

Airplane crashes between Numbers and Regulations

Mohammad Saleh

Abstract—Airplane accidents, still the most dangerous life threatening accident. The main reason for this believe or fearness is the number of survivor across each accident. Nowadays, it is a fact that nobody will survive an airplane accident. The global aviation administrations, mainly FAA (Federal aviation administration) in USA, regularly release new rules to maintain and decrease the number of airplane crashes. Many people show the statistics for these crashes but the results shows that this number is slightly decreasing over the years because they didn't take into account the number of flights yearly. In this letter, we will present the latest analysis for the number of crashes per year globally and among the nations taking into account the increasing number in the flights yearly. We will highlight the main causes of the crashes and the proposed solutions for the fulfillment of the final project of the machine learning course.

Index Terms—crashes, paper, analysis, airplane.

I. INTRODUCTION

AIRPLANE is one the greatest human inventions. Nowadays, the technology for the new inventions in the aerospace field is growing fast, mainly in the drones technology. Lately, the FAA in the USA released a new regulations regarding drone operations to maintain the safety of the commercial airplanes. In this analysis, we wanted to check how effective the regulations and safety on the airplane crashes and the number of survivors. It is very well known that the FAA is one of the main administrations for airplanes safety. We will start with general overview of the crashes yearly and the main reasons and check the proportionally of airplane crashes number in the USA.

II. ANALYSIS

A. Crashes per year

Many people have worked on this analysis in terms of showing the number of crashes yearly without taking into account the number of flights yearly. This analysis has been forked from previous one at the Kaggle website. Figure 1 shows the number of crashes yearly normalized by the number of passengers yearly as we don't have direct access to the number of flights yearly. The results show the number of crashes till 2009 because of the limitation of the dataset. The numbers on the y-axis are up to a constant value, the most important thing is the general trend in the number of crashes. It shows that the number of crashes decrease as times go, except couple years (1988, 1991) that had some reasons for having large number of crashes. In the following paragraphs, we will look in the causes of the crashes and the main reason for this decrease in the number of crashes.

B. Cause of crashes

Figure 2 shows the main reason for the crashes where the y-axis is in log-scale for comparison. The highest contribution

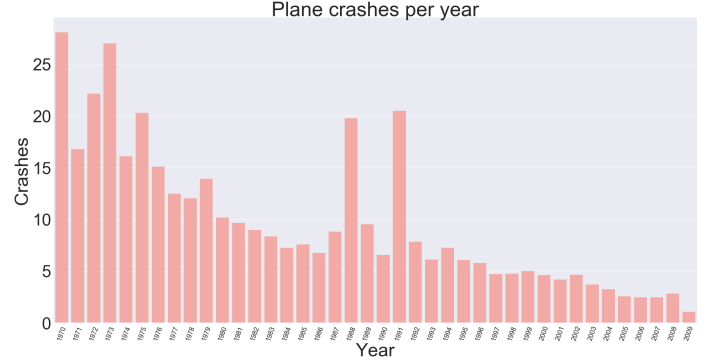


Fig. 1. Proportion of the number of crashes of airplanes (commercial and non-commercial) yearly from 1996 to 2004

is coming from the engine failure then bad weather and many different reasons shown in the figure. Principal component analysis (PCA) using the first two features scatter plot for the six different main causes shows that the engine failure is slightly correlated with other causes which is not surprising as engine failure is mostly from the engine mechanics. Engine failure is also not correlated with the FAA regulations and safety. To investigate more in the matter of the main cause of airplane causes, we looked for the number of crashes across the nations without normalizing these number of crashes. This investigation will show us the effect of FAA regulations on the number of crashes in the USA yearly.

C. Crashes in USA

Figure 3 shows the number for crashes yearly across different nations. The y-axis plotted in log-scale for comparison. It shows that the USA has the largest number of accidents followed by Russia, whereas the USA is one of the leaders in the aviation regulations and safety and the third world country doesn't show much contribution in the figure. The large number of accidents in the USA is directly correlated to the number of flights. It is clear that this large number is directly related to the number of flights and engine failures as the cause of the accidents. Moreover, the crashes are mainly from single-engine airplanes which doesn't have emergency plans for engine failure. Even for multi-engine airplanes, the emergency plans for passengers still not effective these days, but this raise more wide research in terms of what new designs of airplanes that can lead to larger number of survivors in case of airplane crashes, many proposal are done but they still far from reality, the difficulty in this subject is due to the high altitude of the airplane.

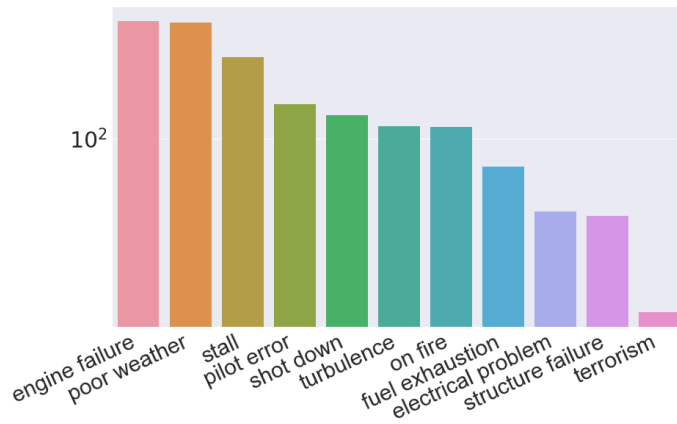


Fig. 2. The main causes of airplane crashes

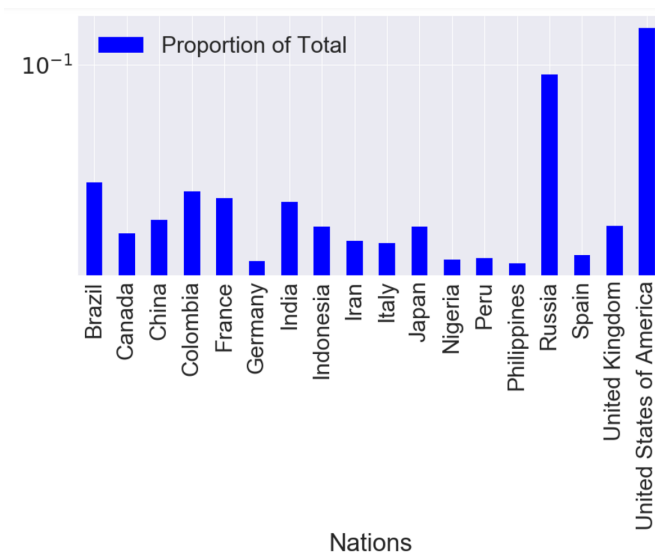


Fig. 3. The main causes of airplane crashes

III. CONCLUSION

We have shown the numbers of accidents and the main causes. The number of accidents are mainly from technology side (engine). The main cause lead to very small number of survivor. The best approach for increasing the number of survivors is also not effective due to the safety followed these days. To decrease the number of crashes and increase the number of survivor, the rules should be regarding airplanes manufacturing and emergency plans that can be followed by passengers other than the current ones. Lastly, to take you far from airplane fearness, the probability to die from car accident is much smaller than the probability to die from airplane accident.

APPENDIX A CODE

```
# coding: utf-8

##### The plan is to normalize these accidents by gathering dataset
for number of among nations yearly, these dataset are not
directly available and we are working on gathering thses data
from the link provide later.

# In[1]:

import numpy as np
import pandas as pd
import matplotlib
from sklearn.linear_model import LinearRegression
from matplotlib import pyplot as plt
get_ipython().magic('matplotlib inline')
from collections import Counter
import seaborn as sns
import re
import operator

# In[2]:

frame=pd.read_csv('Airplane_Crashes_and_Fatalities_Since_1908_1.csv',sep=',')
flights=pd.read_csv('flights.csv') # my own dataset for
normalization, more details will be later

# In[3]:

frame.head()
flights.head()

# In[4]:

frame['Date'] = pd.to_datetime(frame['Date'])
flights_year = flights['YearFlight']
flights_year_crashes = flights['Flights']

# In[5]:

frame['year'] = frame['Date'].dt.year
flights_year.head()
flights_year_norm=flights_year.tolist()
flights_crashes_norm=flights_year_crashes.tolist()
print(flights_crashes_norm)

# In[6]:

frame['Date'] = pd.to_datetime(frame['Date'])
frame['Day'] = frame['Date'].map(lambda x: x.day)
frame['Year'] = frame['Date'].map(lambda x: x.year)
frame['Month'] = frame['Date'].map(lambda x: x.month)

# In[7]:

crashes_per_year = Counter(frame['Year'])
#crashes_per_year = Counter(frame[frame['Year']>2000])
years = list(crashes_per_year.keys())
yearsnew=list()
crashes_year = list(crashes_per_year.values())
crashes_yearnew=list()
for i in range(len(years)):
    if (years[i]) > 1969:
        # print(years[i])
        yearsnew.append(years[i])
        crashes_yearnew.append(crashes_year[i])

# In[8]:

def get_season(month):
    if month >= 3 and month <= 5:
        return 'spring'
    elif month >= 6 and month <= 8:
        return 'summer'
    elif month >= 9 and month <= 11:
        return 'autumn'
    else:
        return 'winter'

frame['Season'] = frame['Month'].apply(get_season)

# In[9]:

crashes_per_season = Counter(frame['Season'])
seasons = list(crashes_per_season.keys())
crashes_season = list(crashes_per_season.values())

# In[10]:

fig = plt.figure(figsize=(30, 30))

sub1 = fig.add_subplot(211)
sns.barplot(x=yearsnew, y=crashes_yearnew, color='r', ax=sub1)
sub1.set(ylabel="Crashes", xlabel="Year", title="Plane_crashes_per_
year")
plt.setp(sub1.patches, linewidth=0)
plt.setp(sub1.get_xticklabels(), rotation=70, fontsize=20)

sub2 = fig.add_subplot(223)
sns.barplot(x=seasons, y=crashes_season, color='g', ax=sub2)
```

```

texts = sub2.set(ylabel="Crashes", xlabel="Season", title="Plane_
crashes_per_season")

plt.tight_layout(w_pad=4, h_pad=1)

# The above plots shows the number of airplanes crashes per year and
season wise. This plots many people have shown these. My main
concern that these plots as it is now, doesnot show much useful
information unless it will be normalized by the number of
flights yearly. Please take a look at the following plots for
this matter

# In[11]:

for i in range(len(yearsnew)):
    for j in range(len(flights_year_norm)):
        if yearsnew[i] == flights_year_norm[j]:
            crashes_yearnew[i] =
            crashes_yearnew[i]/flights_crashes_norm[i]
sns.set_color_codes("pastel")

fig = plt.figure(figsize=(30, 30))
sns.set(font_scale=3)
sub1 = fig.add_subplot(211)
sns.barplot(x=yearsnew, y=crashes_yearnew, color='r', ax=sub1)
sub1.set(ylabel="Crashes", xlabel="Year", title="Plane_crashes_per_
year")
plt.setp(sub1.patches, linewidth=0)
plt.setp(sub1.get_xticklabels(), rotation=70, fontsize=20)

#sub2 = fig.add_subplot(223)
#sns.barplot(x=seasons, y=crashes_season, color='g', ax=sub2)
#texts = sub2.set(ylabel="Crashes", xlabel="Season", title="Plane
crashes per season")

plt.tight_layout(w_pad=4, h_pad=1)

# The above plots, mainly the crashes per year plot, shows the
number of crashes per year nomalized by the number of
passengers yearly as an indication of the number of flight
yearly, the number of passenger yearly was obtained from this
link
http://data.worldbank.org/indicator/IS.AIR.PSGR?end=2015&name_desc=true&start=1970
# The results are from 1970 only as the number of passengers yearly
dataset is only from 1970 and later

# In[12]:

failures = {
    'pilot_error': '(pilot|crew)|(error|fatigue)',
    'engine_failure': 'engine.*(fire|fail)',
    'structure_failure': '(structural|fail)|(fuel|leak)|(langing_
gear)',
    'electrical_problem': 'electrical',
    'poor_weather':
        '(poor|bad).*(weather|visibility)|thunderstorm',
    'stall': 'stall',
    'on_fire': '(caught_fire)|(caught_on_fire)',
    'turbulence': 'turbulence',
    'fuel_exhaustion': '(out_of_fuel)|(fuel.*exhaust)',
    'terrorism': 'terrorist|terrorism',
    'shot_down': 'shot_down',
}

failure_counts = {'other':0}

for s in frame.Summary.dropna():
    other = True
    for failure, exp in failures.items():
        if re.search(exp, s.lower()):
            other = False
            if failure in failure_counts:
                failure_counts[failure] += 1
            else:
                failure_counts[failure] = 1
    if other:
        failure_counts['other'] += 1

nan_counts = len(frame.Summary.isnull())

del failure_counts['other']

sortedcauses = sorted(failure_counts.items(),
    key=operator.itemgetter(1), reverse=True)
for k, v in sortedcauses:
    print(k, v)

plt.figure(figsize=(14, 8))
x, y = zip(*sortedcauses)
sns.barplot(x=x, y=y)
plt.xticks(rotation=25, horizontalalignment='right')
plt.yscale('log')
plt.show()

# In[13]:

#Splitting out the country from the location to see if we can find
some interesting
#statistics about where the most crashes have taken place.
operator =
    frame[['Operator', 'Fatalities']].groupby('Operator').agg(['sum', 'count'])
fatalities =
    operator['Fatalities', 'sum'].sort_values(ascending=False)
#fig_ops, (ax1, ax2), (ax3, ax4)=plt.subplots(2,2,sharex=True)
#fatalities.head(10).plot(kind='bar', title='Fatalities by
Operator', ax=ax3, grid=True, rot=90)
totalfatal = fatalities.sum()
s = frame['Location'].str[0:].str.split(',', expand=True)
frame['Country'] = s[3].fillna(s[2]).fillna(s[1]).str.strip()
#I've pulled out all the US states so as to be able to assign them a
country

usNames = ['Virginia', 'New Jersey', 'Ohio', 'Pennsylvania',
'Maryland', 'Indiana', 'Iowa',
'Illinois', 'Wyoming', 'Minnesota', 'Wisconsin', 'Nevada',
'NY', 'California',
'WY', 'New York', 'Oregon', 'Idaho',
'Connecticut', 'Nebraska', 'Minnesota', 'Kansas',
'Texas', 'Tennessee', 'West Virginia', 'New Mexico',
'Washington', 'Massachusetts',
'Utah', 'Illinois', 'Florida', 'Michigan',
'Arkansas', 'Colorado', 'Georgia', 'Missouri',
'Montana', 'Mississippi', 'Alaska', 'Jersey', 'California',
'Oklahoma', 'North Carolina',
'Kentucky', 'Delaware', 'D.C.', 'Arizona', 'Arizona', 'South
Dakota', 'New Hampshire', 'Hawaii',
'Washington', 'Massachusetts', 'Washington',
DC', 'Tennessee', 'Delaware', 'Louisiana',
'Massachusetts', 'Louisiana', 'New York',
(Idledild)', 'Oklahoma', 'North Dakota', 'Rhode Island',
'Maine', 'Alaska', 'Wisconsin', 'California', 'Virginia', 'Virginia', 'CA', 'V
HI', 'AK', 'IN', 'GA', 'Colorado', 'Arizona', 'Alabama', 'Alaska'
]

#Decided to try and cleanse the country names.
afNames = ['Afghanistan'] #Afghanistan
anNames = ['Angola'] #Angola
ausNames = ['Qld', 'Australia', 'Queensland', 'Australia', 'Tasmania', 'off_
Australia'] #Australia
argNames = ['Aregentina'] #Argentina
azNames = ['Azores', 'Portugal'] #Azores
baNames = ['Baangladesh'] #Bangladesh
bahNames = ['Great Inagua'] #Bahamas
berNames = ['off_Bermuda'] #Bermuda
bolNames = ['Boliva', 'BO'] #Bolivia
bhNames = ['Bosnia-Herzegovina'] #Bosnia Herzegovina
bulNames = ['Bugaria', 'Bulgaria'] #Bulgaria
canNames = ['British Columbia', 'British Columbia_Canada', 'Canada2',
'Saskatchewan', 'Yukon Territory'] #Canada
camNames = ['Cameroon', 'French Cameroons'] #Cameroon
caNames = ['Cape Verde Islands'] #Cape Verde
chNames = ['Chili'] #Chile
coNames = ['Comoro Islands', 'Comoros Islands'] #Comoros
djNames = ['Djibouti', 'Republic of Djibouti'] #Djibouti
domNames = ['Dominican Republic', 'Dominica'] #Dominican Republic
drcNames = ['Belgian Congo', 'Belgian Congo (Zaire)', 'Belgium Congo',
'Democratic Republic of Congo', 'Democratic Republic of Congo',
'Democratic Republic of Congo', 'Democratic Republic of Congo',
'Democratic Republic of Congo', 'Zaire',
'Zaire'] #Democratic Republic of Congo
faNames = ['French Equatorial Africa'] #French Equatorial Africa
gerNames = ['East Germany', 'West Germany'] #Germany
grNames = ['Crete'] #Greece
haNames = ['Hati'] #Haiti
hunNames = ['Hungary'] #Hungary
inNames = ['Indian'] #India
indNames = ['Indonesia', 'Netherlands Indies'] #Indonesia
jamNames = ['Jamaica'] #Jamaica
malNames = ['Malaya'] #Malaysia
manNames = ['Manmar'] #Myanmar
marNames = ['Mauretania'] #Mauritania
morNames = ['Morocco', 'Morocco'] #Morocco
nedNames = ['Amsterdam', 'The Netherlands'] #Netherlands
niNames = ['Niger'] #Nigeria
phiNames = ['Philippines', 'Philippine Sea', 'Philippines',
'off_the_Philippine_Island_of_Elalat'] #Philippines
romNames = ['Romania'] #Romania
rusNames = ['Russian', 'Soviet Union', 'USSR'] #Russia
saNames = ['Saint Lucia Island'] #Saint Lucia
samNames = ['Western Samoa'] #Samoa
siNames = ['Sierra Leone'] #Sierra Leone
soNames = ['South Africa (Namibia)'] #South Africa
surNames = ['Suriname'] #Surinam
uaeNames = ['United Arab Emirates'] #UAE
ukNames = ['England', 'UK', 'Wales', '110 miles West of Ireland']
#United Kingdom
uvNames = ['US Virgin Islands', 'Virgin Islands'] #U.S. Virgin Islands
wkNames = ['325 miles east of Wake Island'] #Wake Island
yuNames = ['Yugoslavia'] #Yugoslavia
zimNames = ['Rhodesia', 'Rhodesia (Zimbabwe)'] #Zimbabwe

cnames = []
for country in frame['Country'].values:
    if country in afNames:
        cnames.append('Afghanistan')
    elif country in anNames:
        cnames.append('Angola')
    elif country in ausNames:
        cnames.append('Australia')
    elif country in argNames:
        cnames.append('Argentina')
    elif country in azNames:
        cnames.append('Azores')
    elif country in baNames:
        cnames.append('Bangladesh')
    elif country in bahNames:
        cnames.append('Bahamas')
    elif country in berNames:
        cnames.append('Bermuda')
    elif country in bolNames:
        cnames.append('Bolivia')
    elif country in bhNames:
        cnames.append('Bosnia_Herzegovina')
    elif country in bulNames:
        cnames.append('Bulgaria')
    elif country in canNames:
        cnames.append('Canada')
    elif country in camNames:
        cnames.append('Cameroon')
    elif country in caNames:
        cnames.append('Cape Verde')
    elif country in chNames:
        cnames.append('Chile')
    elif country in coNames:
        cnames.append('Comoros')
    elif country in djNames:
        cnames.append('Djibouti')

```

```

elif country in domNames:
    cnames.append('Dominican_Republic')
elif country in drcNames:
    cnames.append('Democratic_Republic_of_Congo')
elif country in faNames:
    cnames.append('French_Equatorial_Africa')
elif country in gerNames:
    cnames.append('Germany')
elif country in grNames:
    cnames.append('Greece')
elif country in haNames:
    cnames.append('Haiti')
elif country in hunNames:
    cnames.append('Hungary')
elif country in inNames:
    cnames.append('India')
elif country in jamNames:
    cnames.append('Jamaica')
elif country in malNames:
    cnames.append('Malaysia')
elif country in manNames:
    cnames.append('Myanmar')
elif country in marNames:
    cnames.append('Mauritania')
elif country in morNames:
    cnames.append('Morocco')
elif country in nedNames:
    cnames.append('Netherlands')
elif country in niNames:
    cnames.append('Nigeria')
elif country in philNames:
    cnames.append('Philippines')
elif country in romNames:
    cnames.append('Romania')
elif country in rusNames:
    cnames.append('Russia')
elif country in saNames:
    cnames.append('Saint_Lucia')
elif country in samNames:
    cnames.append('Samoa')
elif country in siNames:
    cnames.append('Sierra_Leone')
elif country in soNames:
    cnames.append('South_Africa')
elif country in surNames:
    cnames.append('Surinam')
elif country in uaeNames:
    cnames.append('UAE')
elif country in ukNames:
    cnames.append('United_Kingdom')
elif country in usNames:
    cnames.append('United_States_of_America')
elif country in uvNames:
    cnames.append('U.S._Virgin_Islands')
elif country in wkNames:
    cnames.append('Wake_Island')
elif country in yuNames:
    cnames.append('Yugoslavia')
elif country in zimNames:
    cnames.append('Zimbabwe')
else:
    cnames.append(country)

plt.figure(figsize=(14, 8))
frame['Nations'] = cnames
fatalcountries =
    frame[['Fatalities', 'Nations']].groupby(['Nations']).agg('sum')
fatalcountries.reset_index(inplace = 'True')
fatalcountries['Proportion_of_Total'] =
    fatalcountries['Fatalities']/totalfatal

fig_c, (ax1, ax2) = plt.subplots(2,1,sharex = True)

fatalcountries[fatalcountries['Fatalities']>1300].plot(x = 'Nations'
    , Y = 'Fatalities'
    , ax = ax1
    , kind = 'bar'
    , grid = True,
    colormap='winter',
    logy=True,
    figsize=(14,14),
    sort_columns=True)

fatalcountries[fatalcountries['Fatalities']>1300].plot(x = 'Nations'
    , Y = 'Proportion_
    of_Total'
    , ax = ax2
    , kind = 'bar'
    , grid = True,
    colormap='winter',
    logy=True)

# These two plots was also shown before, what we will be doing next,
# is taking the largest two location for accidents, Russia and
# United State of America and gather dataset to normalize these
# by the number of flights in these two countires.

# ##### The plan is to normalize these accidents by gathering
# dataset for number of among nations yearly, these dataset are
# not directly avialable and we are working on gathering thses
# data from the above link provide.

```

REFERENCES

- [1] <https://www.kaggle.com/saurograndi/airplane-crashes-since-1908>.
- [2] <http://data.worldbank.org/indicator/IS.AIR.PSGR?end=2015>