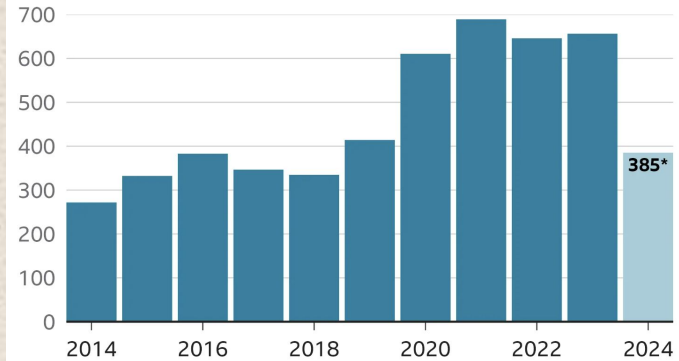# WHY I CHOOSE THIS TOPIC

- **Gun violence is a huge problem within the US**

  - **Over 385 mass shootings in the US alone in 2024 (up to sep 20)**

  - **More than a single mass shooting everyday**

$$\frac{385 \text{ mass shootings}}{262 \text{ days (jan 1 to sep 20)}} = \text{1.4-ish mass shootings a day!}$$

## Mass shootings in the US
Incidents in which four or more people were killed or injured
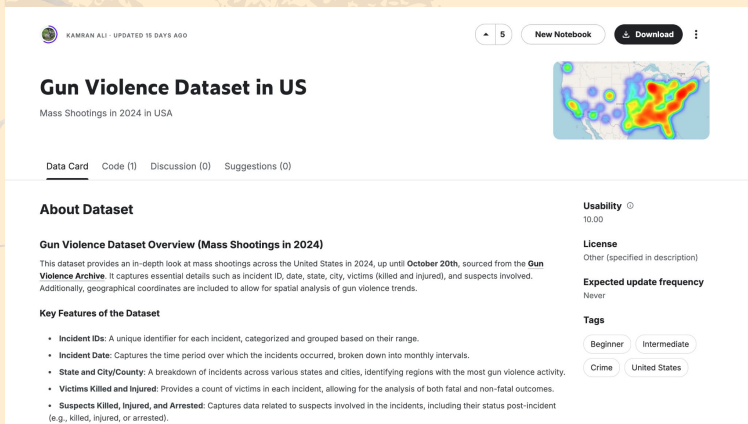
385*

2014  2016  2018  2020  2022  2024

Source: Gun Violence Archive (*data up to 5 September)  BBC

## Mass shootings on the rise

There have been more than 385 mass shootings across the US so far this year, according to the Gun Violence Archive, which defines a mass shooting as an incident in which four or more people are injured or killed. Their figures include shootings that

# THE DATA SET



- **Got this data set from Kaggle**

- **Covers 2024 U.S. mass shootings (up to Oct 20).**

- **Includes incident IDs, dates, locations, victims, and suspects.**

- **Geographical data supports spatial trend analysis.**

# HOW I ANALYZED IT

- **Imported data and verified structure**

- **Cleaned data, addressed missing value**

- **Selected relevant columns**

- **Divided data into training and test sets.**

```python
#loading in the data
file_path = '/content/updated_with_coordinates.csv'
data = pd.read_csv(file_path)

print("Columns in the dataset:", data.columns)

if 'Victims Killed' in data.columns:
    data['fatalities'] = data['Victims Killed'].apply(lambda x: 1 if x > 0 else 0)
    y = data['fatalities']
else:
    raise KeyError("'Victims Killed' column not found. Please check the column names.")

# select features
features = ['Latitude', 'Longitude']
X = data[features]
```

```python
[7] #splitting the data
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

# MODELS USED AND WHY?

```
[8] #KNN
    knn = KNeighborsClassifier(n_neighbors=5)
    knn.fit(X_train, y_train)
    knn_pred = knn.predict(X_test)

    knn_metrics = {
        "Model": "KNN",
        "Accuracy": accuracy_score(y_test, knn_pred),
        "Precision": precision_score(y_test, knn_pred),
        "Recall": recall_score(y_test, knn_pred),
        "F1 Score": f1_score(y_test, knn_pred)
    }

    #Logistic Regression
    log_reg = LogisticRegression()
    log_reg.fit(X_train, y_train)
    log_reg_pred = log_reg.predict(X_test)

    log_reg_metrics = {
        "Model": "Logistic Regression",
        "Accuracy": accuracy_score(y_test, log_reg_pred),
        "Precision": precision_score(y_test, log_reg_pred),
        "Recall": recall_score(y_test, log_reg_pred),
        "F1 Score": f1_score(y_test, log_reg_pred)
    }

[10] #Decision Tree Classifier
    decision_tree = DecisionTreeClassifier(random_state=42)
    decision_tree.fit(X_train, y_train)
    tree_pred = decision_tree.predict(X_test)

    tree_metrics = {
        "Model": "Decision Tree",
        "Accuracy": accuracy_score(y_test, tree_pred),
        "Precision": precision_score(y_test, tree_pred),
        "Recall": recall_score(y_test, tree_pred),
        "F1 Score": f1_score(y_test, tree_pred)
    }

[11] # SVC
    svm = SVC()
    svm.fit(X_train, y_train)
    svm_pred = svm.predict(X_test)

    svm_metrics = {
        "Model": "SVM",
        "Accuracy": accuracy_score(y_test, svm_pred),
        "Precision": precision_score(y_test, svm_pred),
        "Recall": recall_score(y_test, svm_pred),
        "F1 Score": f1_score(y_test, svm_pred)
    }
```

- **K-Nearest Neighbors (KNN)**

- **Logistic Regression**

- **Decision Tree**

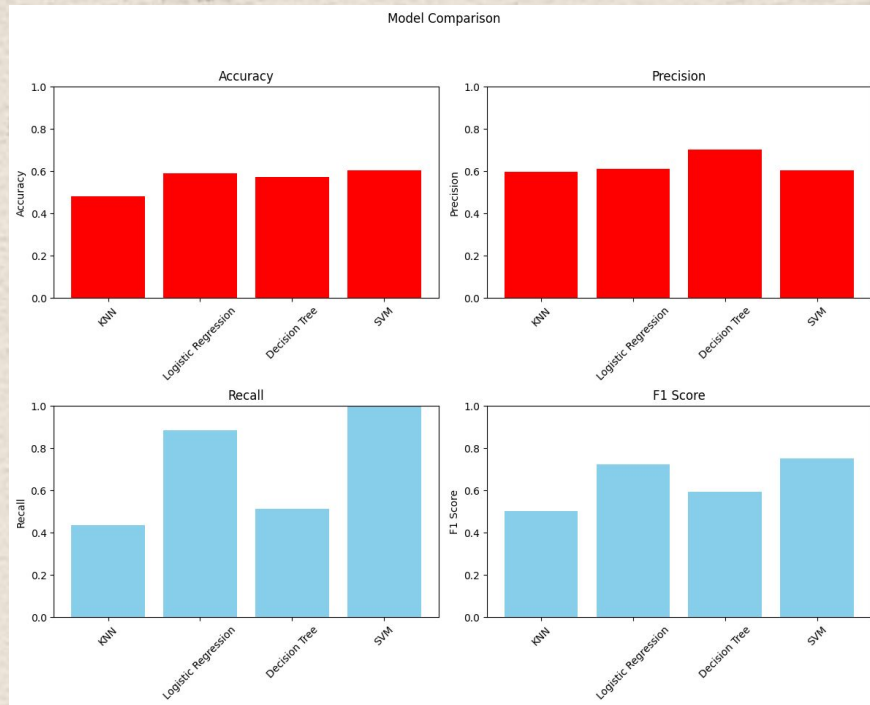- **Support Vector Machine (SVM)**

# CHALLENGES FACED



- **Data quality issues required preprocessing**

- **Precision limitations in models due to imbalanced data**

# THE RESULTS

# NEXT STEPS

**STEP 1** → **STEP 2** → **STEP 3** → **STEP 4**

Explore additional machine learning models

Enhance data preprocessing techniques

Experiment with ensemble methods for improved accuracy

Expand features and test model stability on larger/more datasets