```python
In [1]:    import numpy as np
           import pandas as pd
           import matplotlib.pyplot as plt
           import seaborn as sns
```

```python
In [18]:   df = pd.read_csv("last_two_years_accidents.csv")
```

```python
In [19]:   df.columns
```

```
Out[19]:   Index(['ID', 'Source', 'Severity', 'Start_Time', 'End_Time', 'Start_Lat',
                  'Start_Lng', 'Distance(mi)', 'Street', 'City', 'County', 'State',
                  'Zipcode', 'Country', 'Timezone', 'Airport_Code', 'Weather_Timestamp',
                  'Temperature(F)', 'Humidity(%)', 'Pressure(in)', 'Visibility(mi)',
                  'Wind_Direction', 'Wind_Speed(mph)', 'Weather_Condition', 'Amenity',
                  'Bump', 'Crossing', 'Give_Way', 'Junction', 'No_Exit', 'Railway',
                  'Roundabout', 'Station', 'Stop', 'Traffic_Calming', 'Traffic_Signal',
                  'Turning_Loop', 'Sunrise_Sunset', 'Civil_Twilight', 'Nautical_Twilight',
                  'Astronomical_Twilight'],
                 dtype='object')
```

```python
In [20]:   severity_labels = {
               1: 'Not Severe',
               2: 'Not Severe',
               3: 'Severe',
               4: 'Severe'
           }
```

```python
In [21]:   df['Severity'] = df['Severity'].map(severity_labels)
```

```python
In [22]:   size = df['Severity'].value_counts()['Severe']
```

```python
In [23]:   df_balanced_severity = pd.DataFrame()
```

```python
In [24]:   df_balanced_severity = df.groupby('Severity', group_keys = False).apply(lambda x: x.sample(size, random_state = 30))
```

```python
In [25]:   df_balanced_severity['Severity'].value_counts()
```

```
Out[25]:   Severity
           Not Severe    128292
           Severe        128292
           Name: count, dtype: int64
```

```python
In [26]:   categorical_features = set(["Weather_Condition", "Civil_Twilight", 'Wind_Direction'])
```

```python
In [27]:   for feature in categorical_features:
               df_balanced_severity[feature] = df_balanced_severity[feature].astype("category")
```

```python
In [28]:   bool_columns = df_balanced_severity.select_dtypes(include='bool').columns
```

```python
In [29]:   df_balanced_severity[bool_columns] = df_balanced_severity[bool_columns].replace({True:1, False:0})
```

```python
In [30]:   df2= df_balanced_severity[['Start_Lat','Start_Lng','Distance(mi)', 'Temperature(F)', 'Humidity(%)', 'Pressure(in)',
                   'Visibility(mi)', 'Wind_Speed(mph)','Amenity','Bump','Crossing','Give_Way',
                   'Junction','No_Exit','Railway','Roundabout','Station','Stop','Traffic_Calming','Traffic_Signal',
                   'Civil_Twilight','Weather_Condition','Civil_Twilight',
                   'Wind_Direction','Severity']]]
```

```python
In [31]:   df2 = pd.get_dummies(df2, columns=list(categorical_features), drop_first=True)
```

```python
In [32]:   from sklearn.metrics import classification_report
           from sklearn.model_selection import train_test_split
           from sklearn.tree import DecisionTreeClassifier, plot_tree
           from sklearn.metrics import accuracy_score, confusion_matrix
```

```python
In [33]:   Y = df2['Severity'] # target column
           X = df2.drop(columns = ['Severity']) # features
```

```python
In [34]:   X_train, X_test, y_train, y_test = train_test_split(X, Y,test_size=0.2, random_state=30)
```

```
In [35]:    DTree1 = DecisionTreeClassifier()
            DTree1.fit(X_train, y_train)
            y_test_pred = DTree1.predict(X_test)
            y_train_pred = DTree1.predict(X_train)


            acc = accuracy_score(y_test, y_test_pred)
            CM = confusion_matrix(y_test,y_test_pred)
            print(f"Accuracy of the DTree : {np.round(acc,3)*100}")
            print(f"Confusion Matrix: \n {CM}")

            Accuracy of the DTree : 77.9
            Confusion Matrix:
             [[20010  5680]
             [ 5646 19981]]

In [36]:    DTree1.get_params()

Out[36]:    {'ccp_alpha': 0.0,
             'class_weight': None,
             'criterion': 'gini',
             'max_depth': None,
             'max_features': None,
             'max_leaf_nodes': None,
             'min_impurity_decrease': 0.0,
             'min_samples_leaf': 1,
             'min_samples_split': 2,
             'min_weight_fraction_leaf': 0.0,
             'monotonic_cst': None,
             'random_state': None,
             'splitter': 'best'}

In [37]:    print(classification_report(y_train, y_train_pred))

                          precision    recall  f1-score   support

              Not Severe       0.98      1.00      0.99    102602
                  Severe       1.00      0.98      0.99    102665

                accuracy                           0.99    205267
               macro avg       0.99      0.99      0.99    205267
            weighted avg       0.99      0.99      0.99    205267

In [39]:    print(classification_report(y_test, y_test_pred))

                          precision    recall  f1-score   support

              Not Severe       0.78      0.78      0.78     25690
                  Severe       0.78      0.78      0.78     25627

                accuracy                           0.78     51317
               macro avg       0.78      0.78      0.78     51317
            weighted avg       0.78      0.78      0.78     51317
```
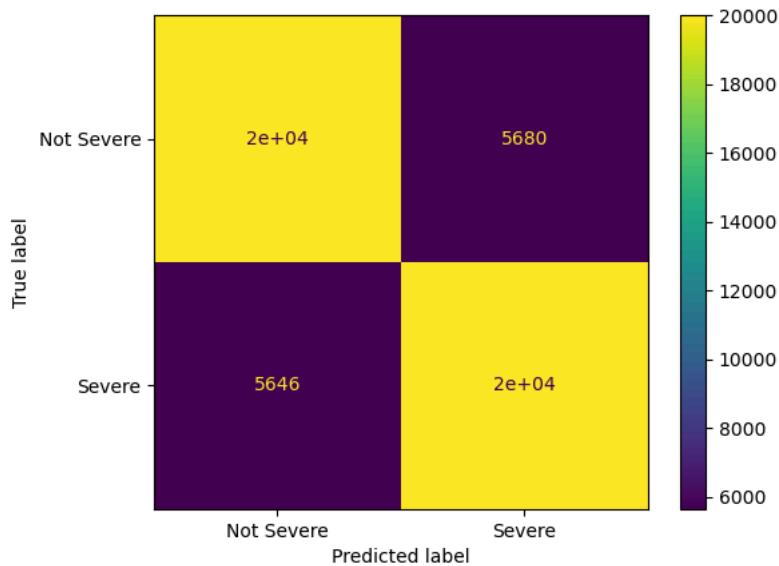
```
In [40]: ▶| from sklearn.metrics import ConfusionMatrixDisplay
             cm = confusion_matrix(y_test, y_test_pred, labels=DTree1.classes_)
             disp = ConfusionMatrixDisplay(confusion_matrix=cm,
                                           display_labels=DTree1.classes_)
             disp.plot()
             plt.show()
```
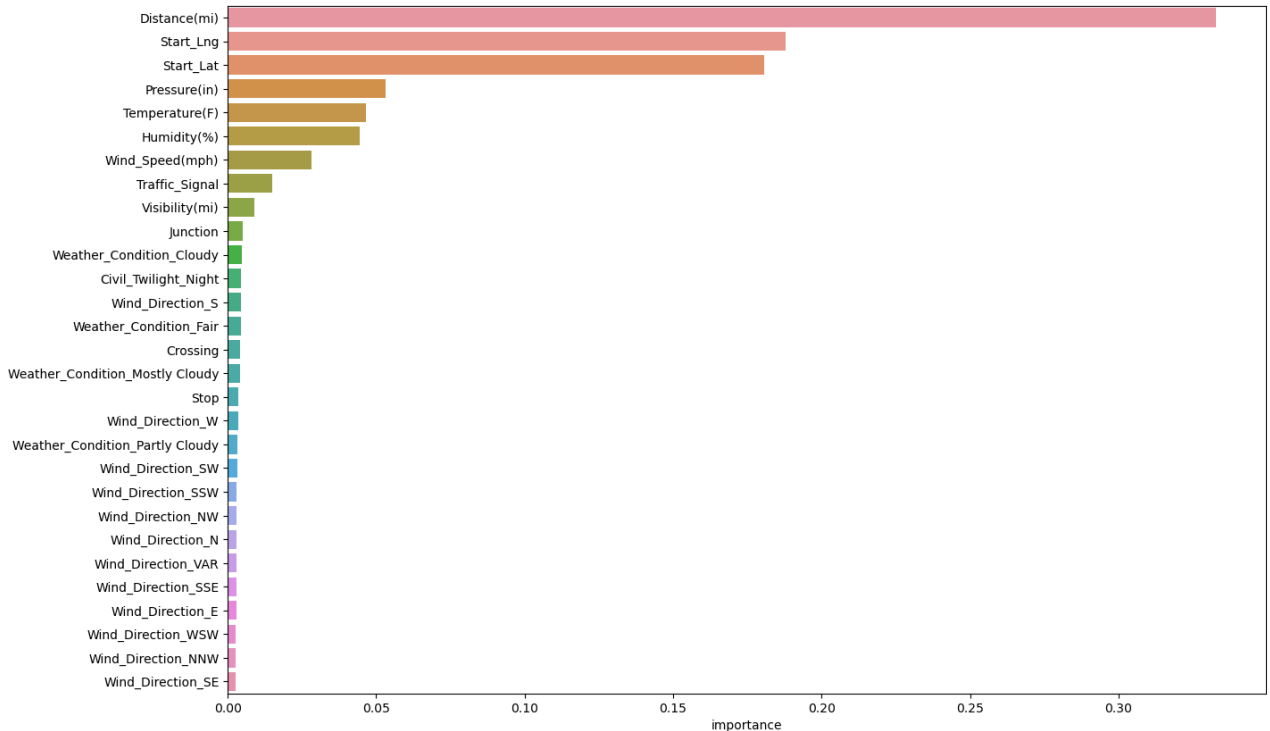


```
In [42]: ▶| importances = pd.DataFrame(np.zeros((X_train.shape[1], 1)), columns=["importance"], index=X_train.columns)
```
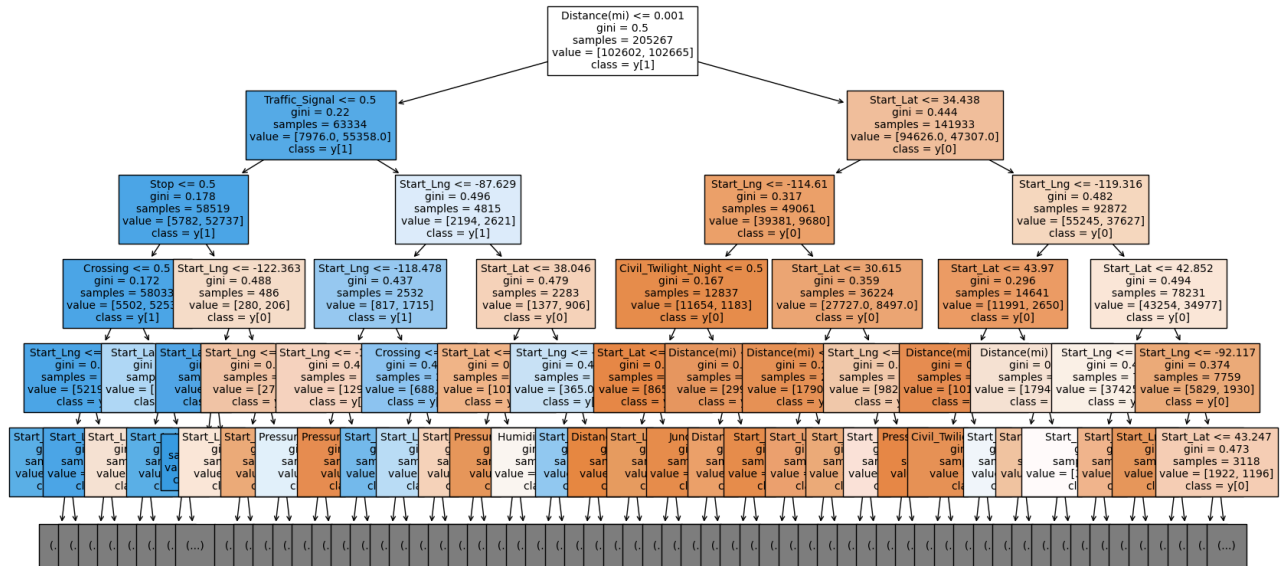
```
In [43]: ▶| importances.iloc[:,0] = DTree1.feature_importances_
```

```
In [44]: ▶| importances = importances.sort_values(by="importance", ascending=False)[:30]
```

```
In [45]: ▶| plt.figure(figsize=(15, 10))
             sns.barplot(x="importance", y=importances.index, data=importances)
             plt.show()
```

```python
In [50]:  fig, ax = plt.subplots(figsize=(20, 10))
          plot_tree(DTree1, max_depth=5, fontsize=10, feature_names=X_train.columns.to_list(), class_names = True, filled=True)
          plt.show()
```



```python
In [51]:  dtc1 = DecisionTreeClassifier(random_state=42, criterion = 'entropy', max_depth = 5, min_samples_leaf=10000, min_samples_spli
```

```python
In [52]:  print("Default scores:")
          dtc1.fit(X_train, y_train)
          print("Train score:", dtc1.score(X_train, y_train))
          print("Validation score:", dtc1.score(X_test, y_test))

          Default scores:
          Train score: 0.7306776052653373
          Validation score: 0.732778611376347
```

```python
In [53]:  y_pred = dtc1.predict(X_test)
          y_pred_train=dtc1.predict(X_train)
```

```python
In [54]:  accuracy_train= accuracy_score(y_train, y_pred_train)
          accuracy = accuracy_score(y_test, y_pred)
          print("Accuracy: {:.2f}%".format(accuracy * 100))
          print("Accuracy_Train: {:.2f}%".format(accuracy_train * 100))

          Accuracy: 73.28%
          Accuracy_Train: 73.07%
```

```python
In [55]:  y_pred = dtc1.predict(X_test)
```

```python
In [56]:  accuracy_score(y_test, y_pred)
```
```
Out[56]:  0.732778611376347
```

```python
In [58]:  print(classification_report(y_test, y_pred))

                        precision    recall  f1-score   support

            Not Severe       0.67      0.93      0.78     25690
                Severe       0.88      0.54      0.67     25627

              accuracy                           0.73     51317
             macro avg       0.77      0.73      0.72     51317
          weighted avg       0.77      0.73      0.72     51317
```

```python
In [59]:  confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
```
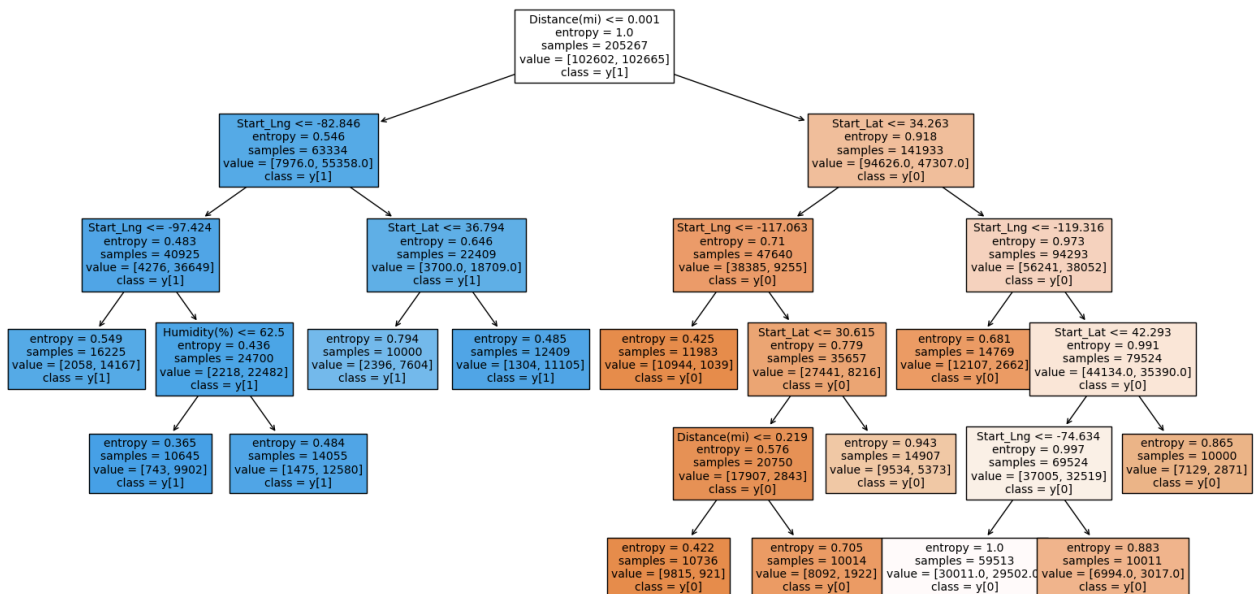
```python
In [60]:  confmat
```
```
Out[60]:  array([[23779,  1911],
                 [11802, 13825]], dtype=int64)
```

```
In [61]:  from sklearn.metrics import ConfusionMatrixDisplay
          cm = confusion_matrix(y_test, y_test_pred, labels=dtc1.classes_)
          disp = ConfusionMatrixDisplay(confusion_matrix=confmat,
                                         display_labels=dtc1.classes_)
          disp.plot()
          plt.show()
```



```
In [62]:  fig, ax = plt.subplots(figsize=(20, 10))
          plot_tree(dtc1, max_depth=5, fontsize=10, feature_names=X_train.columns.to_list(), class_names = True, filled=True)
          plt.show()
```



```
In [63]:  df3= df2.drop('Distance(mi)',axis = 1)
```

```
In [64]:  Y = df3['Severity'] # target column
          X = df3.drop(columns = ['Severity']) # features
```

```
In [65]:  X_train, X_test, y_train, y_test = train_test_split(X, Y,test_size=0.2, random_state=30)
```

```
In [67]:  dtc2 = DecisionTreeClassifier(random_state=42, criterion = 'gini', max_depth = 5,  min_samples_leaf=500, min_samples_split= 1
```

```python
In [68]:  ▶  print("Default scores:")
             dtc2.fit(X_train, y_train)
             print("Train score:", dtc2.score(X_train, y_train))
             print("Validation score:", dtc2.score(X_test, y_test))
```

```
Default scores:
Train score: 0.6305056341253099
Validation score: 0.6334937739930238
```

```python
In [69]:  ▶  y_pred = dtc2.predict(X_test)

             accuracy_score(y_test, y_pred)
```

Out[69]: 0.6334937739930238

```python
In [70]:  ▶  print(classification_report(y_test, y_pred))
```

```
                precision    recall  f1-score   support

  Not Severe       0.68      0.50      0.58     25690
      Severe       0.61      0.76      0.68     25627

    accuracy                           0.63     51317
   macro avg       0.64      0.63      0.63     51317
weighted avg       0.64      0.63      0.63     51317
```

```python
In [71]:  ▶  y_pred = dtc2.predict(X_test)
             confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
```

```python
In [72]:  ▶  confmat
```

Out[72]: array([[12935, 12755],
                [ 6053, 19574]], dtype=int64)

```python
In [73]:  ▶  from sklearn.metrics import ConfusionMatrixDisplay
             cm = confusion_matrix(y_test, y_test_pred, labels=dtc2.classes_)
             disp = ConfusionMatrixDisplay(confusion_matrix=confmat,
                                           display_labels=dtc2.classes_)
             disp.plot()
             plt.show()
```

```
In [74]:  ▶ | fig, ax = plt.subplots(figsize=(20, 10))
             plot_tree(dtc2, max_depth=4, fontsize=10, feature_names=X_train.columns.to_list(), class_names = True, filled=True)
             plt.show()
```

Start_Lat <= 28.641
gini = 0.5
samples = 205267
value = [102602, 102665]
class = y[1]

Start_Lng <= -80.211
gini = 0.374
samples = 16761
value = [12586, 4175]
class = y[0]

Start_Lng <= -108.342
gini = 0.499
samples = 188506
value = [90016.0, 98490.0]
class = y[1]

Start_Lat <= 25.946
gini = 0.322
samples = 13079
value = [10441, 2638]
class = y[0]

gini = 0.486
samples = 3682
value = [2145, 1537]
class = y[0]

Start_Lat <= 47.032
gini = 0.467
samples = 49553
value = [31111.0, 18442.0]
class = y[0]

Start_Lat <= 42.906
gini = 0.488
samples = 138953
value = [58905, 80048]
class = y[1]

gini = 0.082
samples = 3987
value = [3816, 171]
class = y[0]

gini = 0.395
samples = 9092
value = [6625, 2467]
class = y[0]

Start_Lng <= -119.316
gini = 0.461
samples = 47235
value = [30196.0, 17039.0]
class = y[0]

gini = 0.478
samples = 2318
value = [915, 1403]
class = y[1]

Pressure(in) <= 29.545
gini = 0.484
samples = 132211
value = [54290, 77921]
class = y[1]

gini = 0.432
samples = 6742
value = [4615, 2127]
class = y[0]

Start_Lng <= -121.5
gini = 0.411
samples = 17131
value = [12174, 4957]
class = y[0]

Start_Lng <= -119.287
gini = 0.481
samples = 30104
value = [18022, 12082]
class = y[0]

Crossing <= 0.5
gini = 0.458
samples = 70835
value = [25143, 45692]
class = y[1]

Crossing <= 0.5
gini = 0.499
samples = 61376
value = [29147, 32229]
class = y[1]

(...)    (...)    (...)    (...)    (...)    (...)    (...)    (...)

```
In [ ]:  ▶ |
```