**Data Analysis Procedure**

The data was first prepared by transforming the target feature **ACCLASS** into two categories - **Non-Fatal** and **Fatal.** The **DATE** feature was also transformed into datetime datatype to get 5 new features namely; **MONTH, MONTH_NAME, MINUTES, DAY,** and **WEEKDAY.**



```
Data Preparation
  1  def wrangle(path):
  2      # Read csv file into dataframe
  3      df = pd.read_csv(path)
  4
  5      # Select the features
  6      # df = df[['YEAR', 'DATE', 'TIME', 'HOUR', 'STREET1', 'STREET2', 'ROAD_CLASS', 'DISTRICT', 'LOCCOORD', 'ACCLOC', 'TRAFFCTL', 'VISIBILITY', 'LIGHT',
  7      # 'RDSFCOND', 'ACCLASS', 'IMPACTYPE', 'INVTYPE', 'INVAGE', 'INJURY', 'INITDIR', 'VEHTYPE', 'MANOEUVER', 'DRIVACT', 'DRIVCOND', 'PEDTYPE', 'PEDACT',
  8      # 'PEDCOND', 'PEDESTRIAN', 'CYCLIST', 'AUTOMOBILE', 'MOTORCYCLE', 'TRUCK','TRSN_CITY_VEH', 'EMERG_VEH', 'PASSENGER', 'SPEEDING', 'AG_DRIV', 'REDLIGHT',
  9      # 'ALCOHOL', 'DISABILITY', 'POLICE_DIVISION','NEIGHBOURHOOD']]
 10      ## Dropping columns where missing values were greater than 80%
 11      df = df.drop(["PEDTYPE", "PEDACT", "PEDCOND"], axis=1)
 12      # Changing the property damage and non-fatal columns to Non-Fatal
 13      df['ACCLASS'] = np.where(df['ACCLASS'] == 'Property Damage Only', 'Non-Fatal', df['ACCLASS'])
 14      df['ACCLASS'] = np.where(df['ACCLASS'] == 'Non-Fatal Injury', 'Non-Fatal', df['ACCLASS'])
 15
 16      df['MONTH'] = pd.to_datetime(df['DATE']).dt.month
 17      df['MONTH_NAME'] = pd.to_datetime(df['DATE']).dt.month_name()
 18      df['DAY'] = pd.to_datetime(df['DATE']).dt.day
 19      df['MINUTES'] = pd.to_datetime(df['DATE']).dt.minute
 20      df['WEEKDAY'] = pd.to_datetime(df['DATE']).dt.weekday
 21
 22
 23      return df
 24
 25  path = 'data/KSI.csv'
 26  df = wrangle(path)
 27  print(df.shape)
    ✓ 0.8s
  (16860, 59)
```

*Figure 1.1: data preparation function.*

Next, the dataset was cleaned because it contained lots of low and high cardinality features.



```
Data cleaning
  1  df_clean_data = df[['ACCNUM', 'YEAR', 'MONTH', 'DAY', 'HOUR', 'MINUTES', 'WEEKDAY', 'LATITUDE', 'LONGITUDE',
  2  'DISTRICT', 'VISIBILITY', 'LIGHT', 'RDSFCOND', 'PEDESTRIAN', 'CYCLIST', 'AUTOMOBILE', 'MOTORCYCLE', 'TRUCK',
  3      'TRSN_CITY_VEH', 'EMERG_VEH', 'PASSENGER', 'SPEEDING', 'AG_DRIV', 'REDLIGHT', 'ALCOHOL', 'DISABILITY', 'ACCLASS']]
  4
  5  # Transforming the coordinates
  6  df_clean_data['LATITUDE'] = df_clean_data['LATITUDE'].astype('int')
  7  df_clean_data['LONGITUDE'] = df_clean_data['LATITUDE'].astype('int')
```

*Figure 1.2: data preparation function.*

Data analysis of this data was driven using the following questions:

- What story does the data tell regarding the types of collisions?

- What are the accident numbers against years and month?

- What is the fatality heatmap of those that were Fatally and Non-Fatally Injured?

- What is the rate of Fatality over years?

- What areas did accidents occur the most?

- Where are the top 10 Neighborhoods with the most collisions?

- What is the road collision distribution?



*Figure 2.1:* What story does the data tell regarding the types of collisions?



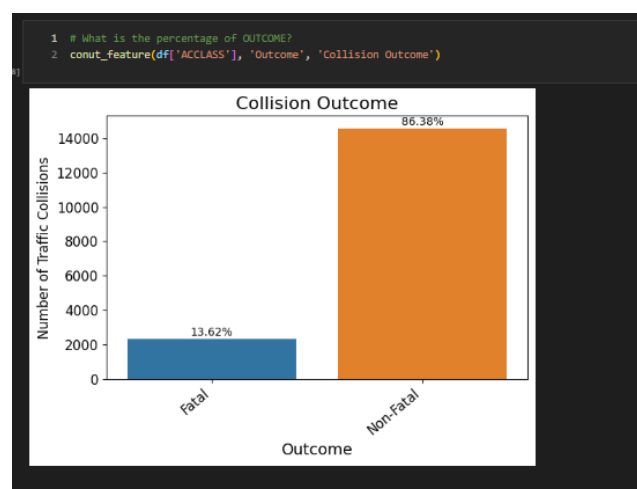*Figure 2.2:* What story does the data tell regarding the types of collisions?
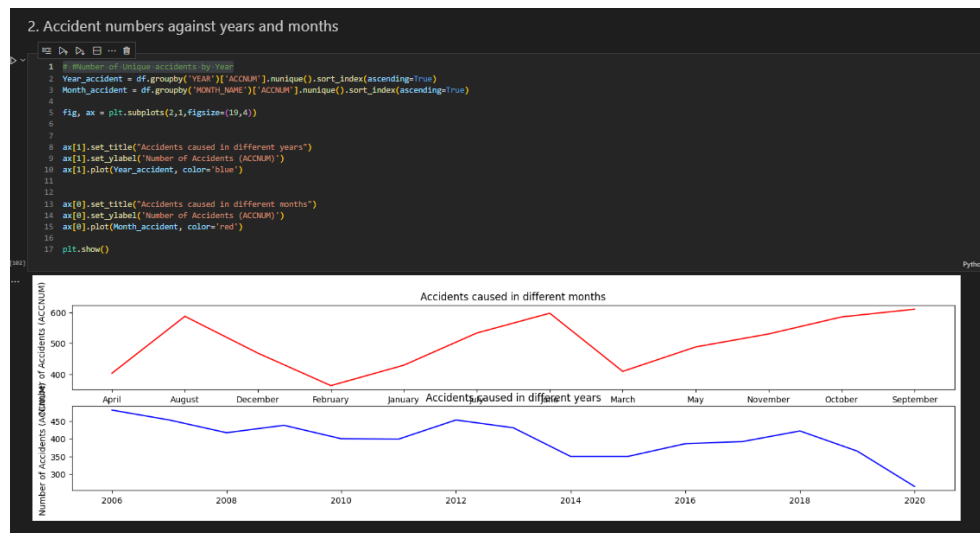
## 2. Accident numbers against years and months

```python
1  # #Number of Unique accidents by Year
2  Year_accident = df.groupby('YEAR')['ACCNUM'].nunique().sort_index(ascending=True)
3  Month_accident = df.groupby('MONTH_NAME')['ACCNUM'].nunique().sort_index(ascending=True)
4
5  fig, ax = plt.subplots(2,1,figsize=(19,4))
6
7
8  ax[1].set_title("Accidents caused in different years")
9  ax[1].set_ylabel('Number of Accidents (ACCNUM)')
10 ax[1].plot(Year_accident, color='blue')
11
12
13 ax[0].set_title("Accidents caused in different months")
14 ax[0].set_ylabel('Number of Accidents (ACCNUM)')
15 ax[0].plot(Month_accident, color='red')
16
17 plt.show()
```



*Figure 3: Accident numbers against years and month*

## 3. Fatality Heatmap of those that where Fatally Injured

```python
1  df_Fatal = df[df['INJURY'] == 'Fatal']
2  df_Fatal = df_Fatal[['LATITUDE', 'LONGITUDE']]
3  lat_Toronto_1 = df_Fatal.describe().at['mean','LATITUDE']
4  lng_Toronto_1 = df_Fatal.describe().at['mean','LONGITUDE']
5  Toronto_location_F = [lat_Toronto_1, lng_Toronto_1]
6  Fatal_map_F = folium.Map(Toronto_location_F, zoom_start=10.255)
7  HeatMap(df_Fatal.values, min_opacity =0.3).add_to(Fatal_map_F)
8  Fatal_map_F
9
10
```



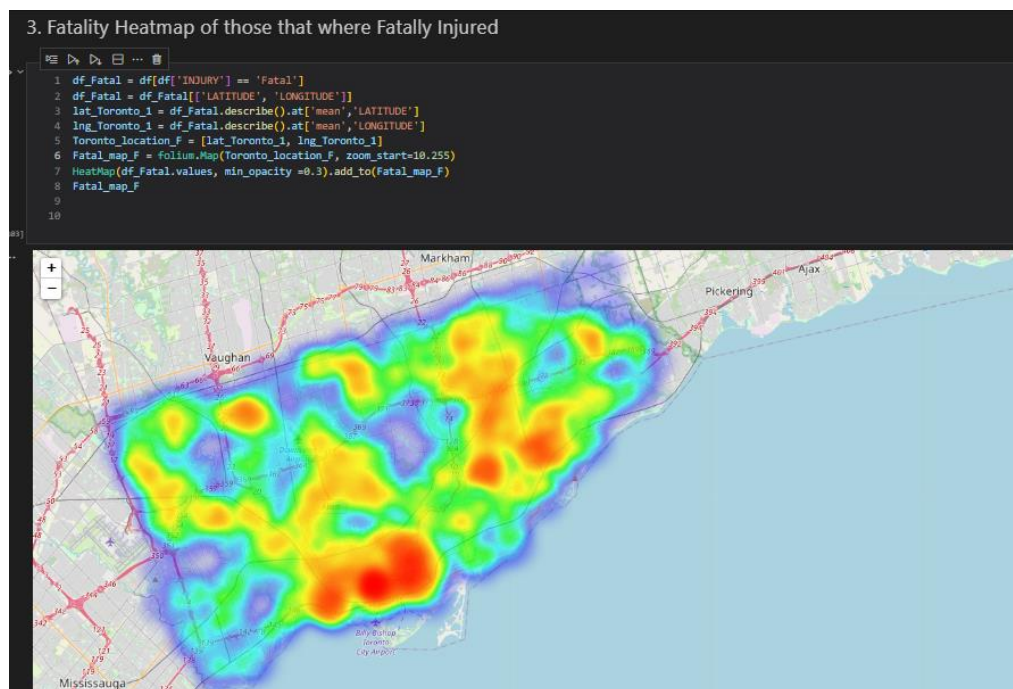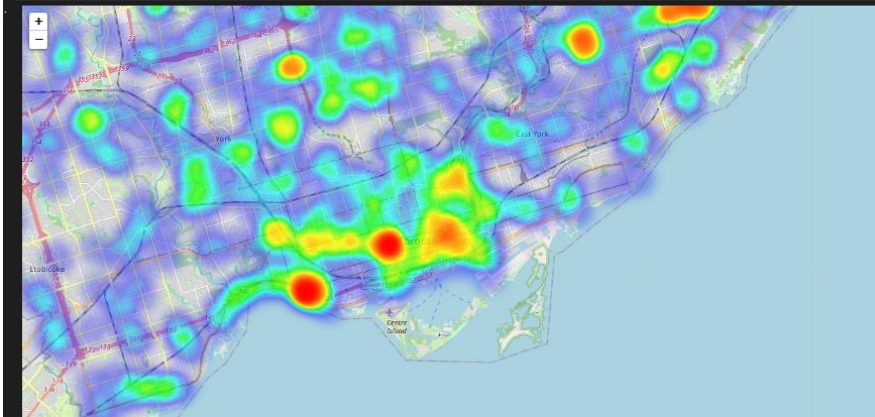*Figure 4: Fatality Heatmap of those that were Fatally Injured*

```
4. Fatality Heatmap of those that where not Fatally Injured

1   df_Non_Fatal = df[df['INJURY'] != 'Fatal']
2   df_Non_Fatal = df_Non_Fatal[['LATITUDE', 'LONGITUDE']]
3   lat_Toronto_2 = df_Non_Fatal.describe().at['mean','LATITUDE']
4   lng_Toronto_2 = df_Non_Fatal.describe().at['mean','LONGITUDE']
5   Toronto_location_N = [lat_Toronto_2, lng_Toronto_2]
6   Fatal_map_N = folium.Map(Toronto_location_F, zoom_start=10.255)
7   HeatMap(df_Fatal.values, min_opacity =0.3).add_to(Fatal_map_N)
8   Fatal_map_N
```

*Figure 5: Fatality Heatmap of those that were Non-Fatally Injured.*



```
5. Fatality over years (# of people died)

1   #Lets look at Fatality over years (# of people died)
2   Fatality = df[df['INJURY'] =='Fatal']
3   Fatality = Fatality.groupby(df['YEAR']).count()
4   plt.figure(figsize=(12,6))
5
6
7   plt.ylabel('Number of Injury=Fatal')
8   Fatality['INJURY'].plot(kind='bar',color="y" , edgecolor='black')
9
10  plt.show()
```
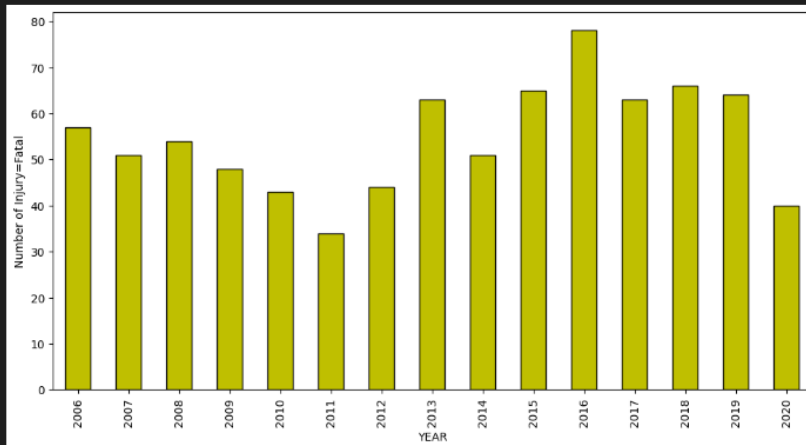
*Figure 6: The rate of Fatality over years (# of people died)*

## 6. Looking at area where accident happens

```
1  Region_df = df['DISTRICT'].value_counts()
2  plt.figure(figsize=(12,6))
3  plt.ylabel('Number of Accidents')
4  Region_df.plot(kind='bar',color=list('rgbkmc') )
5  plt.show()
```
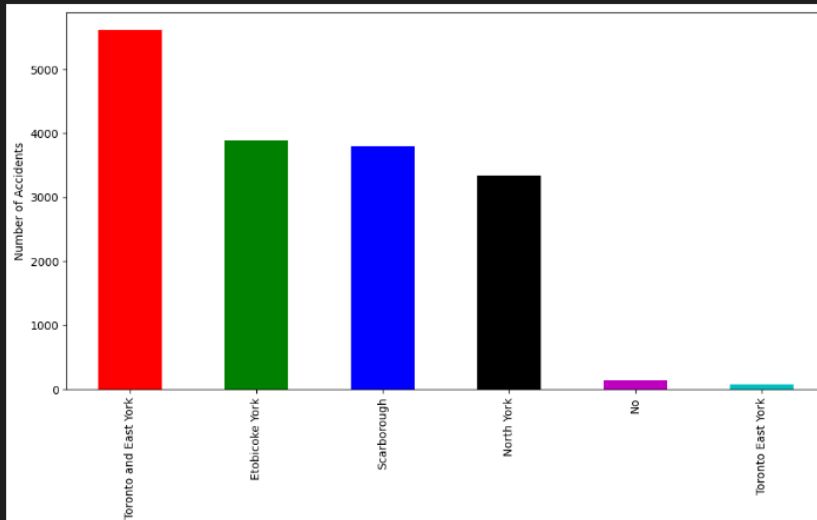


*Figure 7: Areas where accidents occur the most*

## 7. Top 10 Neighbourhood with Most collisions

```
1  Hood_df = df['NEIGHBOURHOOD'].value_counts()
2  plt.figure(figsize=(12,6))
3  plt.ylabel('Number of Accidents')
4  Hood_df.nlargest(10).plot(kind='bar',color=list('rgbkmc') )
5  plt.show()
```
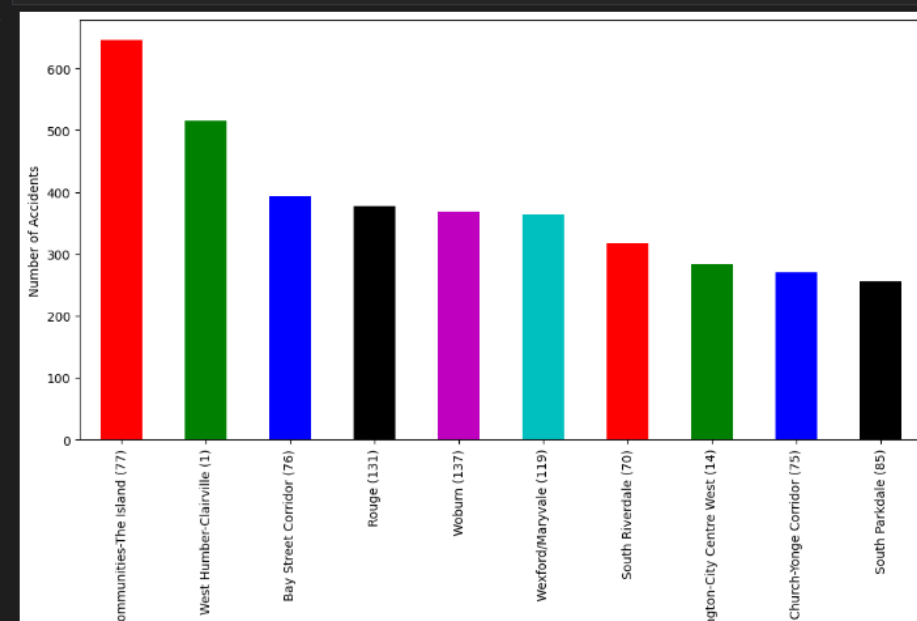


*Figure 8: Areas where accidents occur the most*

```
1
2  conut_feature(df['ROAD_CLASS'], 'level of road classes',"Road Collision")
3
```
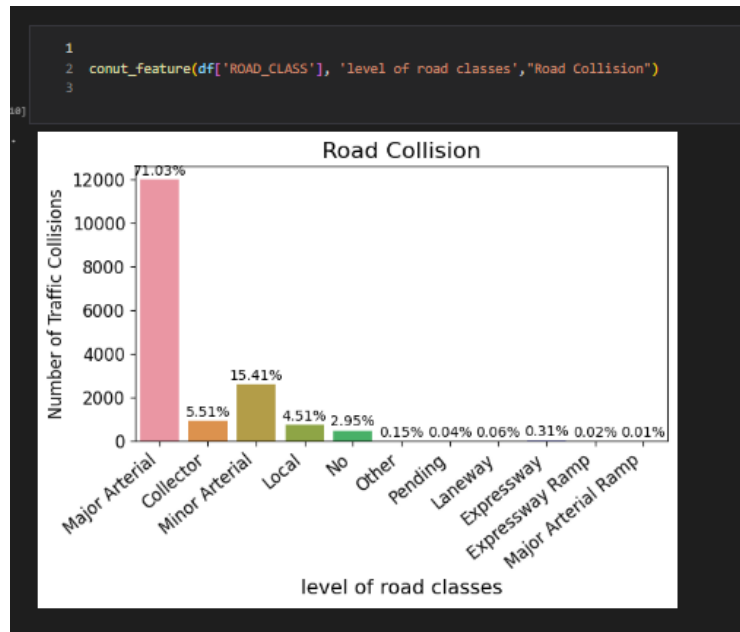


*Figure 9:* Road collision distribution

The data was encoded by getting the dummies of all categorical features in the clean data, to build a

scalable machine learning model. As seen below.

```
Encoding
1  df_clean_data = pd.get_dummies(df_clean_data, columns=['VISIBILITY','RDSFCOND','LIGHT','DISTRICT','PEDESTRIAN','CYCLIST', 'AUTOMOBILE', 'MOTORCYCLE', 'TRUCK',
2      'TRSN_CITY_VEH', 'EMERG_VEH', 'PASSENGER', 'SPEEDING', 'AG_DRIV', 'REDLIGHT', 'ALCOHOL', 'DISABILITY'])
3
4  df_clean_data.info()

<class 'pandas.core.frame.DataFrame'>
```

*Figure 10:* Data Encoding

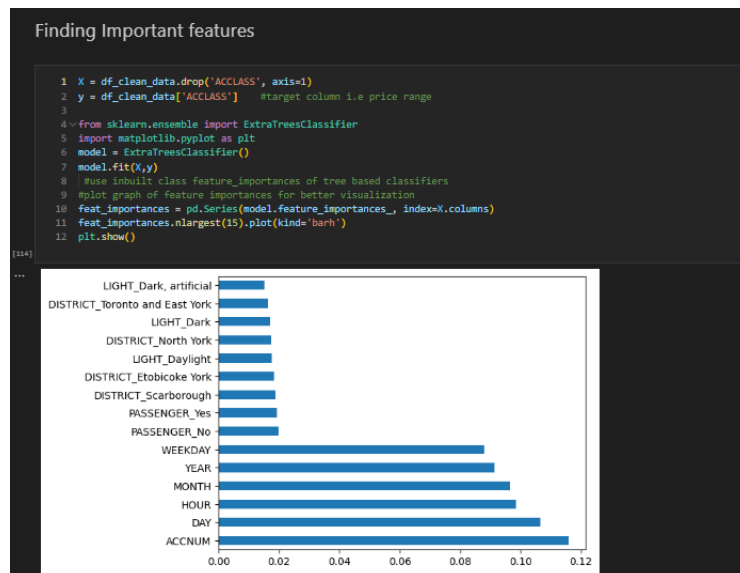The most important feature was gotten by applying the *ExtraTreeClassifier* class from the *sklearn*

module.

```
Finding Important features

1   X = df_clean_data.drop('ACCLASS', axis=1)
2   y = df_clean_data['ACCLASS']    #target column i.e price range
3
4 v from sklearn.ensemble import ExtraTreesClassifier
5   import matplotlib.pyplot as plt
6   model = ExtraTreesClassifier()
7   model.fit(X,y)
8   #use inbuilt class feature_importances of tree based classifiers
9   #plot graph of feature importances for better visualization
10  feat_importances = pd.Series(model.feature_importances_, index=X.columns)
11  feat_importances.nlargest(15).plot(kind='barh')
12  plt.show()
```

*Figure 11:* Feature Extraction

The dataset was then split into a training and a test set using a 70:30 ratio.



```
Split

1   target = 'ACCLASS'
2   X = df_clean_data.drop(target, axis=1)
3   y = df_clean_data[target].map({'Fatal':1,'Non-Fatal':0})


1   from sklearn.model_selection import train_test_split
2
3   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
```
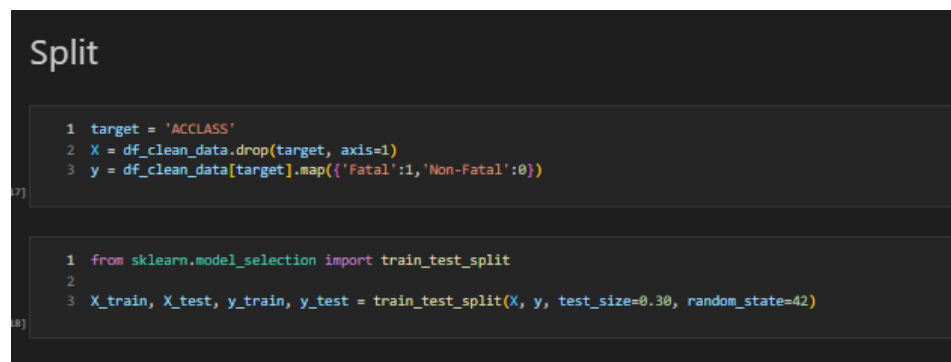
*Figure 12:* Road collision distribution

Afterwards, the training set was used to train each model and the test set was used to evaluate the performance of each model. The training and test data were used on two different machine learning model which are the 1.  Xtreeme boosting classifier and the Support Vector Machine Classifier.
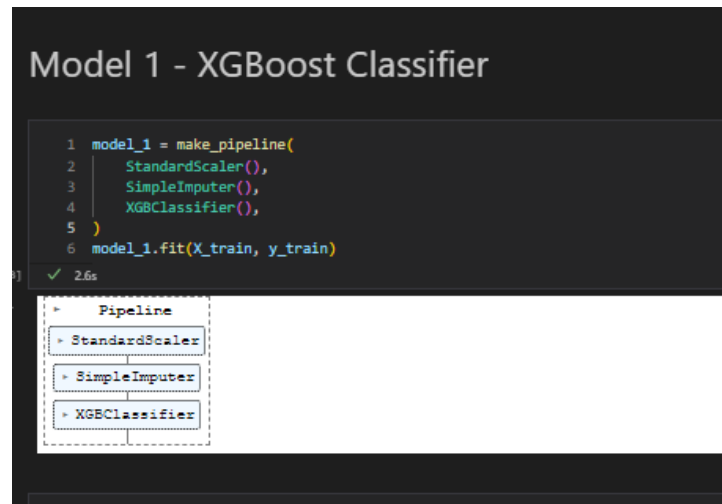
1. **Xtreeme boosting classifier**



*Figure 13:* XGBoost Classifier

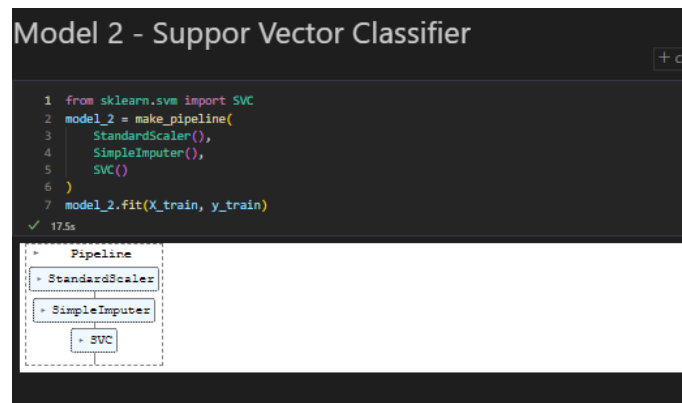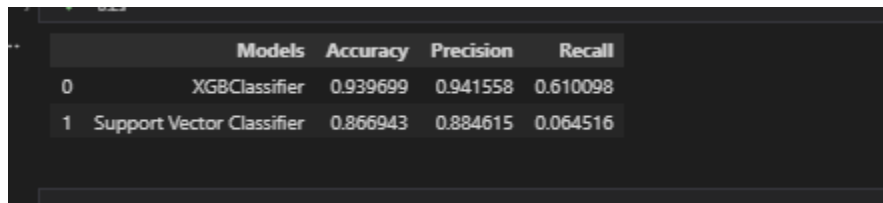2. **Support Vector Machine Classifier**



*Figure 14:* SVM Classifier

Both models were evaluated based on the following metrics.

- Accuracy

- Precision

- Recall



| | Models | Accuracy | Precision | Recall |
|---|---|---|---|---|
| 0 | XGBClassifier | 0.939699 | 0.941558 | 0.610098 |
| 1 | Support Vector Classifier | 0.866943 | 0.884615 | 0.064516 |

*Figure 15:* Model Evaluations

## Conclusion

The test performance for the XGBoost Classifier model was better than the SVM Classifier model with a difference of **0.057144.** The **XGboost** took less time to run compared to the **SVM** Classifier. The Xgboost model generated robust and more significant results than the SVM Classifier as shown in ***figure 15***, that's why the XGBoost model is the best to use for this case study.