4 1.1.3.1 Mod

1.1.4 Best LR

▼ 1.2 Random Fo

▼ 1.2.1 Model 1

1.2.1.1 Mod ▼ 1.2.2 Model 2

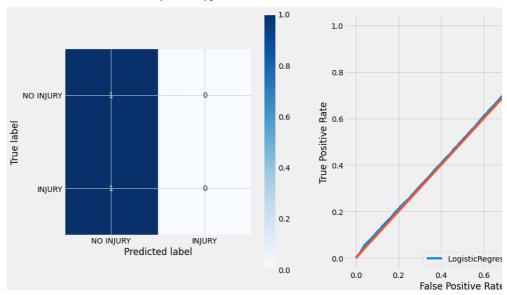
1.2.2.1 Mod

2 Data Interpretat3 Recommendation

localhost:8888/notebooks/student-complete.ipynb#

## Contents **₽** ❖

- 0.1 Global Func
- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg ▼ 1.1.1 Model 2
    - 1.1.1.1 Model 2
    - ▼ 1.1.2 Model 3
    - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
      - 1.1.3.1 Mod 1.1.4 Best LR
  - ▼ 1.2 Random Fo
  - ▼ 1.2.1 Model 1
    - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
      - 1.2.2.1 Mod
  - 2 Data Interpretat
  - 3 Recommendation



<Figure size 432x288 with 0 Axes>

**Observations** Feature selection did not seem to neagtvely impact the n feature count by 21%.

Actions See if fizing class imbalance will help the model

#### 1.1.3 Model 4

```
In [18]:
```

```
1 lr_4 = LogisticRegression(class_weight="balanced")
```

## In [19]:

localhost:8888/notebooks/student-complete.ipynb#

Contents 2 5

▼ 1 Complete jupyte

0.1 Global Func

```
2 lr_4_gridsearch.fit(X_train_lr_sel, y_train)
```

c:\python\python37\lib\site-packages\sklearn\utils\validation.py versionWarning: A column-vector y was passed when a 1d array was lease change the shape of y to (n\_samples, ), for example using return f(\*\*kwargs)

1 lr\_4\_gridsearch = GridSearchCV(lr\_4, params\_4, scoring="re

```
Out[20]:
```

In [20]:

```
GridSearchCV(cv=5, estimator=LogisticRegression(class_weight='bal
             n_jobs=-1,
             param_grid={'C': [0.01, 0.1, 1, 10000000.0, 10000000
                               1000000000.0, 10000000000.0],
                         'penalty': ['l1', 'l2'], 'solver': ['lil
             scoring='recall')
```

#### In [21]:

```
1 lr_4_gridsearch.best_params_
```

#### Out[21]:

```
{'C': 0.01, 'penalty': 'l1', 'solver': 'liblinear'}
```

#### 1.1.3.1 Model Evaluation

#### In [22]:

<b>v</b> 1	<pre>model_eval(lr_4_gridsearch.best_estimator_, X_train_lr_sel</pre>
<b>▼</b> 2	<pre>X_test_lr_sel,y_test,</pre>
3	<pre>prev_model=lr_3_gridsearch.best_estimator_ )</pre>

MODEL EVAL VS PREVIOUS (TEST)

\_\_\_\_\_\_

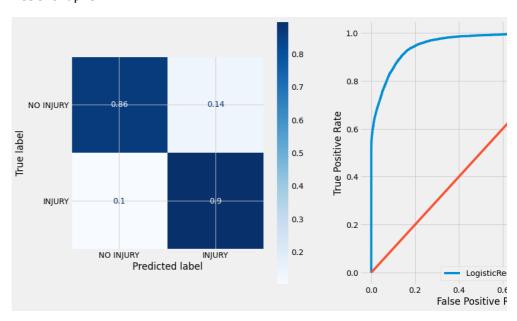
	Previous Model	<b>Current Model</b>	Delta
Recall	0.8	0.9	0.1
F1	0.83	0.83	0
Accuracy	0.89	0.88	-0.01
AUC	0.96	0.96	0

CURRENT MODEL: Not Overfit (Recall)

Recall on Training: 0.9 Recall on Test: 0.9

Classification	n Reports				
	precision	recall	f1-score	support	
0	0.95	0.86	0.90	23838	
1	0.76	0.90	0.83	11648	
accuracy			0.88	35486	
macro avg	0.86	0.88	0.86	35486	
weighted avg	0.89	0.88	0.88	35486	

Test Graphs------



PREVIOUS MODEL

## Contents 2 ♥

0.1 Global Func▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2 1.1.1.1 Model

▼ 1.1.2 Model 3

1.1.2.1 Mod

▼ 1.1.3 Model 4 1.1.3.1 Mod

1.1.4 Best LR
▼ 1.2 Random Fo

▼ 1.2.1 Model 1

1.2.1.1 Mod

▼ 1.2.2 Model 2

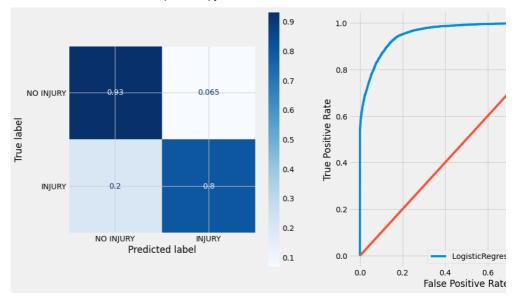
1.2.2.1 Mod

2 Data Interpretat

3 Recommendation

## Contents **₽** ♥

- o.1 Global Func
- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg
    - ▼ 1.1.1 Model 2
      - 1.1.1.1 Mod
    - ▼ 1.1.2 Model 3
      - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
      - 1.1.3.1 Mod
      - 1.1.4 Best LR
  - ▼ 1.2 Random Fo
    - ▼ 1.2.1 Model 1
      - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
      - 1.2.2.1 Mod
  - 2 Data Interpretat
  - 3 Recommendation



<Figure size 432x288 with 0 Axes>

**Observations** Balancing the classes improved the model by increasing and increased the accuracy by 1%

## 1.1.4 Best LR Model

The besst Logistic Regression model was model #4 which produced the followin important features:

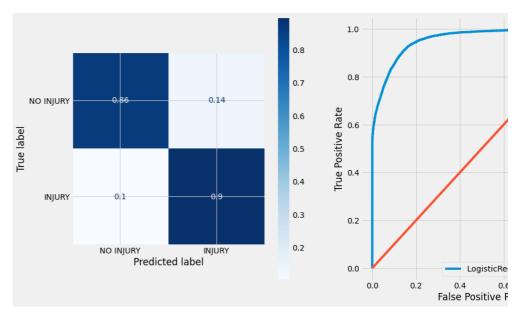
#### In [46]:

<b>v</b> 1	<pre>model_eval(lr_4_gridsearch.best_estimator_,</pre>
2	<pre>X_train_lr_sel, y_train, X_test_lr_sel,</pre>
3	y_test)

#### CURRENT MODEL: Not Overfit (Recall)

Classification Reports					
Classificacio	precision	recall	f1-score	support	
0	0.95	0.86	0.90	23838	
1	0.76	0.90	0.83	11648	
accuracy			0.88	35486	
macro avg	0.86	0.88	0.86	35486	
weighted avg	0.89	0.88	0.88	35486	

Test Graphs-----



#### In [52]:

#find top coefficients
best\_coeffs = get\_coefficients(lr\_4\_gridsearch.best\_estimate)
best\_coeffs\_pos = best\_coeffs.loc[best\_coeffs >= 0].sort\_vatable
best\_coeffs\_neg = best\_coeffs.loc[best\_coeffs < 0].sort\_vatable
best\_coeffs\_overall = best\_coeffs.abs().sort\_values(ascendate)</pre>

#### Contents 2 &

o.1 Global Func

▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2

1.1.1.1 Model

▼ 1.1.2 Model 3 1.1.2.1 Mod

▼ 1.1.3 Model 4 1.1.3.1 Mod 1.1.4 Best LR

▼ 1.2 Random Fo

▼ 1.2.1 Model 1 1.2.1.1 Mod

▼ 1.2.2 Model 2 1.2.2.1 Moc

2 Data Interpretat

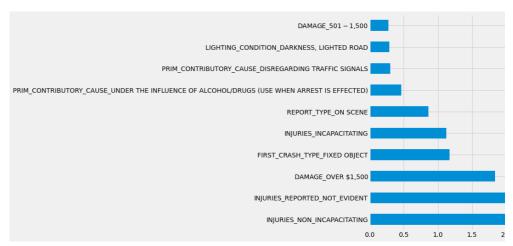
3 Recommendation

#### In [56]:

- 1 fig = plt.subplots(figsize=(8, 8))
- 2 best\_coeffs\_pos[:10].plot.barh()

#### Out[56]:

#### <AxesSubplot:>



#### 

0.1 Global Func▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2

1.1.1.1 Mod

▼ 1.1.2 Model 3

1.1.2.1 Mod

▼ 1.1.3 Model 4

1.1.3.1 Mod 1.1.4 Best LR

▼ 1.2 Random Fo

▼ 1.2.1 Model 1

1.2.1.1 Mod

▼ 1.2.2 Model 2

1.2.2.1 Mod

2 Data Interpretat

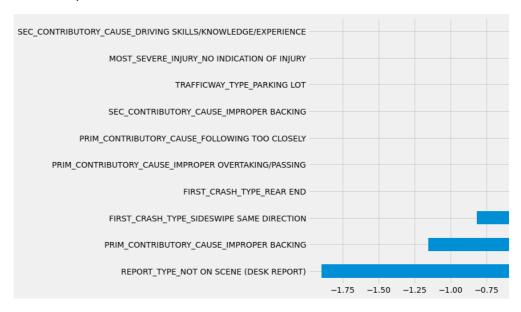
3 Recommendation

#### In [53]:

- 1 fig = plt.subplots(figsize=(8, 8))
- 2 best\_coeffs\_neg[:10].plot.barh()

#### Out[53]:

#### <AxesSubplot:>



## 1.2 Random Forest

## 1.2.1 Model 1



#### In [23]:

1 X\_train\_dt

#### Out[23]:

Conte	ents	$\mathcal{C}$	ť
0.1	Glob	al F	u

- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg
  - ▼ 1.1.1 Model 2
    - 1.1.1.1 Model 3 ▼ 1.1.2 Model 3
    - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
    - 1.1.3.1 Mod
    - 1.1.4 Best LR
  - ▼ 1.2 Random Fo
    - ▼ 1.2.1 Model 1
      - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
      - 1.2.2.1 Moc
  - 2 Data Interpretat
  - 3 Recommendatio

	LIGHTING_CONDITION_DARKNESS, LIGHTED ROAD	LIGHTING_CONDITION_DAWN	LIGHT
0	0.0	0.0	
1	1.0	0.0	
2	0.0	0.0	
3	0.0	0.0	
4	0.0	0.0	
82794	0.0	0.0	
82795	1.0	0.0	
82796	0.0	0.0	
82797	0.0	0.0	

1.0

82799 rows × 360 columns

## In [24]:

82798

```
1 rf_1 = RandomForestClassifier()
```

2 rf\_1.fit(X\_train\_dt, y\_train)

c:\python\python37\lib\site-packages\ipykernel\_launcher.py:2: Dawarning: A column-vector y was passed when a 1d array was expected hange the shape of y to (n\_samples,), for example using ravel().

#### Out[24]:

RandomForestClassifier()

#### 1.2.1.1 Model Evaluation

localhost:8888/notebooks/student-complete.ipynb#

0.0

## Contents 2 &

0.1 Global Func

- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg
    - ▼ 1.1.1 Model 2 1.1.1.1 Model
    - ▼ 1.1.2 Model 3
    - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
    - 1.1.3.1 Mod
    - 1.1.4 Best LR
  - ▼ 1.2 Random Fo
    - ▼ 1.2.1 Model 1
    - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
    - 1.2.2.1 Mod
  - 2 Data Interpretat
  - 3 Recommendation

#### In [25]:

<b>v</b> 1	<pre>model_eval(rf_1, X_train_dt, y_train, X_test_dt, y_test,</pre>
2	<pre>prev_model=lr_4_gridsearch.best_estimator_,</pre>
3	<pre>prev_X_train=X_train_lr_sel,</pre>
<b>▼</b> 4	<pre>prev_y_train=y_train,</pre>
5	<pre>prev_X_test=X_test_lr_sel, prev_y_test=y_test)</pre>

MODEL EVAL VS PREVIOUS (TEST)

\_\_\_\_\_

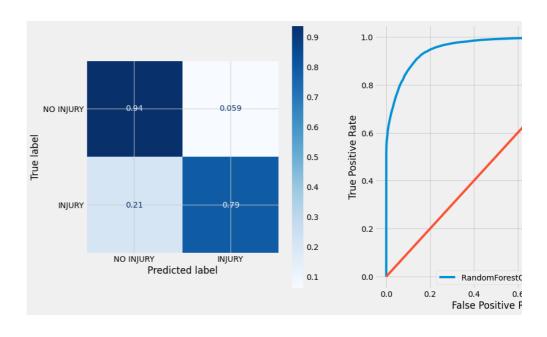
	Previous Model	<b>Current Model</b>	Delta
Recall	0.9	0.79	-0.11
F1	0.83	0.83	0
Accuracy	0.88	0.89	0.01
AUC	0.96	0.96	0

CURRENT MODEL: Overfit (Recall)

-----

Recall on Training: 1.0 Recall on Test: 0.79

Classification Reports						
Clussificacio	precision	recall	f1-score	support		
0	0.90	0.94	0.92	23838		
1	0.87	0.79	0.83	11648		
accuracy			0.89	35486		
macro avg	0.88	0.87	0.87	35486		
weighted avg	0.89	0.89	0.89	35486		



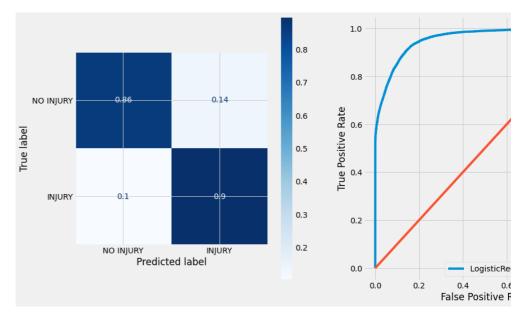
## Contents **₽** ♦

0.1 Global Func

- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg
    - ▼ 1.1.1 Model 2
      - 1.1.1.1 Mod
    - ▼ 1.1.2 Model 3
      - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
    - 1.1.3.1 Mod
    - 1.1.4 Best LR
  - ▼ 1.2 Random Fo
    - ▼ 1.2.1 Model 1
      - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
      - 1.2.2.1 Moc
  - 2 Data Interpretat
  - 3 Recommendation

#### PREVIOUS MODEL





<Figure size 432x288 with 0 Axes>

## 1.2.2 Model 2

#### In [26]:

1 rf\_2 = RandomForestClassifier()



#### Contents 2 &

```
0.1 Global Func
▼ 1 Complete jupyte
  ▼ 1.1 Logistic Reg
    ▼ 1.1.1 Model 2
        1.1.1.1 Mod
    ▼ 1.1.2 Model 3
        1.1.2.1 Mod
    ▼ 1.1.3 Model 4
        1.1.3.1 Mod
      1.1.4 Best LR
  ▼ 1.2 Random Fo
     ▼ 1.2.1 Model 1
        1.2.1.1 Mod
    ▼ 1.2.2 Model 2
        1.2.2.1 Moc
  2 Data Interpretat
```

3 Recommendation

```
In [27]:
```

#### In [28]:

```
1 rf_2_gridsearch = GridSearchCV(rf_2, params, scoring="reca
2 rf_2_gridsearch.fit(X_train_dt, y_train)
```

c:\python\python37\lib\site-packages\sklearn\model\_selection\\_sea
DataConversionWarning: A column-vector y was passed when a 1d ar
cted. Please change the shape of y to (n\_samples,), for example ().

```
self.best_estimator_.fit(X, y, **fit_params)
```

#### Out[28]:

#### 1.2.2.1 Model Evaluation

#### In [29]:

```
1 rf_2_gridsearch.best_params_
```

#### Out[29]:

```
{'class_weight': 'balanced',
  'criterion': 'entropy',
  'max_depth': 20,
  'max_features': 'auto',
  'min_samples_leaf': 3,
  'min_samples_split': 2,
  'n_estimators': 100}
```



#### In [30]:

▼ 1 2	mod	<pre>del_eval(rf_2_gridsearch.best_estimato</pre>	r_, X_train_dt, y_
	4		

MODEL EVAL VS PREVIOUS (TEST)

\_\_\_\_\_

	Previous Model	<b>Current Model</b>	Delta
Recall	0.79	0.9	0.11
F1	0.83	0.83	0
Accuracy	0.89	0.88	-0.01
AUC	0.96	0.96	0

CURRENT MODEL: Not Overfit (Recall)

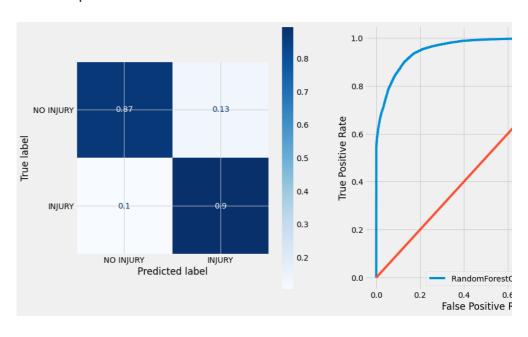
\_\_\_\_\_

Recall on Training: 0.91

Recall on Test: 0.9

Classification	n Reports				
	precision	recall	f1-score	support	
0	0.95	0.87	0.91	23838	
1	0.78	0.90	0.83	11648	
accuracy	2.25		0.88	35486	
macro avg weighted avg	0.86 0.89	0.89 0.88	0.87 0.88	35486 35486	

Test Graphs------



## Contents 2 ❖

o.1 Global Func

▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2 1.1.1.1 Model

▼ 1.1.2 Model 3 1.1.2.1 Mod ▼ 1.1.3 Model 4

▼ 1.1.3 Model 4 1.1.3.1 Mod 1.1.4 Best LR

▼ 1.2 Random Fo

▼ 1.2.1 Model 1 1.2.1.1 Mod ▼ 1.2.2 Model 2

1.2.2.1 Moc 2 Data Interpretat

3 Recommendation

#### Contents 2 &

0.1 Global Func

▼ 1 Complete jupyte

▼ 1.1 Logistic Reg

▼ 1.1.1 Model 2

1.1.1.1 Model 3

1.1.2.1 Model 3

1.1.2.1 Model 4

1.1.3.1 Model 4

1.1.3.1 Model 4

1.1.4 Best LR

▼ 1.2 Random Fo

▼ 1.2.1 Model 1

1.2.1.1 Model 1

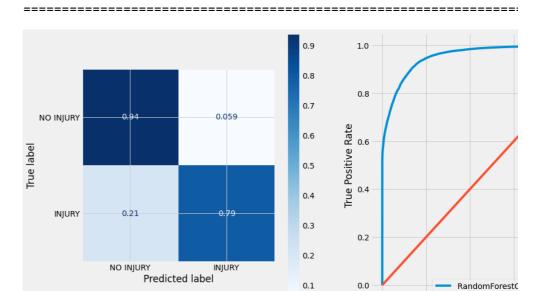
▼ 1.2.2 Model 2

2 Data Interpretat

3 Recommendation

1.2.2.1 Mod

#### PREVIOUS MODEL



0.2

False Positive F

<Figure size 432x288 with 0 Axes>

#### **Observations**

• The model increased recall by 12% but decreased accucacy by 1%

## In [33]:

1 rf\_2\_fi = rf\_2\_gridsearch.best\_estimator\_.feature\_importan

#### 

```
0.1 Global Func
▼ 1 Complete jupyte
  ▼ 1.1 Logistic Reg
    ▼ 1.1.1 Model 2
        1.1.1.1 Mod
    ▼ 1.1.2 Model 3
        1.1.2.1 Mod
    ▼ 1.1.3 Model 4
        1.1.3.1 Mod
      1.1.4 Best LR
  ▼ 1.2 Random Fo
     ▼ 1.2.1 Model 1
        1.2.1.1 Mod
    ▼ 1.2.2 Model 2
        1.2.2.1 Mod
  2 Data Interpretat
  3 Recommendation
```

## Out[42]:

#### Feature I

```
LIGHTING_CONDITION_DARKNESS, LIGHTED ROAD

LIGHTING_CONDITION_DAYLIGHT

SEC_CONTRIBUTORY_CAUSE_DISREGARDING TRAFFIC SIGNALS

LIGHTING_CONDITION_DAWN

LIGHTING_CONDITION_DUSK

SEC_CONTRIBUTORY_CAUSE_DISREGARDING STOP SIGN

SEC_CONTRIBUTORY_CAUSE_DISTRACTION - FROM INSIDE VEHICLE

SEC_CONTRIBUTORY_CAUSE_DISREGARDING OTHER TRAFFIC SIGNS

SEC_CONTRIBUTORY_CAUSE_DISREGARDING ROAD MARKINGS

SEC_CONTRIBUTORY_CAUSE_DISREGARDING ROAD MARKINGS
```

```
In [37]:
    1 rf_2_features_sorted
Out[37]:
```

#### Feature I

LIGHTING\_CONDITION\_DARKNESS, LIGHTED ROAD

LIGHTING\_CONDITION\_DAYLIGHT

SEC\_CONTRIBUTORY\_CAUSE\_DISREGARDING TRAFFIC SIGNALS

LIGHTING\_CONDITION\_DAWN

LIGHTING\_CONDITION\_DUSK

SEC\_CONTRIBUTORY\_CAUSE\_DISREGARDING STOP SIGN

SEC\_CONTRIBUTORY\_CAUSE\_DISTRACTION - FROM INSIDE VEHICLE

SEC\_CONTRIBUTORY\_CAUSE\_DISREGARDING OTHER TRAFFIC SIGNS

SEC\_CONTRIBUTORY\_CAUSE\_DISREGARDING ROAD MARKINGS

SEC\_CONTRIBUTORY\_CAUSE\_CELL PHONE USE OTHER THAN TEXTING

## Contents 2 &

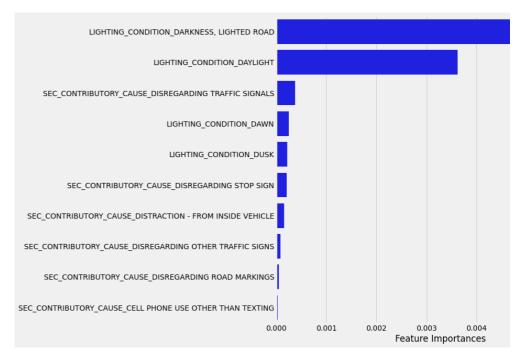
o.1 Global Func ▼ 1 Complete jupyte ▼ 1.1 Logistic Reg ▼ 1.1.1 Model 2

- ▼ 1.1.1 Model 2 1.1.1.1 Model
- ▼ 1.1.2 Model 3 1.1.2.1 Mod
- ▼ 1.1.3 Model 4 1.1.3.1 Mod 1.1.4 Best LR
- ▼ 1.2 Random Fo
  - ▼ 1.2.1 Model 1 1.2.1.1 Mod
  - ▼ 1.2.2 Model 2
- 1.2.2.1 Moc 2 Data Interpretat
- 3 Recommendation

#### In [39]:

#### Out[39]:

<AxesSubplot:xlabel='Feature Importances'>



# 2 Data Interpretation

This dataset took a lot of cleaning, especially involving unknown values original dataset. After cleaning, a logistic regression and random forest trained to predict civilians in Chicago being injured or not. Classification and optimized for finding an injured civilian due to car crash.

#### In [ ]:

1



4

## Contents **₽** ♥

- 0.1 Global Func
- ▼ 1 Complete jupyte
  - ▼ 1.1 Logistic Reg
    - ▼ 1.1.1 Model 2
      - 1.1.1.1 Mod
    - ▼ 1.1.2 Model 3
      - 1.1.2.1 Mod
    - ▼ 1.1.3 Model 4
      - 1.1.3.1 Mod
      - 1.1.4 Best LR
  - ▼ 1.2 Random Fo
    - ▼ 1.2.1 Model 1
      - 1.2.1.1 Mod
    - ▼ 1.2.2 Model 2
      - 1.2.2.1 Mod
  - 2 Data Interpretat
  - 3 Recommendation