

# Exploring Aircrafts And Their Risks: A Data-Driven Approach

Author: Caleb Kinondo

## Overview

Choosing aircrafts for investment purposes is difficult due to the many risks that a particular aircraft encompasses. The aim of this project is to explore various aircrafts and the accompanied risks, in order to make a balanced approach of choosing an aircraft with the lowest risk. The dataset is from the National Transportation Safety Board that includes aviation accident data from 1962 to 2023 about 90,948 aviation accidents and selected incidents in the United States and international waters.

## Data Understanding

In [3]:	<pre>aviation_data.shape</pre>
Out[3]:	<pre>(90348, 31)</pre>
From above, we can see that the dataset has 90348 rows and 31 columns	
In [4]:	<pre># Enables the viewing of column names, their data types and the number of values in them aviation_data.info()</pre> <div>&lt;class 'pandas.core.frame.DataFrame'&gt; RangeIndex: 90348 entries, 0 to 90347 Data columns (total 31 columns): # Column Non-Null Count Dtype -- -- -- 0 Event.Id 88889 non-null object 1 Investigation.Type 90348 non-null object 2 Accident.Number 88889 non-null object 3 Event.Date 88889 non-null object 4 Location 88837 non-null object 5 Country 86963 non-null object 6 Latitude 34382 non-null object 7 Longitude 34373 non-null object 8 Airport.Code 89248 non-null object 9 Airport.Name 52780 non-null object 10 Injury.Severity 87889 non-null object 11 Aircraft.damage 85695 non-null object 12 Aircraft.Category 32287 non-null object 13 Registration.Number 87872 non-null object 14 Make 88826 non-null object 15 Model 88787 non-null object 16 Amateur.Built 7543 non-null object 17 Number.of.Engines 82895 non-null float64 18 Engine.Type 81812 non-null object 19 FAR.Description 32823 non-null object 20 Schedule 12582 non-null object 21 Purpose.of.flight 82697 non-null object 22 Air.carrier 16648 non-null object 23 Total.Fatal.Injuries 77485 non-null float64 24 Total.Serious.Injuries 76379 non-null float64 25 Total.Minor.Injuries 76956 non-null float64 26 Total.Uninjured 82077 non-null object 27 Weather.Condition 84397 non-null object 28 Broad.phase.of.flight 81724 non-null object 29 Report.Status 82508 non-null object 30 Publication.Date 73659 non-null object dtypes: float64(5), object(26) memory usage: 21.4+ MB</div>

In [5]:	<pre># Returns the number of null values in each column aviation_data.isna().sum()</pre>
Out[5]:	<pre>Event.Id      11459 Investigation.Type      0 Accident.Number      1459 Event.Date      1459 Location      1511 Country      1685 Latitude      5695 Longitude      55975 Airport.Code      40999 Airport.Name      37558 Injury.Severity      2459 Aircraft.damage      4683 Aircraft.Category      58951 Registration.Number      2776 Model      1551 Amateur.Built      1561 Number.of.Engines      7543 Engine.Type      8536 FAR.Description      58295 Schedule      77766 Purpose.of.Flight      7651 Air.carrier      73780 Total.Fatal.Injuries      12860 Total.Serious.Injuries      13969 Total.Minor.Injuries      13282 Total.Uninjured      7971 Weather.Condition      5951 Broad.phase.of.flight      28624 Report.Status      7840 Publication.Date      16689 dtype: int64</pre>
From above, we can see that some columns have a lot of missing values, for example the Schedule column, Air.carrier column and FAR.Description.	

In [6]:	<pre># Returns the percentage of missing values in each column percent_missing = aviation_data.isnull().sum() * 100 / len(aviation_data) percent_missing = percent_missing.sort_values(ascending = False) percent_missing</pre>
Out[6]:	<pre>Schedule      86.873848 Air.carrier      81.573471 FAR.Description      64.595939 Aircraft.Category      64.263736 Longitude      61.944886 Latitude      61.954886 Airport.Code      61.944924 Airport.Name      44.382831 Air.carrier      41.570372 Broad.phase.of.flight      31.683341 Publication.Date      18.473909 Total.Serious.Injuries      15.463227 Total.Minor.Injuries      14.822686 Total.Fatal.Injuries      14.233851 Engine.Type      9.447913 Report.Status      8.677558 Purpose.of.flight      8.408367 Number.of.Engines      8.348529 Total.Uninjured      8.158454 Weather.Condition      6.586753 Aircraft.damage      5.159866 Registration.Number      3.072564 Injury.Severity      2.721698 Country      1.859511 Amateur.Built      1.727764 Model      1.716695 Make      1.684597 Location      1.672422 Event.Date      1.614867 Accident.Number      1.614867 Event.Id      1.614867 Investigation.Type      0.008080</pre>

By finding the percentage of missing values in each column, I can choose the best way to deal with the missing values.

Country1.885611  
Aviation.Type1.727754  
Model1.736695  
Make1.484587  
Location1.672422  
Event.Date1.614867  
Accident.Number1.614867  
Event.Id1.614867  
Investigation.Type0.000000  
dtype: float64

By finding the percentage of missing values in each column, I can choose the best way to deal with the missing values.

In [7]:

```
# Finding out the number of duplicate events based on the Event.Id column  
duplicates = aviation_data[aviation_data.duplicated(subset='Event.Id')]  
print(len(duplicates))  
duplicates
```

2396

Out [7]:

	Event.Id	Investigation.Type	Accident.Number	Event.Date	Location	Country	Latitude	Longitude	Airport.Code	Airport.Name	...	Purpose.of.flight	Air.carrier	Total.Fatal.Injuries
118	20000917X01908	Accident	DCA82FA012A	1982-01-19	ROCKPORT, TX	United States	NaN	NaN	RKP	ARANSAS COUNTY AIRPORT	...	Executive/Corporate	NaN	
159	20000917X02400	Accident	MA82FA038A	1982-01-23	NEWPORT RICHEY, FL	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	
160	20000917X02259	Accident	LAX82FA049B	1982-01-23	VICTORVILLE, CA	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	
245	20000917X02568	Accident	SEA82DA020B	1982-02-06	MEDFORD, OR	United States	NaN	NaN	MFR	MEDFORD-JACKSON COUNTY	...	Personal	NaN	
248	20000917X02173	Accident	LAX82DA065A	1982-02-06	SAN JOSE, CA	United States	NaN	NaN	RHV	RED HILLVIEW	...	Personal	NaN	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
90097	NaN	20-12-2002	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
90286	20021121102676	Accident	CEN23MA034	2002-11-12	Dallas, TX	United States	324020N	0965140W	RBD	Dallas Executive	...	AS/HC Commemorative Air Force	0	0
90295	20021121100636	Accident	WPR23LA041	2002-11-18	Las Vegas, NV	United States	361238N	1151140W	VGT	NORTH LAS VEGAS	...	Instructional HELICOPTER (HC)	0	0
90297	20021121100630	Incident	DCA23WA071	2002-11-18	Marrakech, Morocco	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
90275	20021210X03654	Accident	WPR23LA045	2002-11-22	San Diego, CA	United States	323414N	1166825W	SDM	Brown Field Municipal Airport	...	Public Aircraft - Federal	0	0

2396 rows x 31 columns

## Data Preparation

### Dealing with duplicate values

First, I will drop the duplicates values.

Out[8]:

	Event.Id	Investigation.Type	Accident.Number	Event.Date	Location	Country	Latitude	Longitude	Airport.Code	Airport.Name	...	Purpose.of.flight	Air.carrier	Total.Fatal.Injuries	Total.Minor.Injuries
0	20001218X45444	Accident	SEAB7LA080	1948-10-24	MOOSE CREEK, ID	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	2.0	0.0
1	20001218X45447	Accident	LAX9LA336	1962-07-19	BRIDGEPORT, CA	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	4.0	0.0
2	2000105X01555	Accident	NYC07LA065	1974-08-30	Salville, VA	United States	36.9222	-81.8781	NaN	NaN	...	Personal	NaN	3.0	0.0
3	20001218X45448	Accident	LAX9LA321	1977-06-19	EUREKA, CA	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	2.0	0.0
4	20001105X01764	Accident	CHY7FA064	1979-08-02	Canton, OH	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	1.0	0.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
90343	20021227106491	Accident	ERA23LA093	2002-12-26	Annapolis, MD	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	0.0	0.0
90344	20021227106494	Accident	ERA23LA095	2002-12-26	Hampton, NH	United States	NaN	NaN	NaN	NaN	...	NaN	NaN	0.0	0.0
90345	20021227106497	Accident	WPR23LA075	2002-12-26	Payson, AZ	United States	341525N	1110201W	PAN	PAVSON	...	Personal	NaN	0.0	0.0
90346	20021227106498	Accident	WPR23LA076	2002-12-26	Morgan, UT	United States	NaN	NaN	NaN	NaN	...	Personal	CESNA 210N LLC	0.0	0.0
90347	20021227106513	Accident	ERA23LA097	2002-12-29	Athens, GA	United States	NaN	NaN	NaN	NaN	...	Personal	NaN	0.0	0.0

87952 rows x 31 columns

Dealing with missing values

Columns

Let's revisit the percentage of missing values in each column

In [8]:

	percent_missing
Schedule	86.873848
Air.carrier	81.873471
FAA.description	64.555939
Aircraft.Category	4.263736
Location	61.554866

Out[8]:

### Dealing with missing values

#### Columns

Let's revisit the percentage of missing values in each column

```
percent_missing
Schedule      86.873848
Air.carrier      81.573471
FAR.Description      64.559339
Aircraft.Category      64.263736
Longitude      61.944886
Latitude      61.944924
Airport.Code      44.382831
Airport.Name      41.570372
Broad.phase.of.flight      31.683341
Publication.Date      18.473909
Total.Serious.Injuries      15.463227
Total.Minor.Injuries      14.822686
Total.Fatal.Injuries      14.233851
Engine.Type      9.447913
Purpose.of.flight      8.408367
Number.of.Engines      8.348529
Total.Uninjured      8.158454
Weather.Condition      6.586753
Aircraft.damage      5.159866
Registration.Number      3.072564
Injury.Severity      2.721698
Country      1.859511
Amateur.Built      1.727764
Model      1.716695
Make      1.684597
Location      1.672422
Event.Date      1.614867
Accident.Number      1.614867
Event.Id      1.614867
Investigation.Type      0.008080
dtype: float64
```

In one line, I will:

- Drop the Air.carrier and Schedule columns since they have the most missing values: over 80%. I will drop all irrelevant columns.
- Drop the Aircraft.Category, Airport.Code and Airport.Name columns because they are irrelevant and have a high percentage of missing values(over 40%).
- Drop the Longitude, FAR.Description and Latitude columns because they have a high percentage of missing values (60%) and is irrelevant.
- Drop any other irrelevant column

```
In [10]: aviation_data = aviation_data.drop(['Air.carrier', 'Schedule', 'Aircraft.Category', 'Airport.Code', 'Airport.Name', 'Longitude', 'Latitude', 'FAR.Description'], axis=1)
aviation_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 87952 entries, 0 to 87951
Data columns (total 17 columns):
# Column Non-Null Count Dtype
-- -- --
0 Event.Id 87951 non-null object
1 Investigation.Type 87951 non-null object
2 Accident.Number 87951 non-null object
3 Event.Date 87951 non-null object
4 Injury.Severity 86963 non-null object
5 Aircraft.damage 84848 non-null object
6 Make 87888 non-null object
7 Model 87859 non-null object
8 Amateur.Built 7543 non-null object
9 Number.of.Engines 83024 non-null float64
10 Engine.Type 80927 non-null object
11 Purpose.of.flight 81829 non-null object
12 Total.Fatal.Injuries 76684 non-null float64
13 Total.Serious.Injuries 75629 non-null float64
14 Total.Minor.Injuries 76331 non-null float64
15 Total.Uninjured 82888 non-null float64
16 Broad.phase.of.flight 60837 non-null object
dtypes: float64(5), object(12)
memory usage: 18.4+ MB
```

The Broad.phase.of.flight column has a missing value percentage of 30%. Since this column consists of categorical data and is important in my analysis, I decided to drop the rows that consist of NaN values in the column.

```
In [11]: aviation_data = aviation_data.dropna(subset=['Broad.phase.of.flight'])
aviation_data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 60837 entries, 0 to 63911
Data columns (total 17 columns):
# Column Non-Null Count Dtype
-- -- --
0 Event.Id 60837 non-null object
1 Investigation.Type 60837 non-null object
2 Accident.Number 60837 non-null object
3 Event.Date 60837 non-null object
4 Aircraft.damage 59459 non-null object
5 Make 66887 non-null object
6 Model 66819 non-null object
7 Number.of.Engines 59942 non-null float64
8 Engine.Type 46946 non-null object
9 Purpose.of.flight 59780 non-null object
10 Total.Fatal.Injuries 50299 non-null float64
11 Total.Serious.Injuries 49818 non-null float64
12 Total.Minor.Injuries 50248 non-null float64
13 Total.Uninjured 55871 non-null float64
14 Broad.phase.of.flight 60837 non-null object
15 Fatality 66832 non-null object
dtypes: float64(4), object(12)
memory usage: 12.0+ MB
```

<ipython-input-12-ed8534a48bb>:17: SettingWithCopyWarning:  
A value is being set on a copy of a slice from a DataFrame  
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy  
aviation\_data['Fatality'].fillna(aviation\_data['Fatality'].isna() == 0, inplace=True)

In [14]:	<pre># Check the percentage of missing values per column for the updated dataset tmp = aviation_data percent_missing = tmp.isnull().sum() * 100 / len(aviation_data) percent_missing = percent_missing.sort_values(ascending = False) percent_missing</pre>
Out[14]:	<pre>Total.Serious.Injuries      18 Total.Minor.Injuries      17 Total.Uninjured      8 Aircraft.damage      2 Purpose.of.flight      2 Number.of.Engines      1 Engine.Type      1 Model      0 Amateur.Built      0 Make      0 Fatality      0 Year      0 Broad.phase.of.flight      0 Accident.Number      0 Investigation.Type      0 Event.Id      0 dtype: float64</pre>

There are still some missing values, notably in the 'Total.Serious.Injuries' and 'Total.Minor.Injuries' columns but I do not find them significant in my analysis.

## Data Modeling

Let's investigate the relationship of make, models and accidents.

```

# Replaces the categorical values with 0
aviation_data['Fatality'].replace({'Non-Fatal': 0, 'Minor': 0, 'Serious': 0, 'Incident': 0, inplace=True)

#
aviation_data['Fatality'] = aviation_data.apply(lambda row: row['Total.Fatal.Injuries'] if row['Fatality'] == 'Fatal' else row['Fatality'], axis=1)

# Replaces the 'Unavailable' values with NaN
aviation_data['Fatality'].replace('Unavailable', np.nan, inplace=True)

# Changes the values from strings to integers
aviation_data['Fatality'] = aviation_data['Fatality'].fillna(0)
aviation_data['Fatality'] = aviation_data['Fatality'].astype(int)

# Changes to float
pd.options.display.float_format = '{:.0f}'.format

aviation_data = aviation_data.drop(['Injury.Severity', 'Total.Fatal.Injuries'], axis = 1)
aviation_data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 68687 entries, 0 to 68691
Data columns (total 26 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Event.Id              68687 non-null  object
 1   Investigation.Type     68687 non-null  object
 2   Accident.Number       68687 non-null  object
 3   Event.Date            59459 non-null  object
 4   Aircraft.Damage       59459 non-null  object
 5   Make                  68686 non-null  object
 6   Model                 68687 non-null  object
 7   Manufacturer          68619 non-null  object
 8   Number of Engines     59942 non-null  float64
 9   Engine.Type           64556 non-null  object
10   Purpose of Flight     59798 non-null  object
11   Total.Serious.Injuries 46616 non-null  float64
12   Total.Minor.Injuries  59246 non-null  float64
13   Total.Uninjured       55871 non-null  float64
14   Broad.Phase.of.Flight  68687 non-null  object
15   Fatality               68682 non-null  object
dtypes: float64(4), object(12)
memory usage: 10.4+ MB

<python-input-12-e5b534a46bb>:17: 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/user_guide/indexing.html#returning-a-view-versus-a-copy
aviation_data['Fatality'] = aviation_data['Fatality'].fillna(0)
aviation_data['Fatality'] = aviation_data['Fatality'].astype(int)

```