

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.ensemble import GradientBoostingClassifier
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.neural_network import MLPClassifier
11 from xgboost import XGBClassifier
12 from lightgbm import LGBMClassifier
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call dr

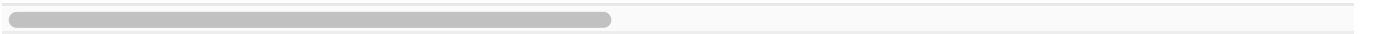


```
1 # reading the data
2
3 train_df = pd.read_csv("/content/drive/MyDrive/CSV files/Hackathon- Mac1
4 test_df = pd.read_csv("/content/drive/MyDrive/CSV files/Hackathon- Mach:
5 print(train_df.shape)
6 print(test_df.shape)
```

```
(51490, 42)
```

```
(77235, 42)
```

```
/usr/local/lib/python3.7/dist-packages/IPython/core/interactiveshell.py:2718: Dt
interactivity=interactivity, compiler=compiler, result=result)
```



```
1 # train_df = pd.read_csv('/content/train.csv')
2 # test_df = pd.read_csv('/content/test.csv')
```

```
1 train_df.head()
```

	Report Number	Local Case Number	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name	Cross-Street Type
0	MP060D	10196	Montgomery County Police	Injury Crash	06/05/2017 04:27:00 PM	Maryland (State)	OLD HUNDRED RD	County
1	MP000X	20080	Montgomery County Police	Injury Crash	02/22/2020 10:00:00 AM	County	JANET RD	County
2	MP070N	10103	Montgomery County Police	Property Damage Crash	06/11/2017 08:21:00 AM	Maryland (State)	COLESVILLE RD	Maryland (State)

```
1 test_df.head()
```

	Id	Report Number	Local Case Number	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name
0	0	MP810M	NaN	MONTGOMERY	Property Damage Crash	12/16/2015 05:42:00 PM	County	WHITTIER BLVD
1	1	MP2001	NaN	Montgomery County Police	Property Damage Crash	05/01/2016 07:25:00 PM	County	RAILROAD ST
2	2	MP6404	10125	Montgomery County Police	Property Damage Crash	03/16/2018 02:30:00 PM	Maryland (State)	GEORGIA AVE
3	3	MP0405	10270	Montgomery County Police	Property Damage Crash	05/25/2018 05:00:00 PM	NaN	NaN
4	4	MP090D	10374	Montgomery County Police	Property Damage Crash	09/23/2017 05:40:00 PM	Interstate (State)	EISENHOWER MEMORIAL HWY

```
1 # column location can be split with delimiter as space or comma
2
3 new = train_df['Location'].str.split(",| ", expand = True)
4 train_df['Location_0'] = new[0].astype(float)
5 train_df['Location_1'] = new[1].astype(float)
```

```
6 train_df.drop(['Location'] , axis = 'columns', inplace = True)
7 train_df.head()
```

	Report Number	Local Case Number	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name	Crash State
0	MP060D	10196	Montgomery County Police	Injury Crash	06/05/2017 04:27:00 PM	Maryland (State)	OLD HUNDRED RD	Co
1	MP000X	20080	Montgomery County Police	Injury Crash	02/22/2020 10:00:00 AM	County	JANET RD	Co
2	MP070N	10103	Montgomery County Police	Property Damage Crash	06/11/2017 08:21:00 AM	Maryland (State)	COLESVILLE RD	Mar (S
3	MP130Y	10208	Montgomery County Police	Injury Crash	05/10/2019 07:30:00 AM	County	PARKLAND DR	Co
4	MP770Y	1091	Montgomery County Police	Injury Crash	08/02/2016 05:00:00 PM	Maryland (State)	OLD GEORGETOWN RD	Co

```
1 # Repeating same steps on test data
2
3 new = test_df['Location'].str.split(",| " , expand = True)
4 test_df['Location_0'] = new[0].astype(float)
5 test_df['Location_1'] = new[1].astype(float)
6 test_df.drop(['Location'] , axis = 'columns', inplace = True)
7 test_df.head()
```

	Id	Report Number	Local Case Number	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name
0	0	MP810M	NaN	MONTGOMERY	Property Damage Crash	12/16/2015 05:42:00 PM	County	WHITTIE BLV
1	1	MP2001	NaN	Montgomery County Police	Property Damage Crash	05/01/2016 07:25:00 PM	County	RAILROAD S
2	2	MP6404	10125	Montgomery County Police	Property Damage	03/16/2018 02:30:00 PM	Maryland (State)	GEORG AV

```

1 # We assert that both Location_0 and Location_1 coordinates are between
2
3 for index,row in train_df.iterrows():
4     x = row['Location_0']
5     y = row['Location_1']
6     while(x < -10 or x > 10):
7         x /= 10
8     train_df.at[index , 'Location_0'] = x
9
10    while(y < -10 or y > 10):
11        y /= 10
12    train_df.at[index , 'Location_1'] = y
13
14 for index,row in test_df.iterrows():
15     x = row['Location_0']
16     y = row['Location_1']
17     while(x < -10 or x > 10):
18         x /= 10
19     test_df.at[index , 'Location_0'] = x
20
21    while(y < -10 or y > 10):
22        y /= 10
23    test_df.at[index , 'Location_1'] = y
24
25
1 train_df.head()

```

	Report Number	Local Case Number	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name	Crash Status
0	MP060D	10196	Montgomery County Police	Injury Crash	06/05/2017 04:27:00 PM	Maryland (State)	OLD HUNDRED RD	Completed
1	MP000X	20080	Montgomery County Police	Injury Crash	02/22/2020 10:00:00 AM	County	JANET RD	Completed
2	MP070N	10103	Montgomery County Police	Property Damage Crash	06/11/2017 08:21:00 AM	Maryland (State)	COLESVILLE RD	Married (S)
3	MP130Y	10208	Montgomery County Police	Injury Crash	05/10/2019 07:30:00 AM	County	PARKLAND DR	Completed

```

1 # We drop columns Report Number, Local Case Number, Person Id from both
2 # We store column ID in test data for further use and then drop it
3
4 indices = test_df['Id']
5
6 train_df.drop(['Report Number', 'Local Case Number', 'Vehicle ID'], axis=1)
7 test_df.drop(['Report Number', 'Local Case Number', 'Id', 'Vehicle ID'], axis=1)

1 # Since most columns are categorical, the idea is to drop columns with 1
2 # But since many of the columns are having NaNs, replace them with most
3 # Not before dropping columns with more than 20000 Nans
4 pd.isna(train_df).sum()
5
6 for column in train_df.columns:
7     if (pd.isna(train_df).sum()[column] > 20000):
8         train_df.drop([column], axis=1, inplace=True)
9
10 train_df = train_df.apply(lambda x: x.fillna(x.value_counts().index[0]))

1 # retrieving Fault column and storing it as target variable
2 y_train = np.array(train_df['Fault'])
3
4 train_df.drop(['Fault'], axis=1, inplace=True)

1 # rearranging the columns of test_df in the same order as train_df
2 test_df = test_df[train_df.columns]
```

```

3
4 # filling in the NaNs of a column of test_df with most repeated values
5 # same column of train_df
6 for column in test_df.columns:
7     test_df[column].fillna(train_df[column].mode()[0] , inplace = True)
8

```

/usr/local/lib/python3.7/dist-packages/pandas/core/series.py:4536: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/10min/10min_tips.html#downcast=downcast,

```
1 train_df.head()
```

	Agency Name	ACRS Report Type	Crash Date/Time	Route Type	Road Name	Cross-Street Type	Cross-Street Name
0	Montgomery County Police	Injury Crash	06/05/2017 04:27:00 PM	Maryland (State)	OLD HUNDRED RD	County	THURSTON RD
1	Montgomery County Police	Injury Crash	02/22/2020 10:00:00 AM	County	JANET RD	County	FLACK ST
2	Montgomery County Police	Property Damage Crash	06/11/2017 08:21:00 AM	Maryland (State)	COLESVILLE RD	Maryland (State)	STRUC #15082
3	Montgomery County Police	Injury Crash	05/10/2019 07:30:00 AM	County	PARKLAND DR	County	FRANKFORT DR
4	Montgomery County Police	Injury Crash	08/02/2016 05:00:00 PM	Maryland (State)	OLD GEORGETOWN RD	County	MCKINLEY ST

```

1 # we change the crash date/time column into a column which tells us if
2 # the crash was during hours AM or PM
3
4 # extracting AM and PM from Crash Date/Time column
5
6 new = train_df['Crash Date/Time'].str.split(" ", expand = True)
7 train_df['time'] = new[2]
8
9 # drop Crash Date/Time column

```

```

10 train_df.drop(['Crash Date/Time'] , axis = 'columns' , inplace = True)
11
12 # considering more priority to PM
13 train_df['time'] = train_df['time'].replace(['AM' , 'PM'] , [0,1])
14
15
16

```

```

1 # repeating the same steps on Crash Date/Time on test_df
2
3 # we change the crash date/time column into a column which tells us if
4 # the crash was during hours AM or PM
5
6 # extracting AM and PM from Crash Date/Time column
7
8 new = test_df['Crash Date/Time'].str.split(" " , expand = True)
9 test_df['time'] = new[2]
10
11 # drop Crash Date/Time column
12 test_df.drop(['Crash Date/Time'] , axis = 'columns' , inplace = True)
13
14 # considering more priority to PM
15 test_df['time'] = test_df['time'].replace(['AM' , 'PM'] , [0,1])
16
17
18

```

```

1 # finding the number of unique values per column with dtype==object
2
3 for column in train_df.columns:
4     if train_df[column].dtypes == object:
5         print(column , len(pd.unique(train_df[column])))

```

```

Agency Name 10
ACRS Report Type 3
Route Type 10
Road Name 2226
Cross-Street Type 10
Cross-Street Name 4525
Collision Type 18
Weather 12
Surface Condition 11
Light 8
Traffic Control 11
Driver Substance Abuse 11
Person ID 49209
Injury Severity 5
Drivers License State 65

```

```

Vehicle Damage Extent 7
Vehicle First Impact Location 16
Vehicle Second Impact Location 16
Vehicle Body Type 30
Vehicle Movement 22
Vehicle Continuing Dir 5
Vehicle Going Dir 5
Driverless Vehicle 2
Parked Vehicle 2
Vehicle Make 886
Vehicle Model 3128
Equipment Problems 8

```

```

1 # we drop columns with more than 20 unique values
2
3 for column in train_df.columns:
4     if len(pd.unique(train_df[column])) > 50:
5         train_df.drop(column, axis = 'columns' , inplace = True)

1 # keeping same columns in test_df
2 test_df = test_df[train_df.columns]
3
4 # comparing unique values per column in test and train data:
5 for column in test_df.columns:
6     if test_df[column].dtypes == object:
7         print(column , ", train= " ,len(pd.unique(train_df[column]))

Agency Name , train= 10 , test = 10
ACRS Report Type , train= 3 , test = 3
Route Type , train= 10 , test = 10
Cross-Street Type , train= 10 , test = 10
Collision Type , train= 18 , test = 18
Weather , train= 12 , test = 12
Surface Condition , train= 11 , test = 11
Light , train= 8 , test = 8
Traffic Control , train= 11 , test = 11
Driver Substance Abuse , train= 11 , test = 11
Injury Severity , train= 5 , test = 5
Vehicle Damage Extent , train= 7 , test = 7
Vehicle First Impact Location , train= 16 , test = 16
Vehicle Second Impact Location , train= 16 , test = 16
Vehicle Body Type , train= 30 , test = 30
Vehicle Movement , train= 22 , test = 22
Vehicle Continuing Dir , train= 5 , test = 5
Vehicle Going Dir , train= 5 , test = 5
Driverless Vehicle , train= 2 , test = 2
Parked Vehicle , train= 2 , test = 2
Equipment Problems , train= 8 , test = 10

```

```

1 # dropping Equipment problems because testing data has more unique values
2 # training data

```



```
3
4 train_df.drop(['Equipment Problems'] , axis = 'columns' , inplace = True)
5 test_df.drop(['Equipment Problems'] , axis = 'columns' , inplace = True)

1 # one hot encoding
2
3 train_df = pd.get_dummies(train_df)
4 test_df = pd.get_dummies(test_df)
5
6 print(train_df.shape, test_df.shape)

(51490, 216) (77235, 216)

1 X_train = np.array(train_df)
2 X_test = np.array(test_df)

1 model = KNeighborsClassifier(n_neighbors=7)
2 model.fit(X_train,y_train)
3 y_pred = model.predict(X_test)
4
5
6 predictions = []
7 for i in range(len(indices)):
8     predictions.append([indices[i] , y_pred[i]])
9
10 predictions = np.array(predictions)
11
12 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
13 pred_df
14
15 pred_df.to_csv('KNN.csv' , columns = ['Id' , 'Fault'] , header = ['Id' ,
16

1 model = RandomForestClassifier(n_estimators=1000)
2 model.fit(X_train, y_train)
3 y_pred = model.predict(X_test)
4
5 predictions = []
6 for i in range(len(indices)):
7     predictions.append([indices[i] , y_pred[i]])
8
9 predictions = np.array(predictions)
10
```

```
11 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
12 pred_df
13
14 pred_df.to_csv('RFC_1000.csv' , columns = ['Id' , 'Fault'] , header = [

1  model = GradientBoostingClassifier(n_estimators=2500)
2  model.fit(X_train,y_train)
3  y_pred = model.predict(X_test)
4
5  predictions = []
6  for i in range(len(indices)):
7      predictions.append([indices[i] , y_pred[i]])
8
9  predictions = np.array(predictions)
10
11 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
12 pred_df
13
14 pred_df.to_csv('GBC_2500.csv' , columns = ['Id' , 'Fault'] , header =

1 model =  MLPClassifier( max_iter=700)
2 model.fit(X_train,y_train)
3 y_pred = model.predict(X_test)
4
5 predictions = []
6 for i in range(len(indices)):
7     predictions.append([indices[i] , y_pred[i]])
8
9 predictions = np.array(predictions)
10
11 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
12 pred_df
13
14 pred_df.to_csv('MLP_700.csv' , columns = ['Id' , 'Fault'] , header = ['

1 model = XGBClassifier(max_depth = 200, subsample = 0.75)
2 model.fit(X_train, y_train)
3 y_pred = model.predict(X_test)
4
5 predictions = []
6 for i in range(len(indices)):
7     predictions.append([indices[i] , y_pred[i]])
8
9 predictions = np.array(predictions)
10
```

```
11 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
12 pred_df
13
14 pred_df.to_csv('XGB_2.csv' , columns = ['Id' , 'Fault'] , header = ['Id

1 model = LGBMClassifier(n_estimators = 1000)
2 model.fit(X_train, y_train)
3 y_pred = model.predict(X_test)
4
5 predictions = []
6 for i in range(len(indices)):
7     predictions.append([indices[i] , y_pred[i]])
8
9 predictions = np.array(predictions)
10
11 pred_df = pd.DataFrame(predictions , columns = ['Id' , 'Fault'])
12 pred_df
13
14 pred_df.to_csv('LGBM_1000.csv' , columns = ['Id' , 'Fault'] , header = |
```

Conclusion

The final models we used was a Gradient Boosting Classifier with the number of estimators as 2500, and a Light Gradient Boosting Model with n_estimators = 1000.

The reason why we chose Gradient Boosting Classifier is because it combines various weak learners, by improving on the mistakes made by the previous models, and combines the results of all models. It is robust to over-fitting. It is a greedy algorithm, choosing trees based on best split, and the weak learners can be constrained, to ensure they remain weak, but greedy.

Light Gradient Boosting Model is an algorithm, that significantly cuts on memory, and is extremely fast-paced, with note-worthy accuracies. It splits trees leaf-wise, rather than level-wise , and hence, obtains better accuracies compared to many of the other models.

The reason why we chose these models is because of the high number of features in our data due to the one hot encoding done, and an iterative approach to finding the pattern seemed like the best idea.

✓ 42s completed at 1:03 PM

● ✕