In [3]: 
```python
import pandas as pd
import numpy as np
```

In [6]: 
```python
#load the csv data
df = pd.read_csv("AviationData.csv",encoding='latin1',low_memory="false")
df.head()
```

```
C:\Users\HP\anaconda3\envs\learn-env\lib\site-packages\IPython\core\interactive
shell.py:3145: DtypeWarning: Columns (6,7,28) have mixed types.Specify dtype op
tion on import or set low_memory=False.
  has_raised = await self.run_ast_nodes(code_ast.body, cell_name,
```

Out[6]:

|   | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | Latitud |
|---|----------|--------------------|-----------------|------------|----------|---------|---------|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | Na |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | Na |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | 36.92: |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | Na |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | Na |

5 rows × 31 columns

In [72]:
```
df=df.drop_duplicates()
df
```

Out[72]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | L |
|---|---|---|---|---|---|---|---|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | ᛃ |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| ... | ... | ... | ... | ... | ... | ... | |
| 88884 | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| 88885 | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| 88886 | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | 34 |
| 88887 | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| 88888 | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88889 rows × 31 columns

In [73]:
```
# find the number of rows and columns
df.shape
```

Out[73]: (88889, 31)

# Finding sum of the missing values

In [74]: 
```python
#check for missing values
df.isna().sum()
```

Out[74]: 
```
Event.Id                      0
Investigation.Type            0
Accident.Number               0
Event.Date                    0
Location                     52
Country                     226
Latitude                  54507
Longitude                 54516
Airport.Code              38640
Airport.Name              36099
Injury.Severity            1000
Aircraft.damage            3194
Aircraft.Category         56602
Registration.Number        1317
Make                         63
Model                        92
Amateur.Built               102
Number.of.Engines          6084
Engine.Type                7077
FAR.Description           56866
Schedule                  76307
Purpose.of.flight          6192
Air.carrier               72241
Total.Fatal.Injuries      11401
Total.Serious.Injuries    12510
Total.Minor.Injuries      11933
Total.Uninjured            5912
Weather.Condition          4492
Broad.phase.of.flight     27165
Report.Status              6381
Publication.Date          13771
dtype: int64
```

In [75]:
```python
#find sum of the duplicate data
df.duplicated().sum()
df
```

Out[75]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | L |
|---|---|---|---|---|---|---|---|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | 3 |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| ... | ... | ... | ... | ... | ... | ... | |
| 88884 | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| 88885 | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| 88886 | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | 3 |
| 88887 | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| 88888 | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88889 rows × 31 columns

In [76]:
```python
#drop columns that aren't necessary
cols_to_drop =['Latitude', 'Longitude', 'Airport.Code', 'Airport.Name',
               'Aircraft.Category', 'FAR.Description', 'Schedule', 'Air.carrier'
df_cleaned =df.drop(columns=cols_to_drop)
df_cleaned
```

Out[76]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | In |
|---|---|---|---|---|---|---|---|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| ... | ... | ... | ... | ... | ... | ... | |
| 88884 | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| 88885 | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| 88886 | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | |
| 88887 | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| 88888 | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88889 rows × 23 columns

In [77]: 
```python
#drop rows where 'make' or 'model'are missing
df_cleaned.dropna(subset=['Make','Model'],inplace=True)
df_cleaned
```

Out[77]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | In |
|---|---|---|---|---|---|---|---|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| ... | ... | ... | ... | ... | ... | ... | |
| 88884 | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| 88885 | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| 88886 | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | |
| 88887 | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| 88888 | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88777 rows × 23 columns

In [78]:
```python
#fill in missing numerical data
injury_col = ['Total.Fatal.Injuries', 'Total.Serious.Injuries', 'Total.Minor.Inju
df_cleaned[injury_col] =df_cleaned[injury_col].fillna(0)
df_cleaned
```

Out[78]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | In |
|---|---|---|---|---|---|---|---|
| **0** | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| **1** | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| **2** | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | |
| **3** | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| **4** | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| **...** | ... | ... | ... | ... | ... | ... | |
| **88884** | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| **88885** | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| **88886** | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | |
| **88887** | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| **88888** | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88777 rows × 23 columns

In [79]:
```python
# fill in the missing categorical data
catg_col=['Location','Weather.Condition', 'Country', 'Injury.Severity','Purpose.o
          'Engine.Type','Publication.Date', 'Broad.phase.of.flight','Aircraft.o
df_cleaned[catg_col]=df_cleaned[catg_col].fillna('unknown')
df_cleaned
```

Out[79]:

| | Event.Id | Investigation.Type | Accident.Number | Event.Date | Location | Country | In |
|---|---|---|---|---|---|---|---|
| 0 | 20001218X45444 | Accident | SEA87LA080 | 1948-10-24 | MOOSE CREEK, ID | United States | |
| 1 | 20001218X45447 | Accident | LAX94LA336 | 1962-07-19 | BRIDGEPORT, CA | United States | |
| 2 | 20061025X01555 | Accident | NYC07LA005 | 1974-08-30 | Saltville, VA | United States | |
| 3 | 20001218X45448 | Accident | LAX96LA321 | 1977-06-19 | EUREKA, CA | United States | |
| 4 | 20041105X01764 | Accident | CHI79FA064 | 1979-08-02 | Canton, OH | United States | |
| ... | ... | ... | ... | ... | ... | ... | |
| 88884 | 20221227106491 | Accident | ERA23LA093 | 2022-12-26 | Annapolis, MD | United States | |
| 88885 | 20221227106494 | Accident | ERA23LA095 | 2022-12-26 | Hampton, NH | United States | |
| 88886 | 20221227106497 | Accident | WPR23LA075 | 2022-12-26 | Payson, AZ | United States | |
| 88887 | 20221227106498 | Accident | WPR23LA076 | 2022-12-26 | Morgan, UT | United States | |
| 88888 | 20221230106513 | Accident | ERA23LA097 | 2022-12-29 | Athens, GA | United States | |

88777 rows × 23 columns

In [80]:
```python
#save the clean data
df_cleaned.to_csv(r"C:\Users\HP\Documents\PROJECT\AviationData_Cleaned.csv", inde
```
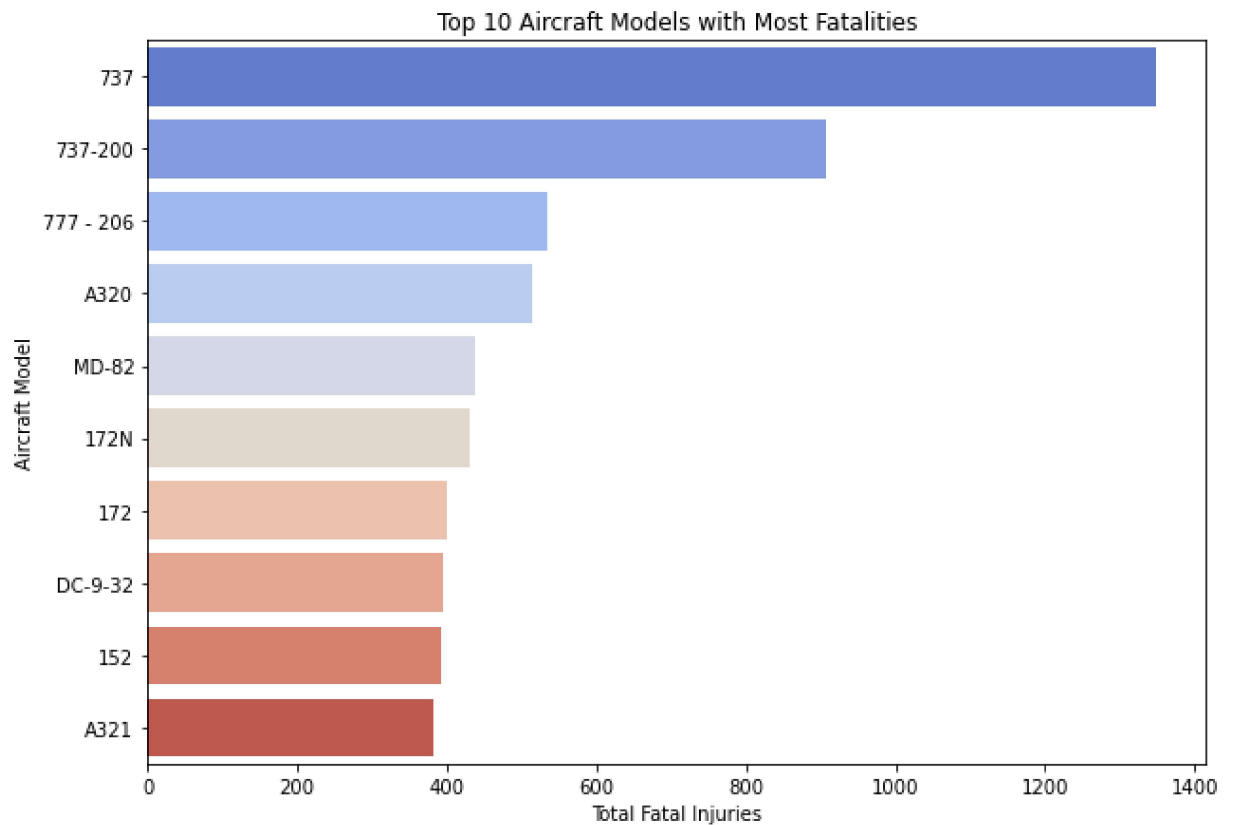
In [81]:
```python
import matplotlib.pyplot as plt
import seaborn as sns
```

In [86]:

```python
fatalities_per_model = df_cleaned.groupby("Model")["Total.Fatal.Injuries"].sum().

plt.figure(figsize=(10,7))
sns.barplot(x=fatalities_per_model.values, y=fatalities_per_model.index, palette=
plt.title("Top 10 Aircraft Models with Most Fatalities")
plt.xlabel("Total Fatal Injuries")
plt.ylabel("Aircraft Model")
plt.show()
```



In [ ]:

In [ ]: