

Aircrafts are becoming safer: An analysis of air crashes and flight safety

Kaiwei Qian

2019/5/4

1 Introduction

There are 6 sections in this report. Data and the ways how they are collected from Internet are described in Section 2. Exploratory analysis is performed in Section 3. In Section 4, the probability of being involved in an air accident is studied for the past 10 years (2009-2018), followed by Section 5, where text mining techniques are tried to extract the keywords from the cause of air accidents. The last section gives the conclusion of the project.

2 Data

2.1 Aircraft accidents dataset

The data used in this study is scraped from the **Bureau of Aircraft Accidents Archives** (<http://www.baaa-acro.com/crash-archives>) for the past 10 years (2009-2018). The Bureau of Aircraft Accidents Archives (BAAA) was established in Geneva in 1990 for the purpose to deal with all information related to aviation accidentology.

Scraping involves two steps.

First, After `read_html()` reads the string-type HTML from the response, `html_node(xpath = '/html/body/div[1]/div[4]/div/section/div[2]/div/div/div[3]/div/table')` locates the table that contains the information of the air accidents, such as Date, A/C Type, Location, Fatalities and Registration. Detailed explanation if these headings are described in Table 1.

If you visit their website, you'll find that there is a plus sign at the end of each row of the table. Clicking on it leads you to the page that briefly summarizes the air crash. It provides more detailed information, such as Flight Phase and Circumstances. So, I choose to keep the hyperlinks to this page. `html_nodes(xpath = '/html/body/div[1]/div[4]/div/section/div[2]/div/div/div[3]/div/table/tbody/tr/td[8]/a/@href')` locates the column of the hyperlinks belonging to and `html_text()` reads the data from that column.

In the second step, the spider visits each page that contains a brief summary of air crashes. I choose to keep the records of Flight Phase, Flight Type, Site, Circumstances and Probable Causes (if exists) for further study. Circumstances and Probable Causes are mainly used to support the text mining part of this project. XPath is used to read the data from the HTML.

Note that instead of air crashes, this dataset include all the records of air accidents, which means the plane may not be necessarily hit land or water and be damaged. For example, some mechanical problems during the flight are also deemed as air accidents.

Table 1: Description of the variables in BAAA Air Accidents dataset.

| Variable Name | Data Type | Description |
|---------------|-----------|-------------------------------------|
| Date | Date | The date when the accident occurred |

| Variable Name | Data Type | Description |
|----------------|-----------|--|
| A/C Type | character | The type of the aircraft |
| Location | character | The location where the accident occurred |
| Fatalities | numeric | Number of deaths caused by the accident |
| Registration | character | A unique ID for the airplane |
| Flight Phase | character | Flight phase during which the accident occurred |
| Circumstances | character | A brief description of the circumstances where the accident occurred |
| Flight Type | character | The type of this flight; can be private, commercial and so on |
| Site | character | Description of the location where the air crash occurred |
| Probable cause | character | Possible causes given by the investigation by NTSB |

2.2 Flight statistics

The flight statistics dataset is scraped from **Statista**. The number of flights performed by the global airline industry from 2009 to 2018 and the number of scheduled passengers boarded by the global airline industry from 2009 to 2018 are collected from <https://www.statista.com/statistics/564769/airline-industry-number-of-flights> and <https://www.statista.com/statistics/564717/airline-industry-passenger-traffic-globally>, respectively.

In each of websites mentioned above, there are two tables, which contain the number of scheduled commercial flights and the number of passengers and the number of scheduled passengers boarded from 2004 to 2019, respectively. For each of the websites, `html_node(xpath = '//table')` locates the table, and `html_table()` reads the table from its HTML structure.

Only the records from 2009 to 2018 are used in this project in order to keep them consistent with the air accidents dataset in terms of time range. Moreover, the numbers from **Statista** are only for the commercial flights while the air accidents dataset also contain the records for private, cargo and other flights. So, it's essential to keep the difference of these two data sources in mind and avoid causing confusion.

3 Exploratory Analysis

3.1 Aircraft Type

As shown in Table 1, the air accidents dataset contains 1343 observations and 10 variables. The 10 variables are Date, A/C Type, Location, Fatalities, Registration, Flight Phase, Circumstances, Flight Type, Site, Probable cause.

Barplots are employed to see what types of aircraft is involved in most air accidents and what types of aircraft cause the most fatalities. After grouping the data by **A/C Type**, I sum the **Fatalities** and count the number of air accidents for each type of aircraft. `gather()` and `facet_wrap()` are used to generate these two plots.

Top 10 aircraft types are shown in the graphs below. **Boeing 777-200** have the most deaths and **PZL-Mielec AN-2** has been involved in the most accidents in the past 10 years. Interestingly, the A/C Type which causes hundreds of deaths are not among the types which are most frequently involved in accidents, except for **Lockheed C-130 Hercules**. One possible reason for this phenomenon is that **Boeing** and **Airbus** produce commercial airplanes that can contain hundreds of people and a handful of fatal air crashes can lead to hundreds of fatalities.

In conclusion, the capacity of different aircrafts must be taken in to consideration when talking about the safety.

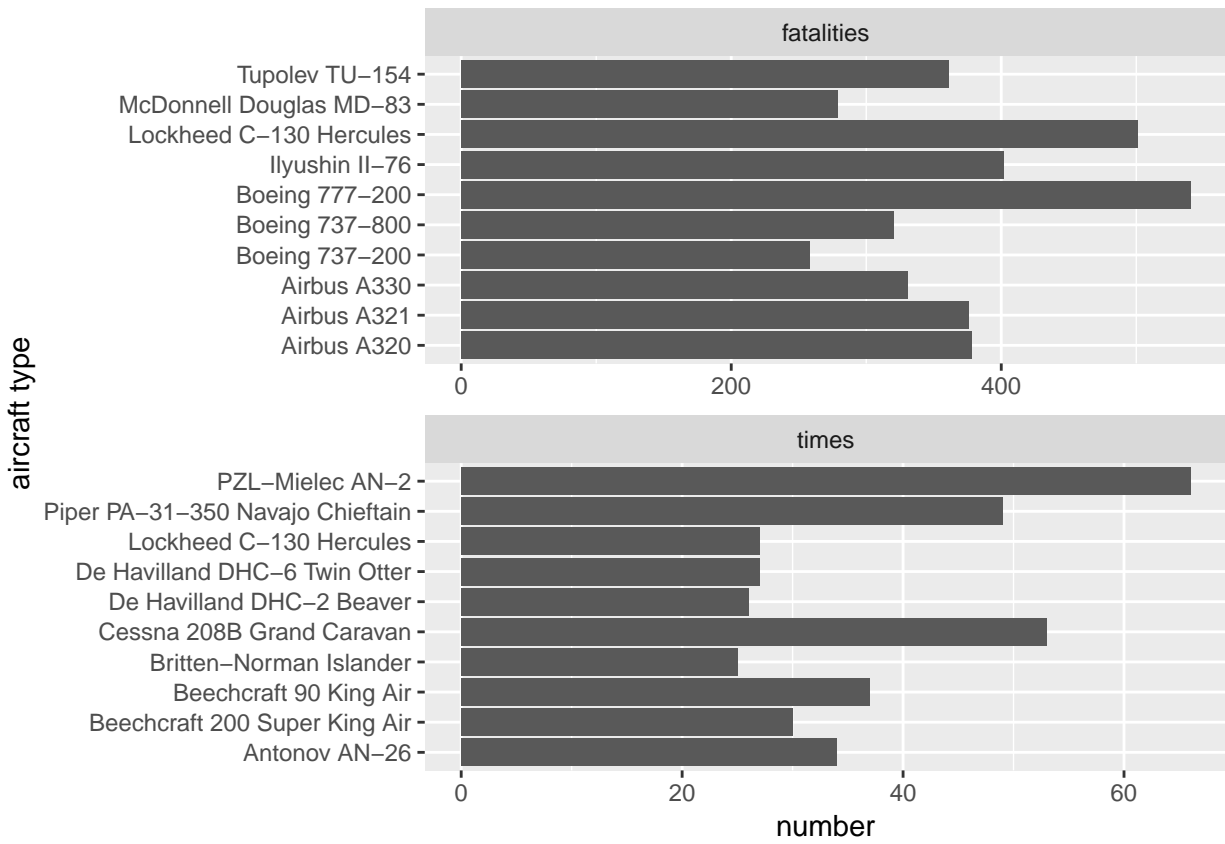


Figure 1: Frequency and Fatalities of Air Accidents by Aircraft Type.

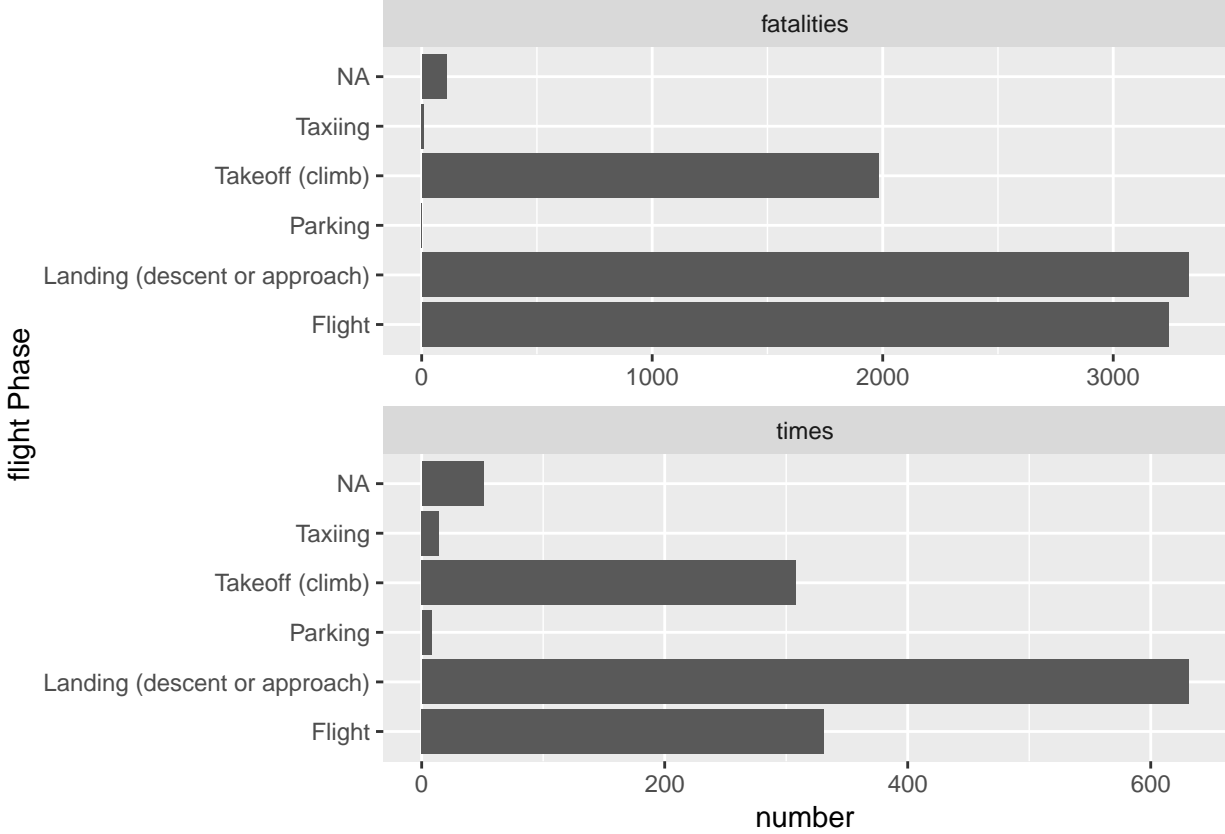


Figure 2: Frequency and Fatalities of Air Accidents by Flight Phase.

3.2 Flight Phase

Barplots in Figure 2 shows that nearly half of the accidents that occur during the landing phase. Though only 1/4 of the accidents happen during the flight, the fatalities caused in this phase are as many as that those caused during the landing phase.

When the plane is landing, it has lower speed and lower altitude. So, even the airplane accidents during the landing or approaching phase occur more often than during flight, the chance of surviving must be higher than the accident in the sky. What's more, landing or approaching means that the plane is near an airport, where the people are ready for any emergency situation.

3.3 Location

The number of air crashes and fatalities by location is computed in this part.

4 Time series

In this section, records are grouped by year, and I'm interested in the trend of the fatalities and number of air accidents from 2009 to 2018. In addition, the number of scheduled flights and the number of passengers boarded are used to normalize the air accidents dataset in order to achieve fair evaluation of flight safety.

Figure 3 indicates that the fatalities related to air accidents fluctuate around 800 per year, and the number of air accidents has decreased from 160 to 110 from 2009 to 2018. The problem is whether fluctuating fatalities

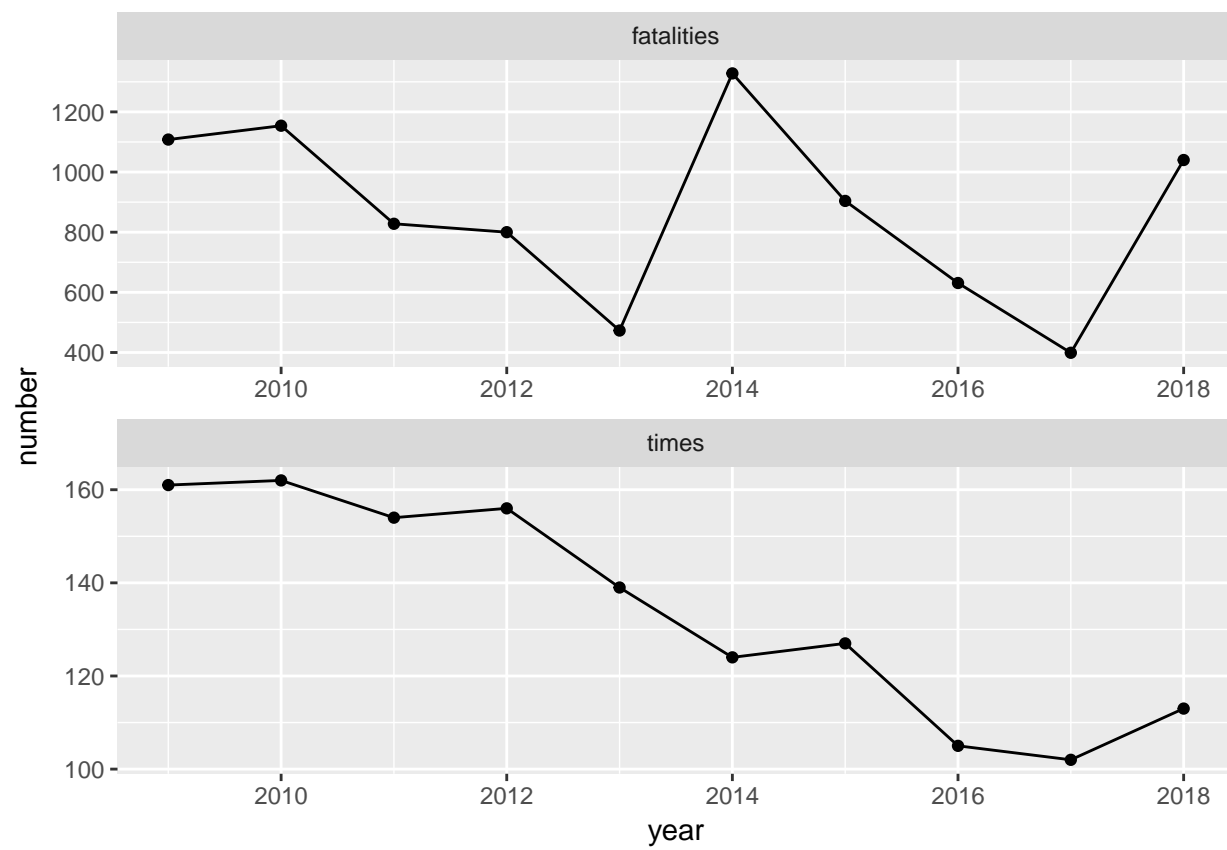


Figure 3: Number of Air Accidents and Fatalities per Year

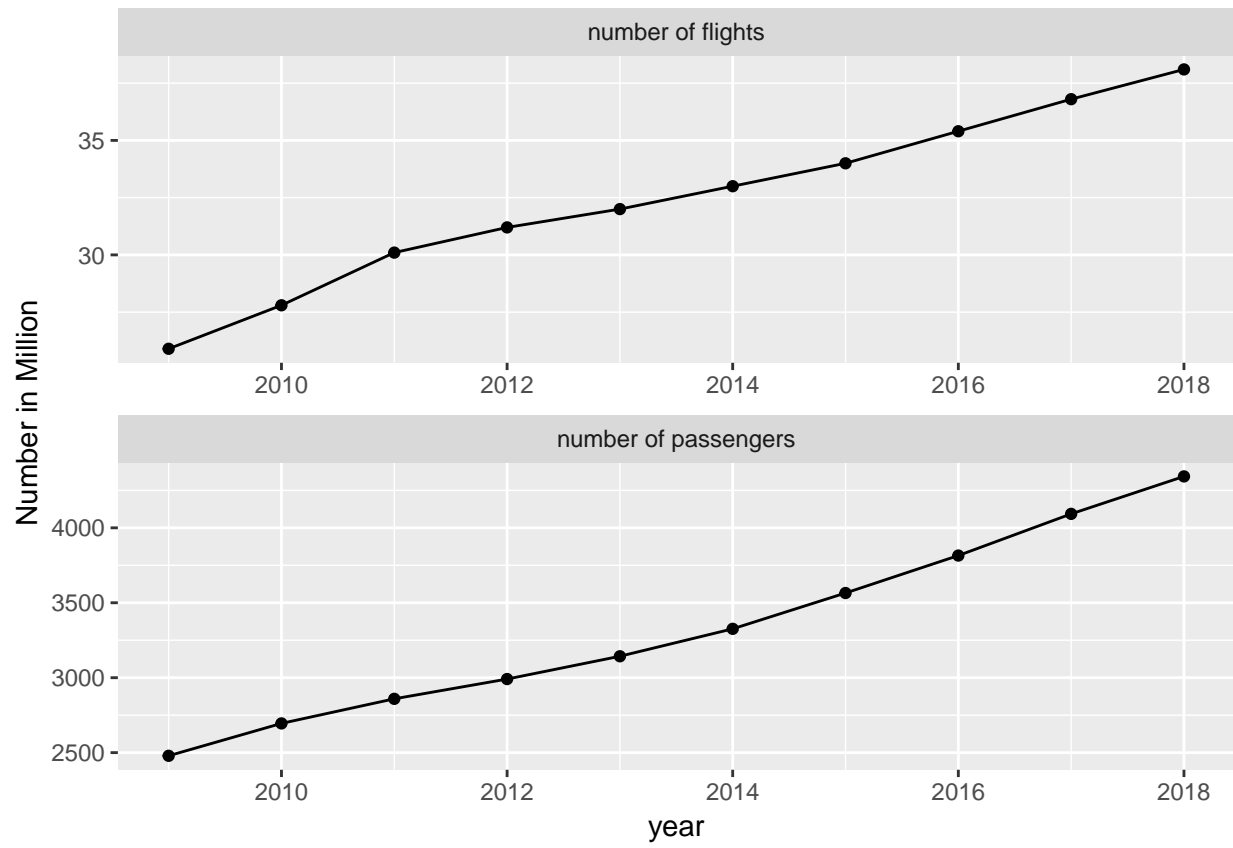


Figure 4: Number of Commercial Flights and Passengers per Year.

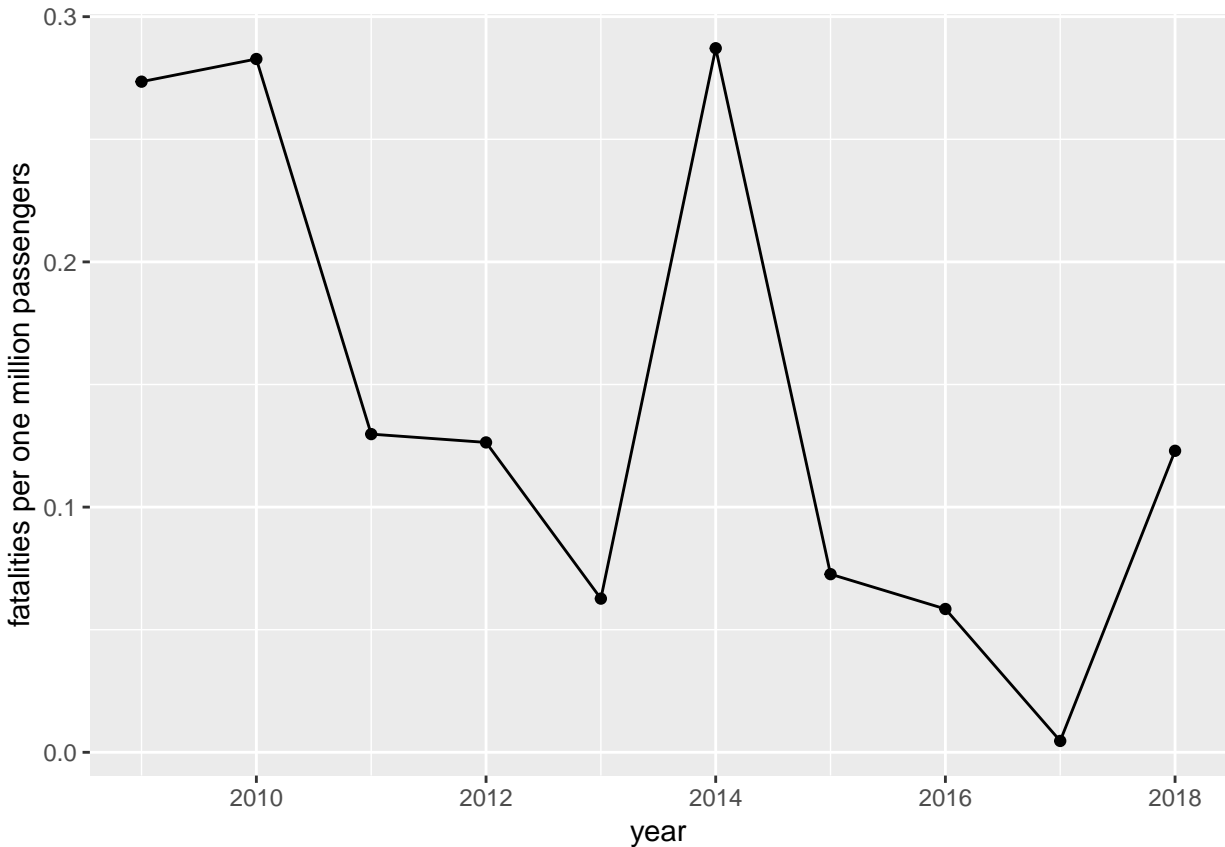


Figure 5: Fatalities per Passengers Boarded per Year.

caused by air crashes each year indicates no improvement in flight safety.

Figure 4 shows that the number of commercial flights scheduled and number of passengers each year has been increasing in the past ten years. Therefore, it's more appropriate to evaluate the flight safety by the ratio of fatalities to the number of passengers aboard or the **fatal accident rate**, which is the ratio of number of fatal air accidents to the total number of flights.

Note that the data of the number of flights and passengers are only for commercial flights. Therefore, only the air accidents with the flight type of "Scheduled Revenue Flight" can be used to calculate those ratios for commercial flights.

4.1 fatalities per passengers boarded

Figure 5 shows a downward trend for the number of victims per one million passengers for commercial flights. 2014 and 2018 are two abnormal years because the number of fatalities do not follow the overall trend for these two years. In 2018, there were two major air crashed that caused more than 300 deaths. One of them is the Boeing 737 MAX 8 that crashed in Jakarta and killed 189 people. In 2014, there were three air crashes that involved Boeing and Airbus. These three major accidents killed 162, 298 and 239, respectively.

4.2 fatal accident rate

Figure 6 indicates that the number of air accidents that caused deaths per one million flights has been decreasing in the past ten years. However, the decreasing fatal accidents rate does not necessarily mean that

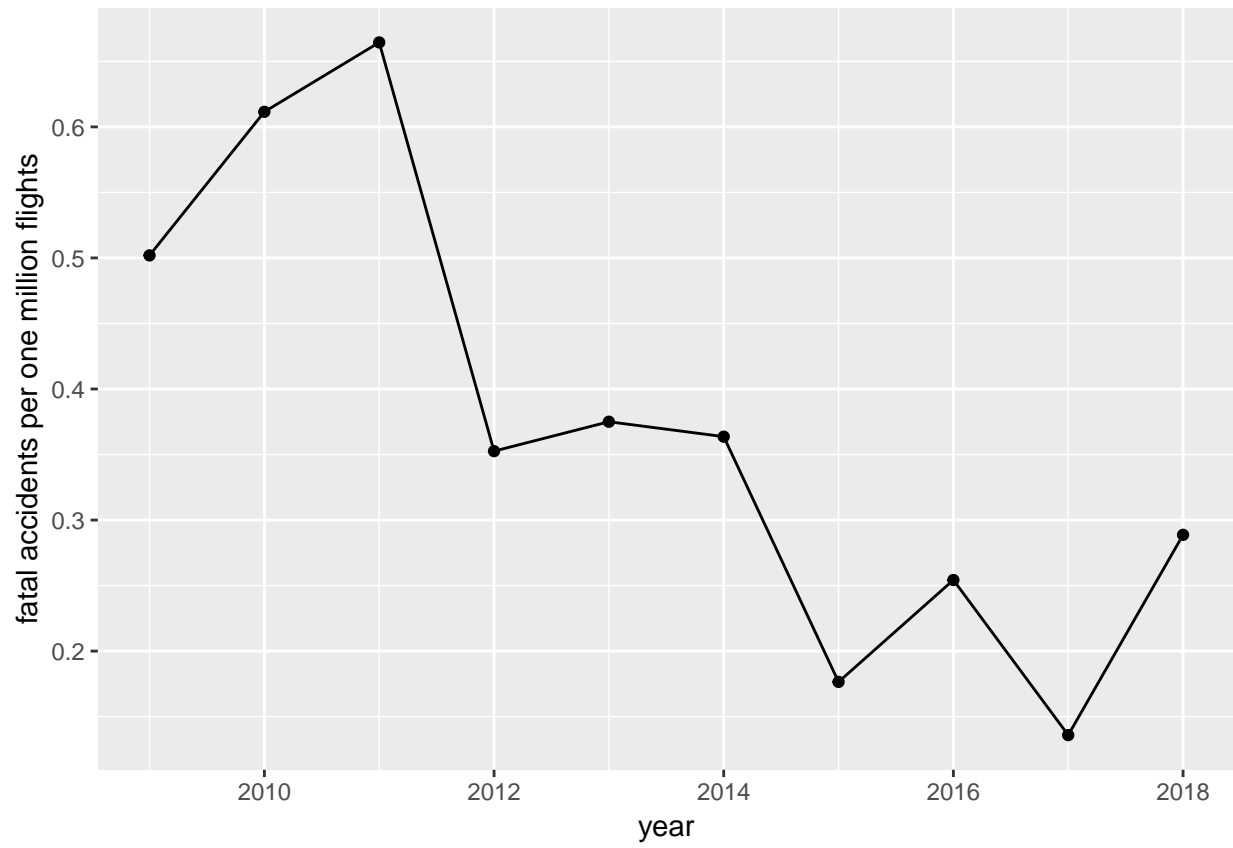


Figure 6: Fatal Accident Rate per Year.

the number of victims related to air crash also goes down. It's clear that Figure 6 fails to capture the extreme high fatalities in 2014 and 2018. Though the number of crashes was small in 2014 and 2018, those accidents involved high-capacity commercial airplanes.

In conclusion, except for some extreme years, the commercial flights have become safer over the past ten years. In addition, the number of fatal accidents keeps going down. Air crashes that involve high-capacity planes, such as Boeing and Airbus, lead to significant increase in the fatalities.

5 Text Mining

BAAA provides information of the circumstances when the air accident happened and there would be one section called "Probable Causes" if the investigation was finished. I'm interested in the key factors that lead to a fatal air accident and decide to use text mining techniques to find them. Naive term frequency and TF-IDF are both tried and the results are compared in this part.

5.1 Term Frequency

5.2 TF-IDF

6 Conclusion