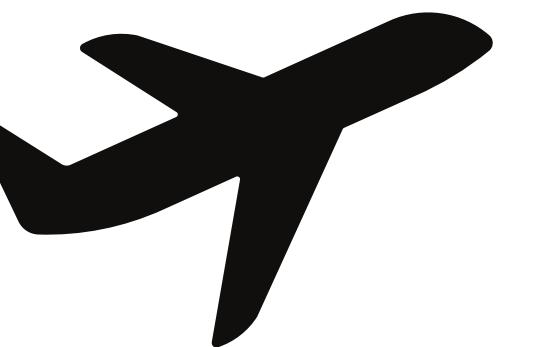


# SC1015 MINI- PROJECT



DATASET: AVIATION INCIDENT

## GROUP 3

MEMBERS: 1. Koh kai jie, Farrel  
2. Wong Xiang Rui  
3. Quek Wei Quan



# Table of Contents

- I Research Background & Motivation
- II Cleaning of Data set
- III Explorative Data Analysis
- IV Machine learning
- V Conclusion & insights  
driven



# BACKGROUND

## Latest plane crash news

[Top Stories](#) [Latest News](#) [Discover](#) [Singapore](#) [Asia](#) [Commentary](#) [Sustainability](#) [CNA Insider](#) [Lifestyle](#) [Watch](#) [Listen](#) [+ All Sections](#)

World

### DHL plane breaks in two during emergency landing in Costa Rica airport



A DHL cargo plane broke in two after an emergency landing at the Juan Santamaria International Airport in Costa Rica, on April 7, 2022 (Photo: AFP/Ezequiel Becerra)

SAN JOSE: A DHL cargo plane carrying mail and packages skidded off the runway and broke in two during an emergency landing in Costa Rica on Thursday, causing the temporary closure of the international airport in San Jose.

**South China Morning Post**

China Eastern Airlines flight MU5735 crash

+ FOLLOW

China / Politics

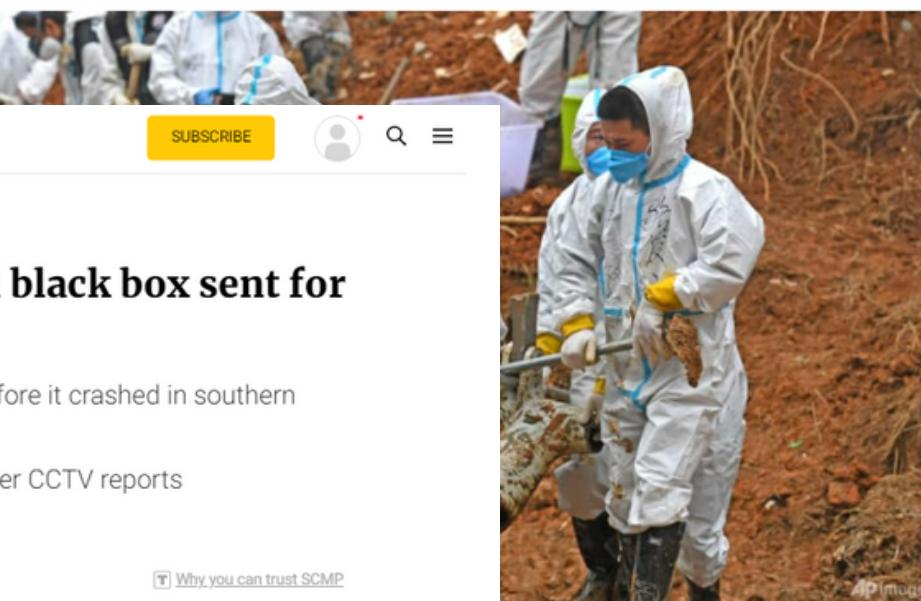
### China Eastern Airlines flight MU5735: second black box sent for decoding

- The plane's data recorder could reveal details of its speed and altitude before it crashed in southern China
- The black box was found buried under 1.5 metres of soil, state broadcaster CCTV reports



Luna Sun and Jack Lau

Published: 11:09am, 27 Mar, 2022



27 Mar 2022 11:41AM  
(Updated: 27 Mar 2022 07:36PM)



TOP PICKS

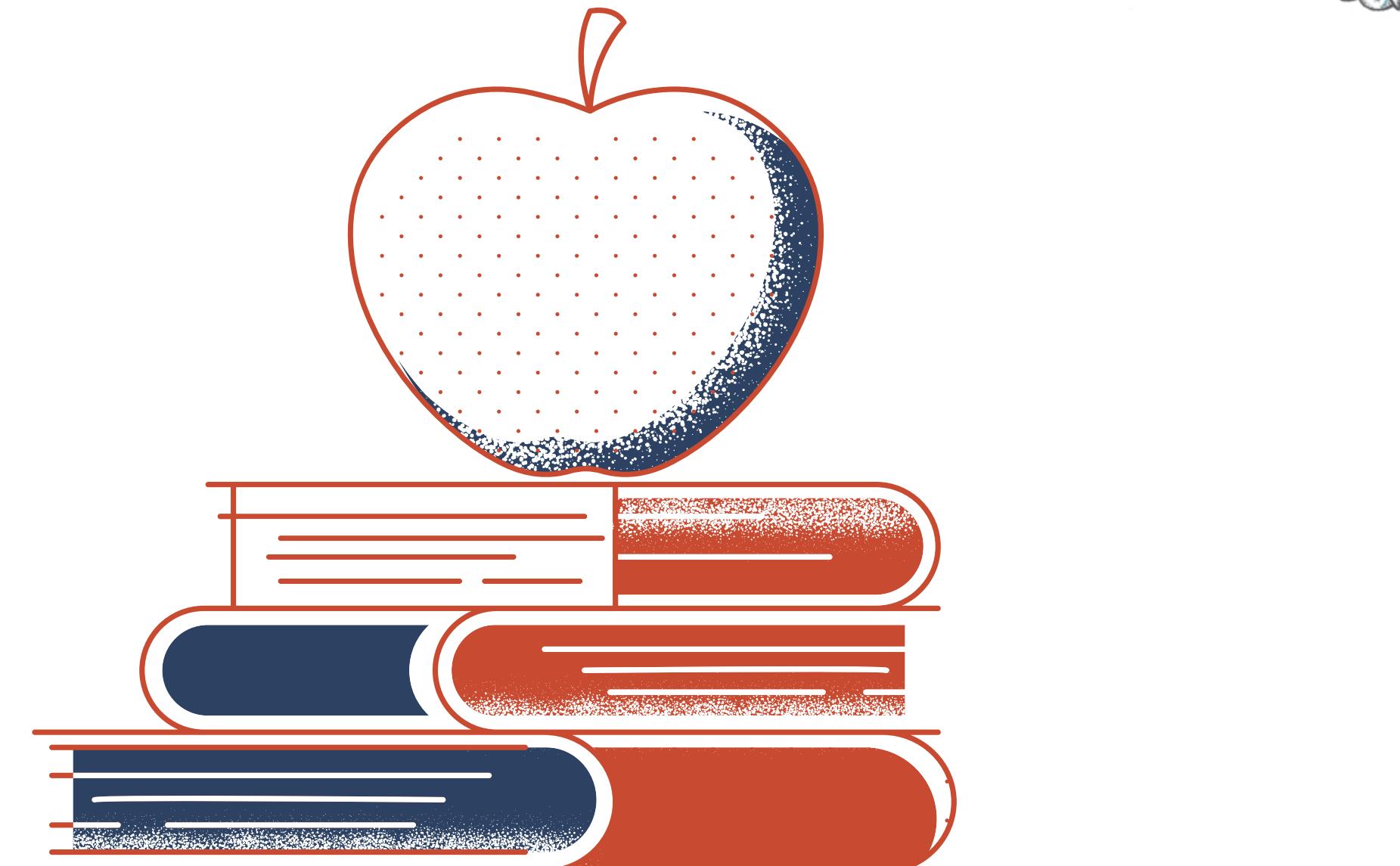
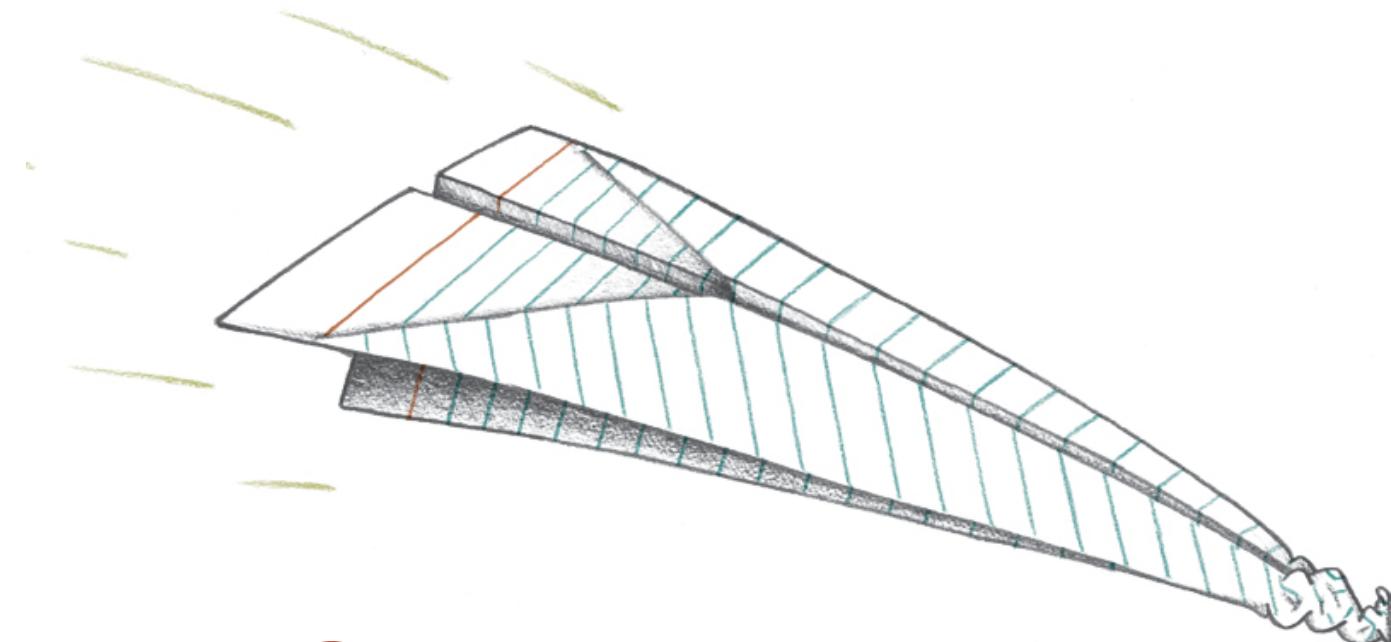
281



# Research Background & Motivation

## **PROBLEM STATEMENT**

To identify the factor that  
has the most impact on  
aviation incident.



# CLEANING DATA SET

## Selecting appropriate data:

- Selecting data columns with the least amount of missing values
- Data that have a great influence on answering our problem statement
- Data inputs that could be used for machine learning

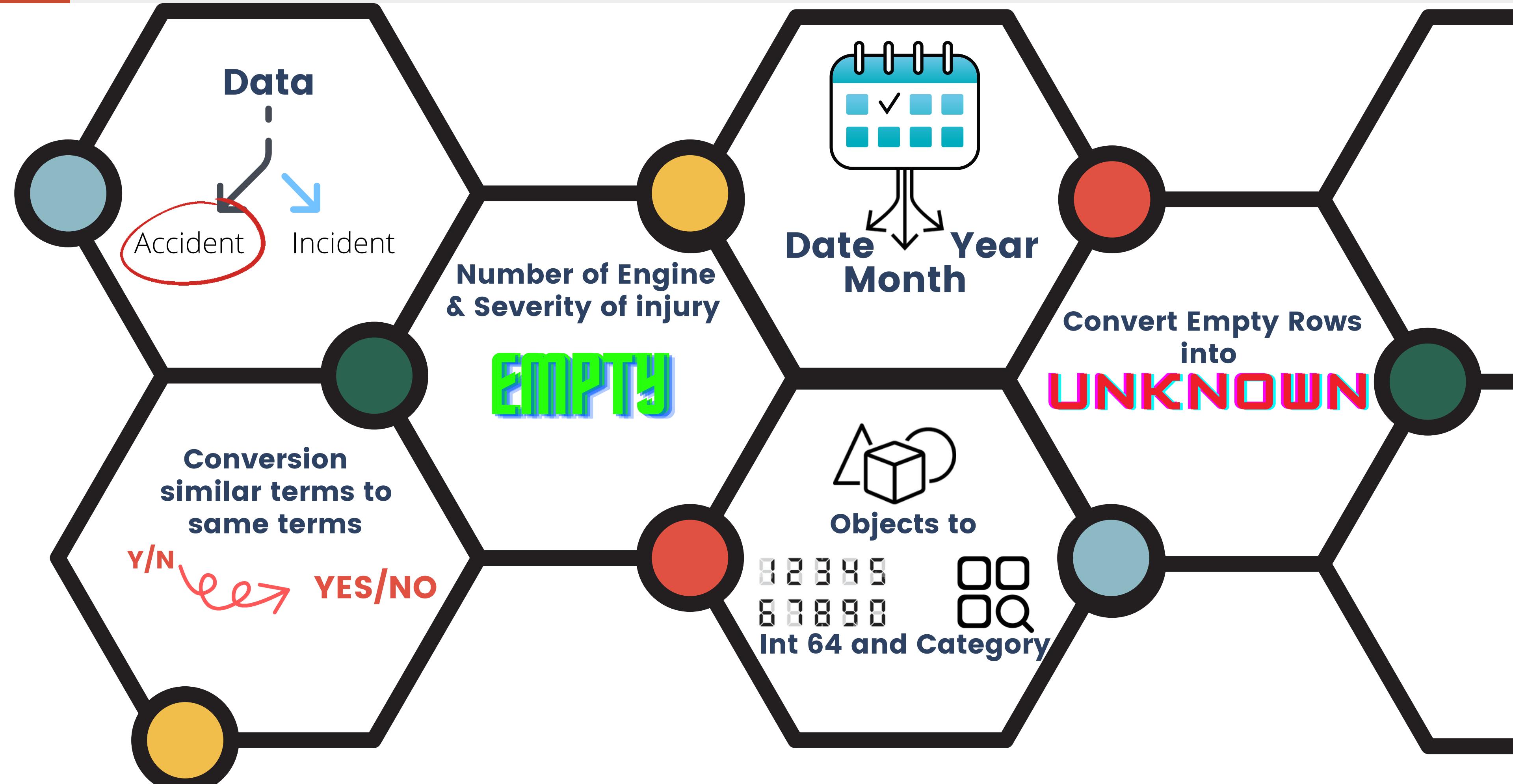


## Selected Data sets

- Event Date
- Location
- Country
- Company
- Number of Engines
- Injury Severity
- Total Fatal Injuries
- Weather Condition
- Broad phase of flight
- Investigation Type



# CLEANING DATA SET



## Over removal of data



- This could result in false data appearing
- Data might appear too "fake" after confusion matrix is done

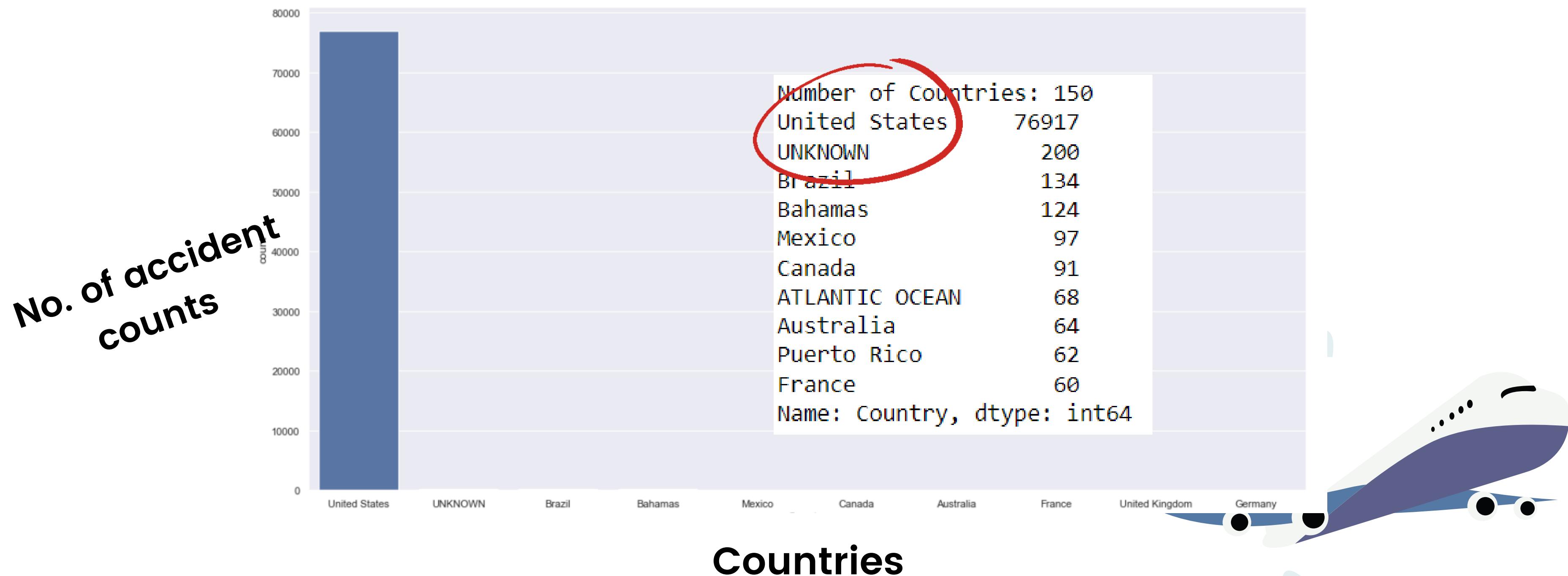


## Removal of selected data set

- Other variables such as the phase of flight were considered.
  - However, they were removed due to the high number of unknowns.

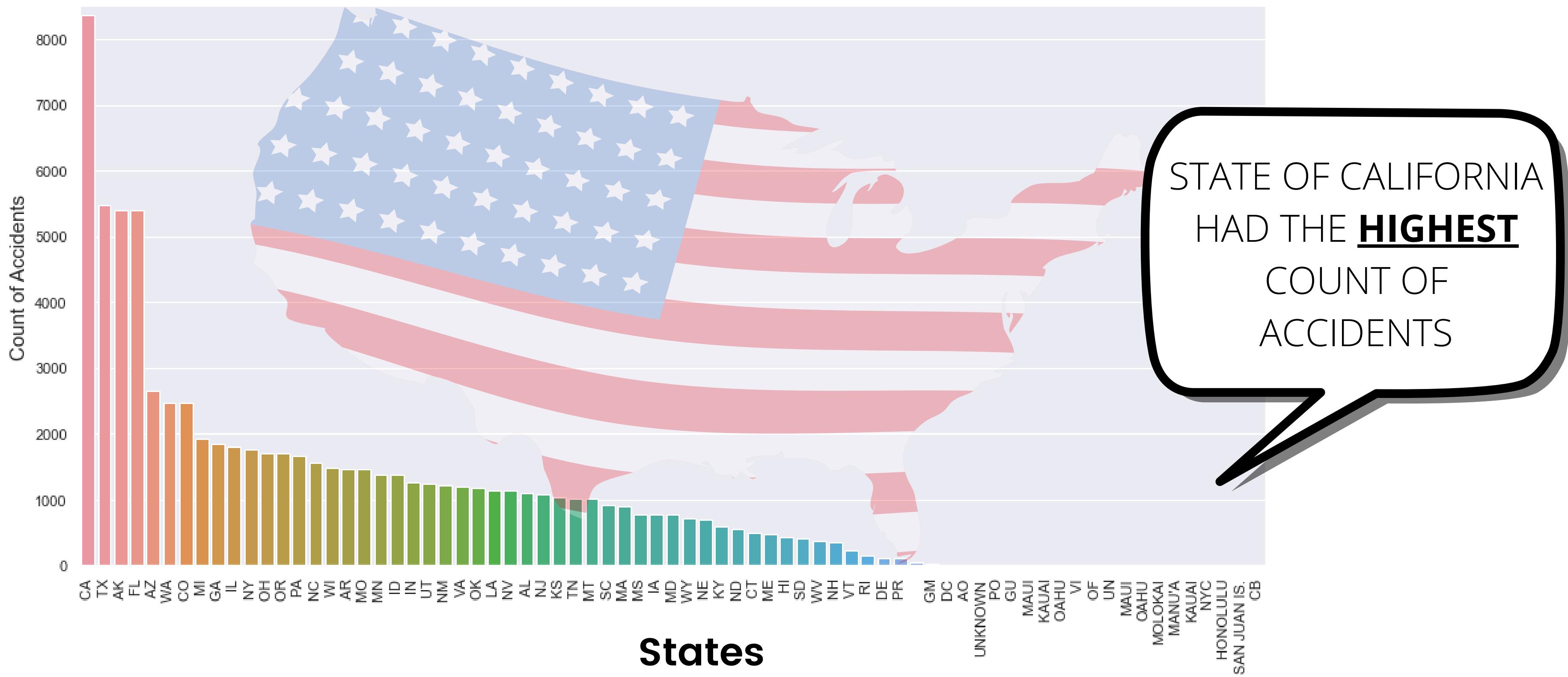


# Top 10 countries with the highest number of accidents

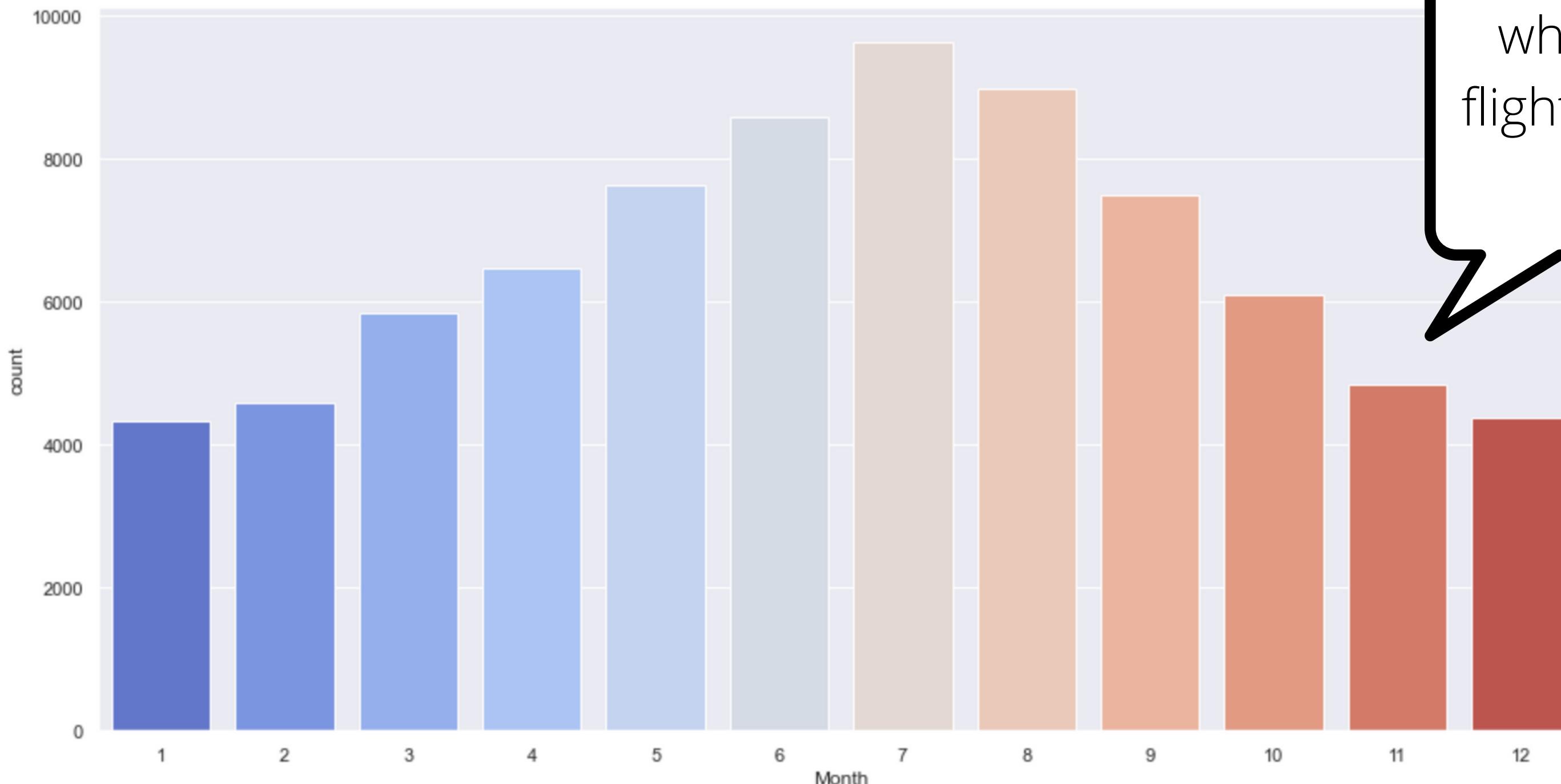


# Explorative Data Analysis

## No. of Accident occurred for each state in the US



## Time Series of Injuries



There is a correlation to the summer season when there are more flights during this period

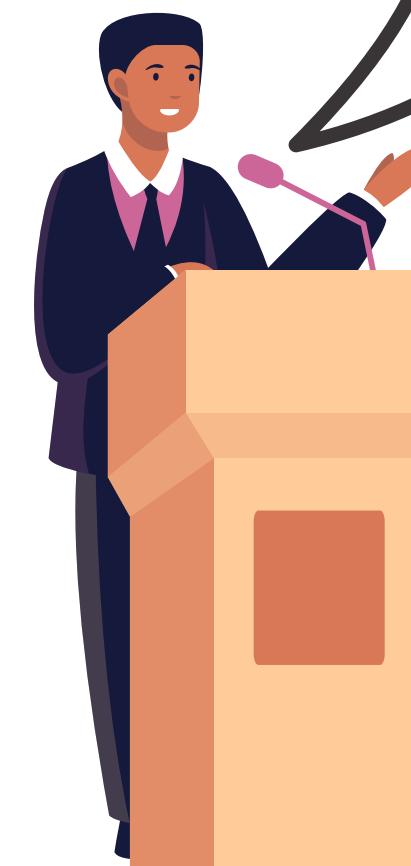


# Explorative Data Analysis

## Top 10 Companies with the highest number of accidents

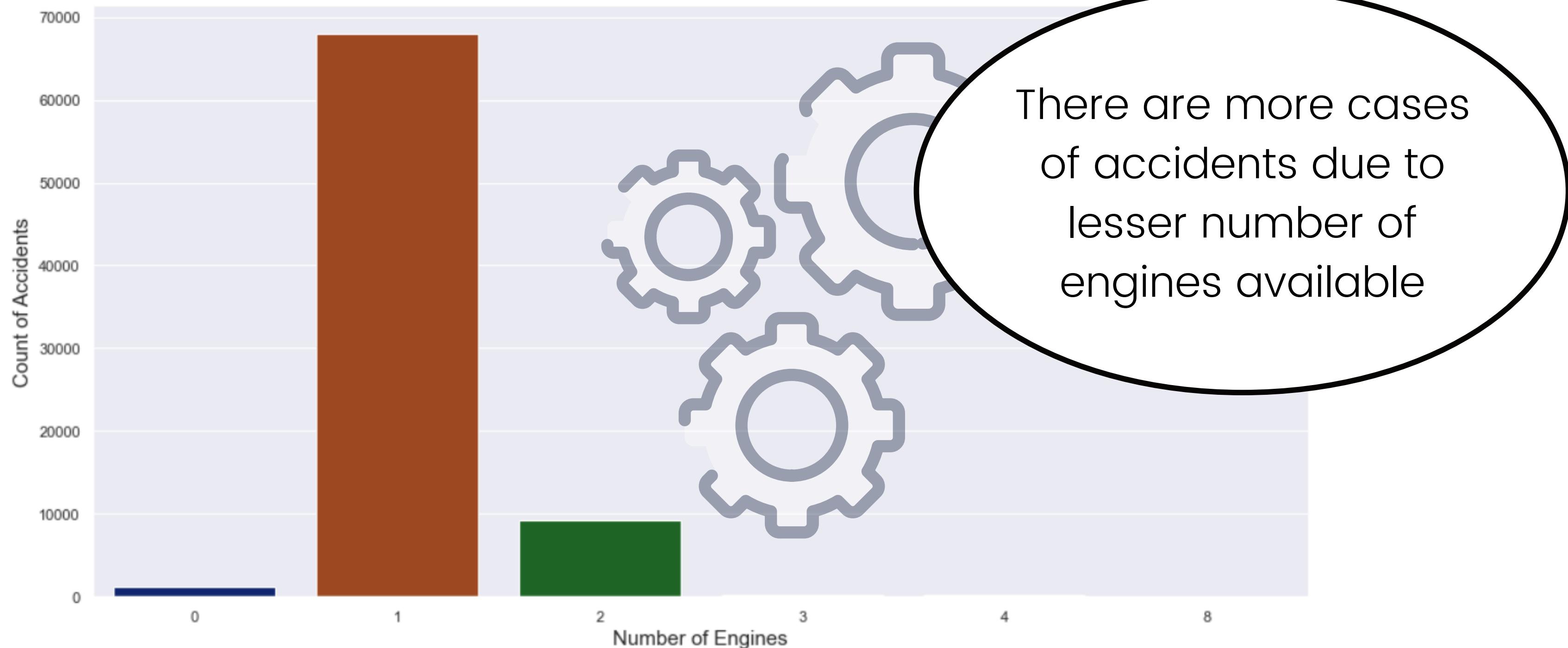


**Cessna** company has both the highest fatal and non-fatal rate accidents

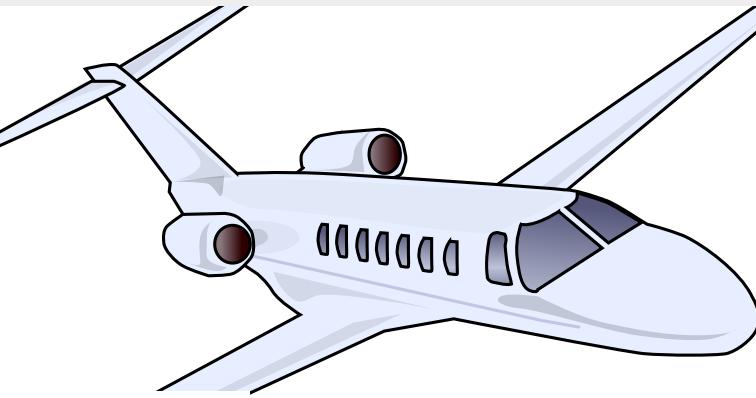


# Explorative Data Analysis

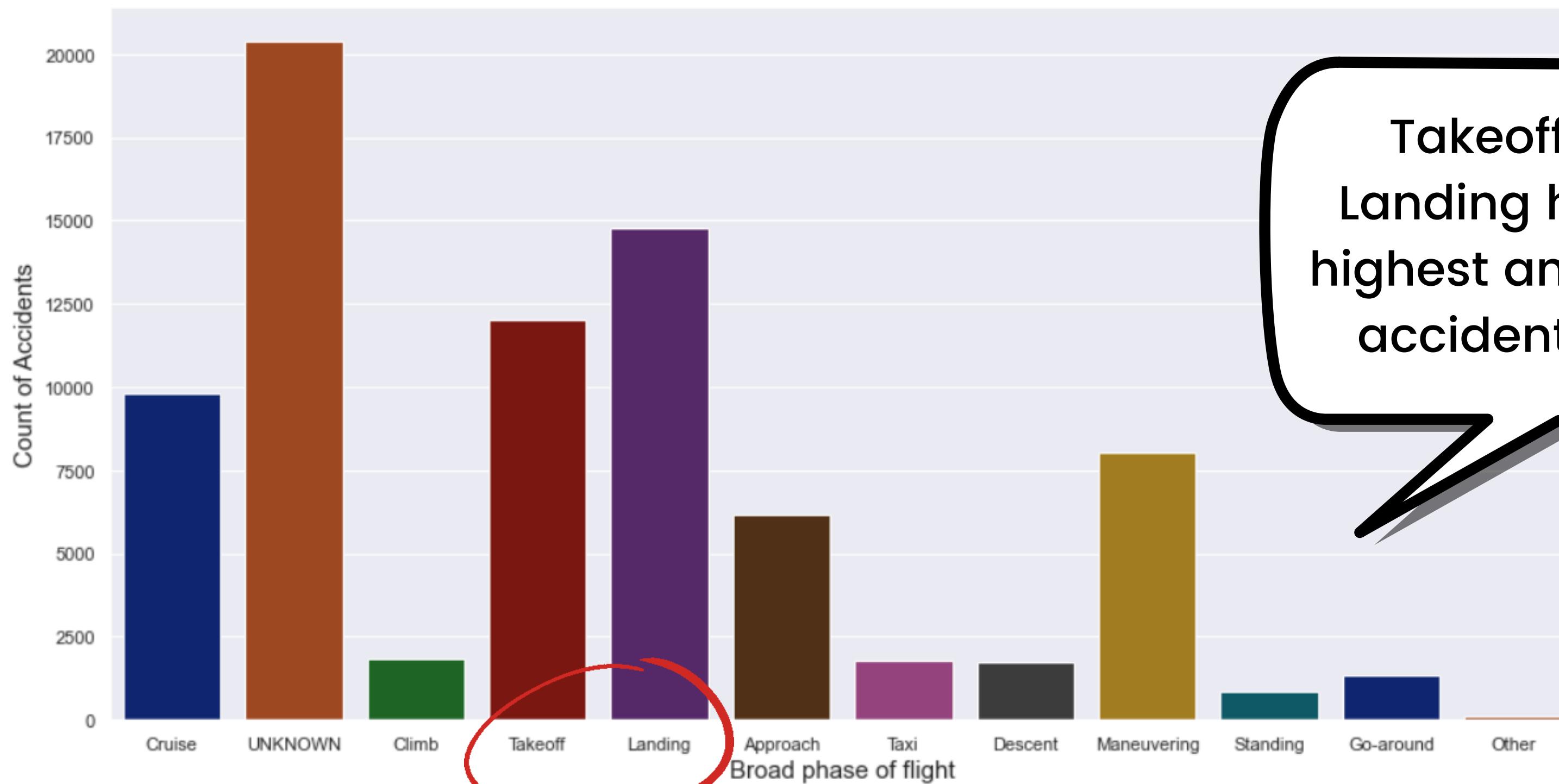
## Number of Engines in relation to the count of accident



# Explorative Data Analysis



## Phase of flight accident scenarios



Takeoff and  
Landing has the  
highest amount of  
accident rates

# CLASSIFICATION TREE



Comparing factors such as the Number of engine, Weather condition, Company & Month with injury severity.

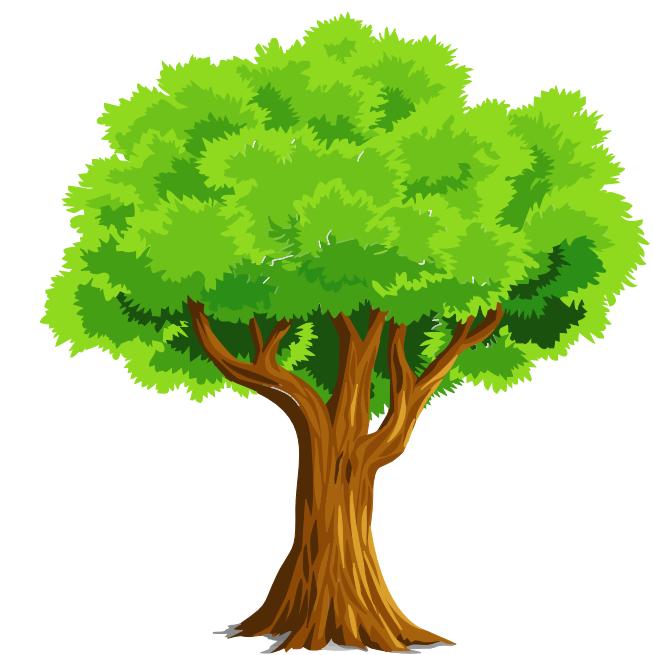
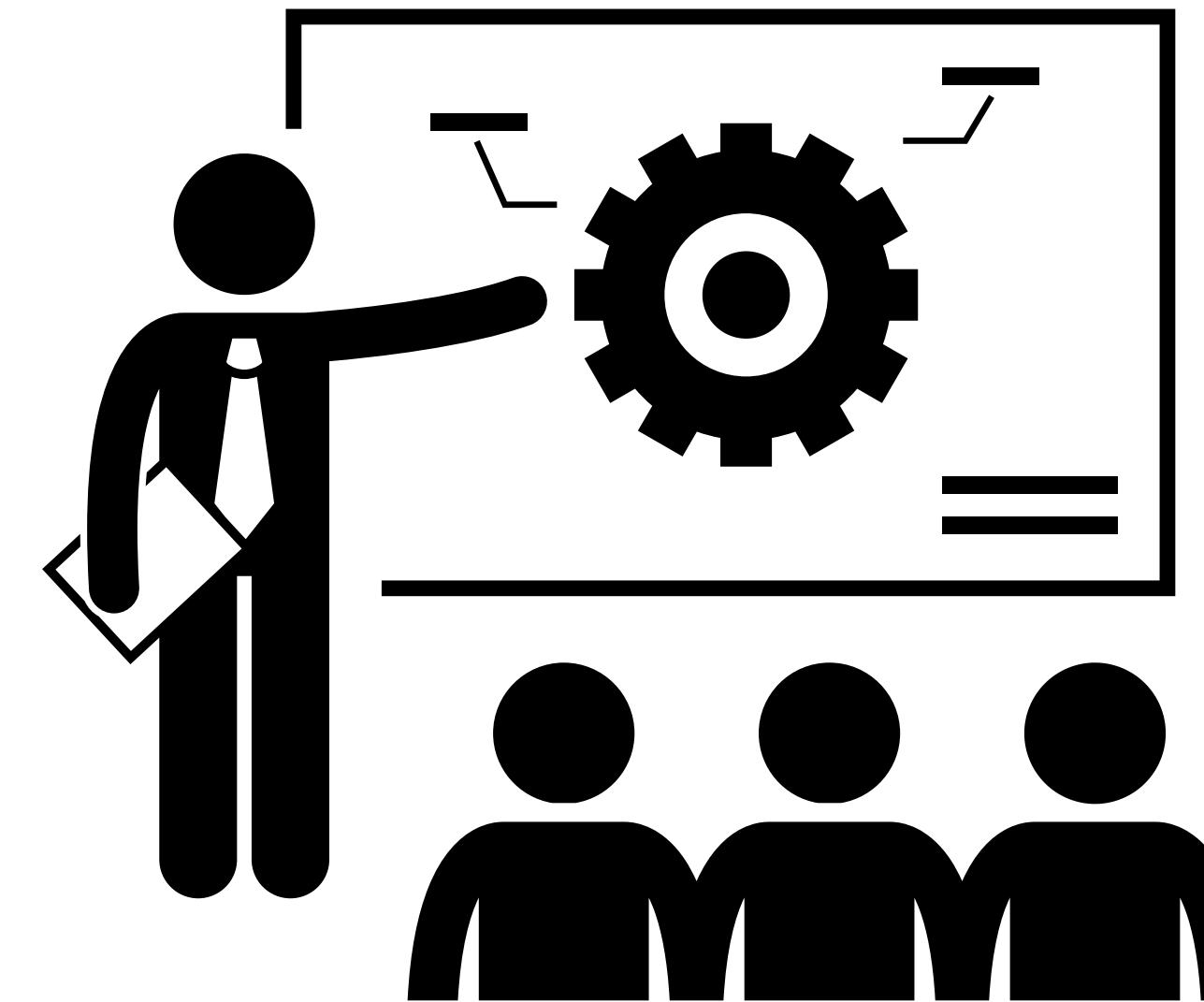
# CLASSIFICATION TREE

# Predictor

y = Injury Severity

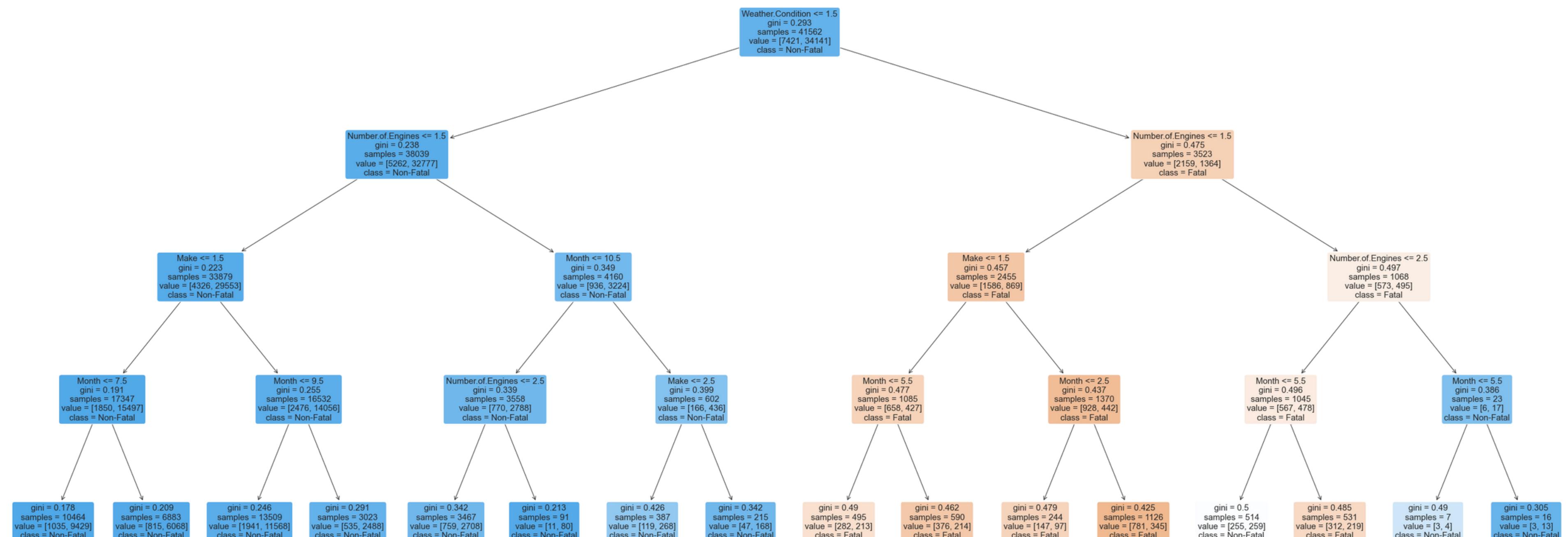
#Response

X = Number.of.Engines,Weather.Condition,Make,Month

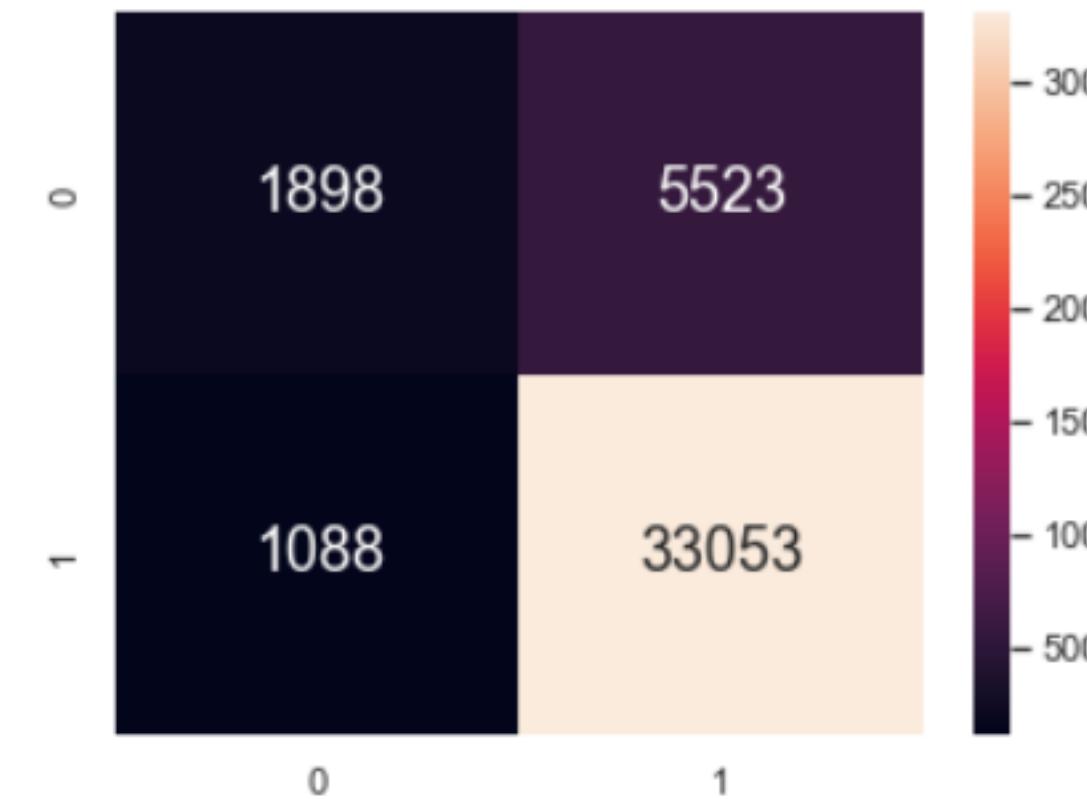
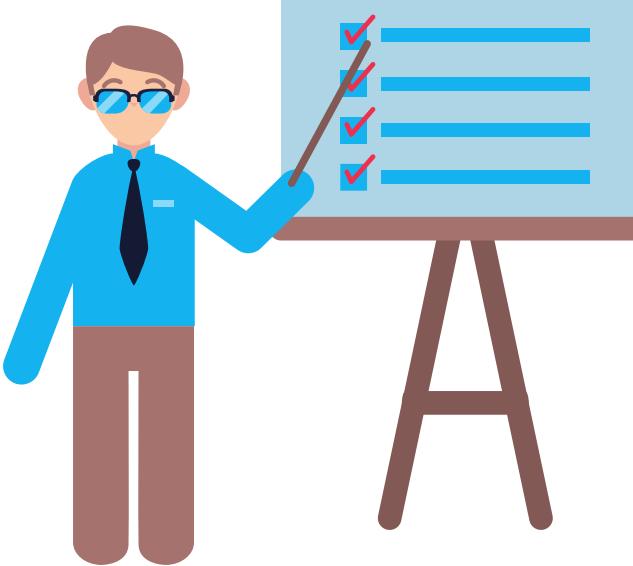


# MACHINE LEARNING(CLASSIFICATION TREE)

- Number of engine, weather condition, make & month with injury severity



# MACHINE LEARNING(CLASSIFICATION TREE)



## TRAIN SAMPLE

- TRUE POSITIVE=33053
- TRUE NEGATIVE=1898
- FALSE POSITIVE=5523
- FALSE NEGATIVE=1088



## TEST SAMPLE

- TRUE POSITIVE=8230
- TRUE NEGATIVE=476
- FALSE POSITIVE=1408
- FALSE NEGATIVE=277

# MACHINE LEARNING(CLASSIFICATION TREE)



## TRAIN DATASET

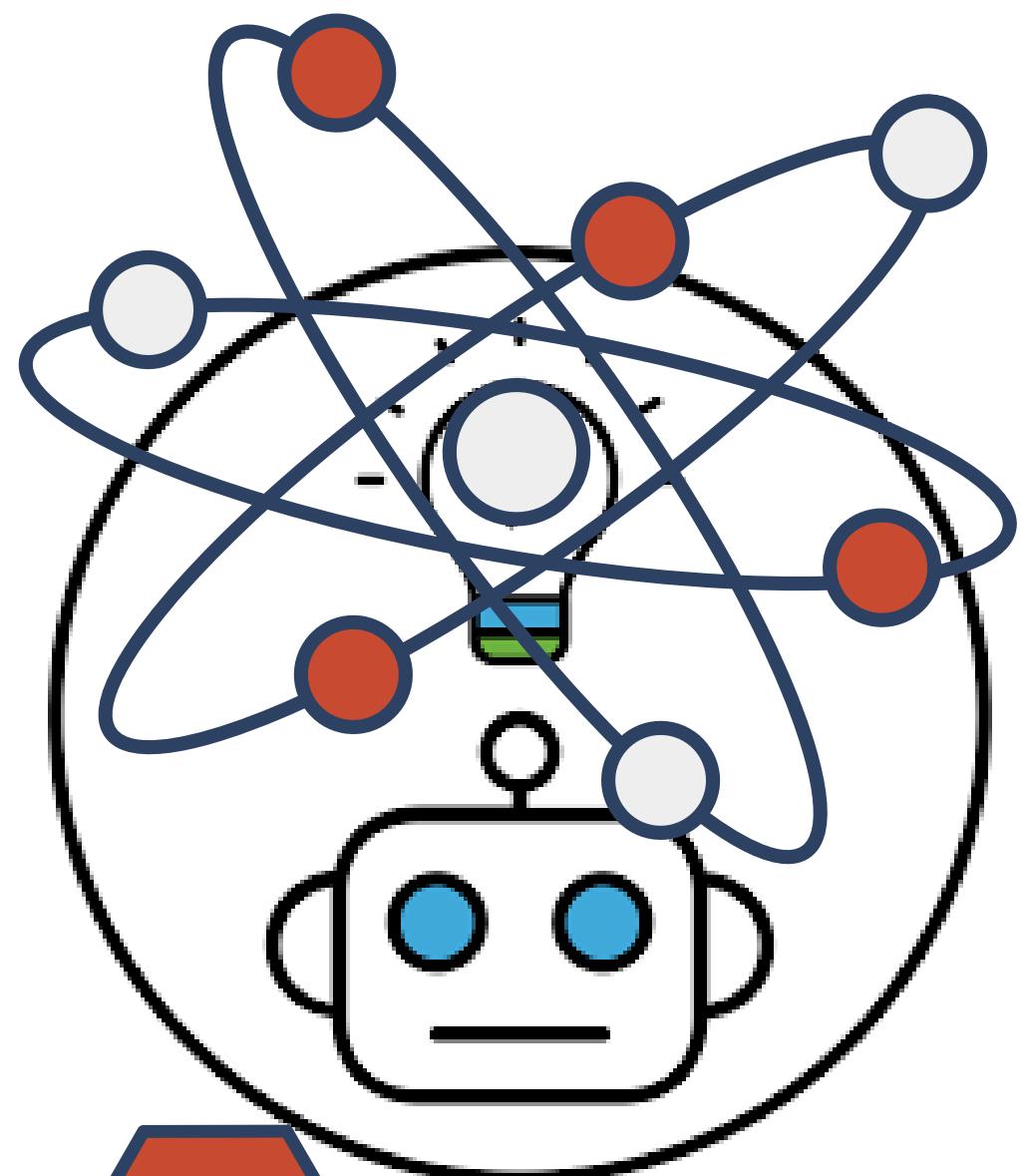
- CLASSIFICATION ACCUARACY=0.825
- TRUE POSITIVE RATE =0.968
- TRUE NEGATIVE RATE=0.255



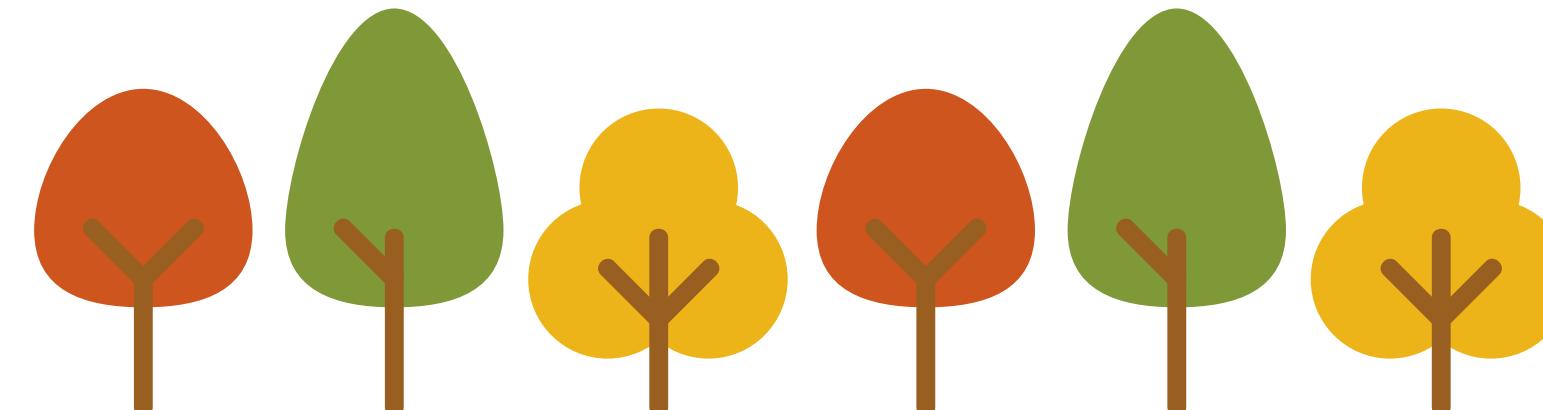
## TEST DATASET

- CLASSIFICATION ACCUARACY=0.824
- TRUE POSITIVE RATE = 0.967
- TRUE NEGATIVE RATE= 0.252

# MACHINE LEARNING

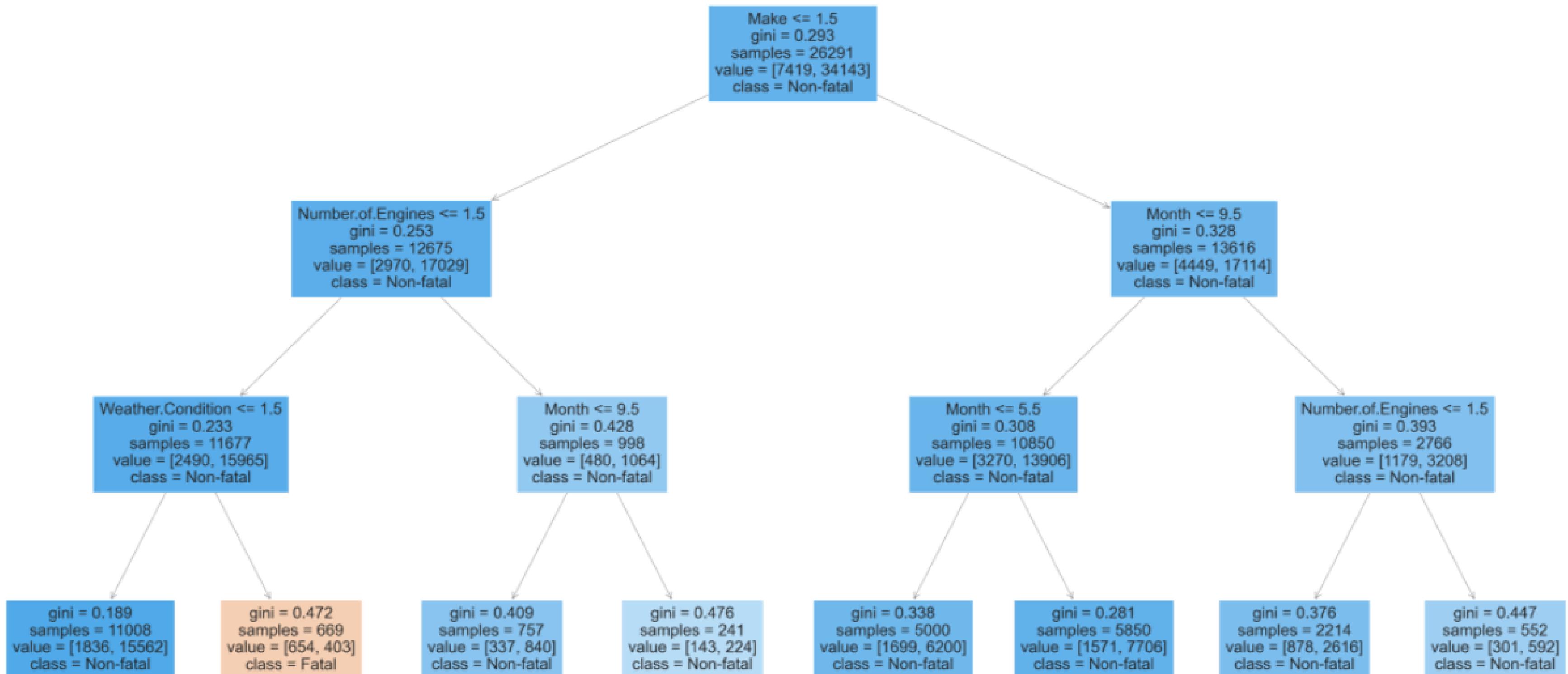


TO FURTHER OPTIMISE OUR  
CLASSIFICATION TREE  
  
WE DECIDED TO USE  
RANDOM FOREST CLASSIFIER



# MACHINE LEARNING(RANDOM FOREST CLASSIFIER)

- Number of engine, weather condition, make & month with injury severity

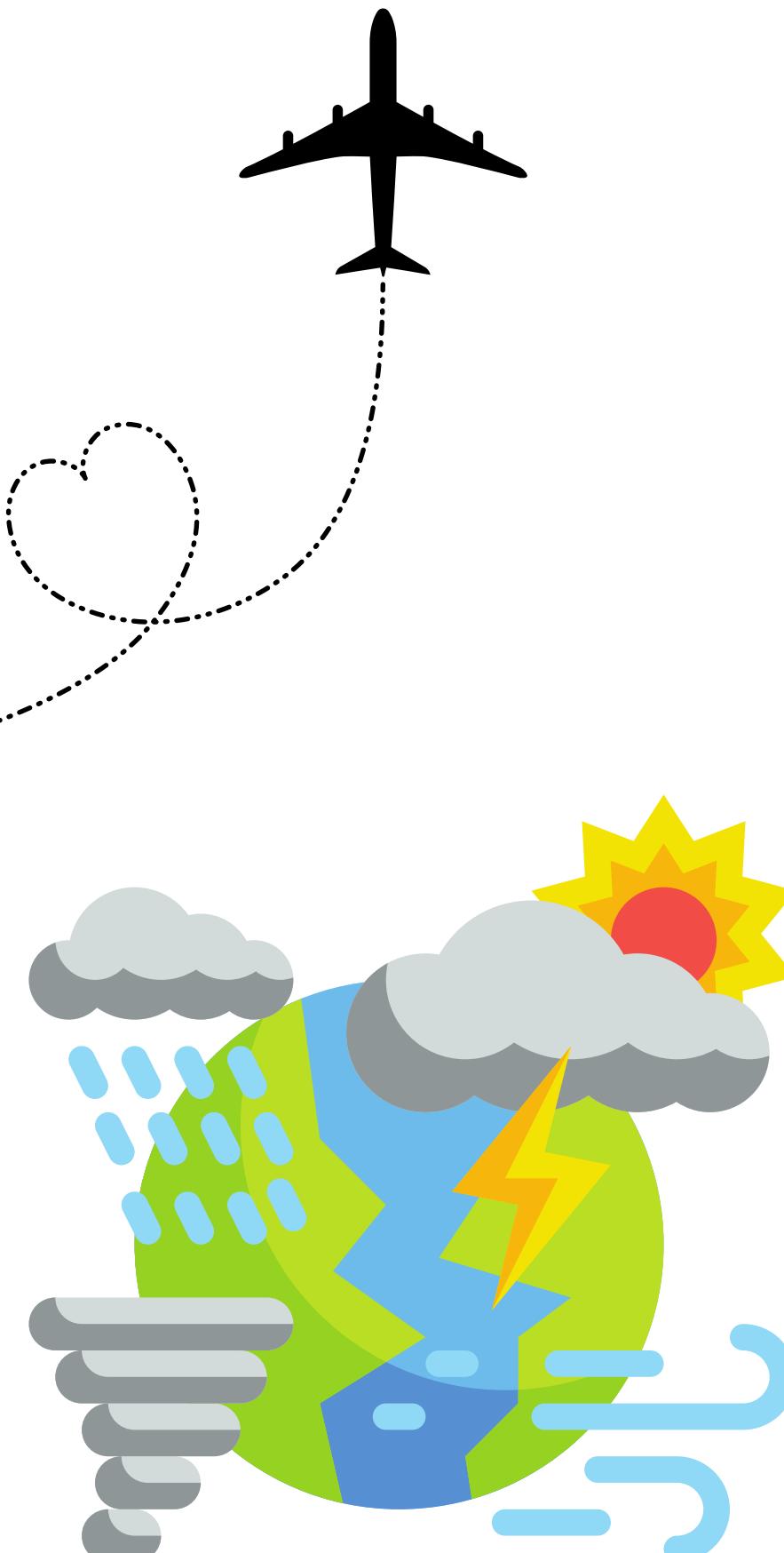


# MACHINE LEARNING(CLASSIFICATION TREE)



	Varname	Imp
1	Weather.Condition	0.849349
0	Number.ofEngines	0.091880
2	Make	0.034458
3	Month	0.024313

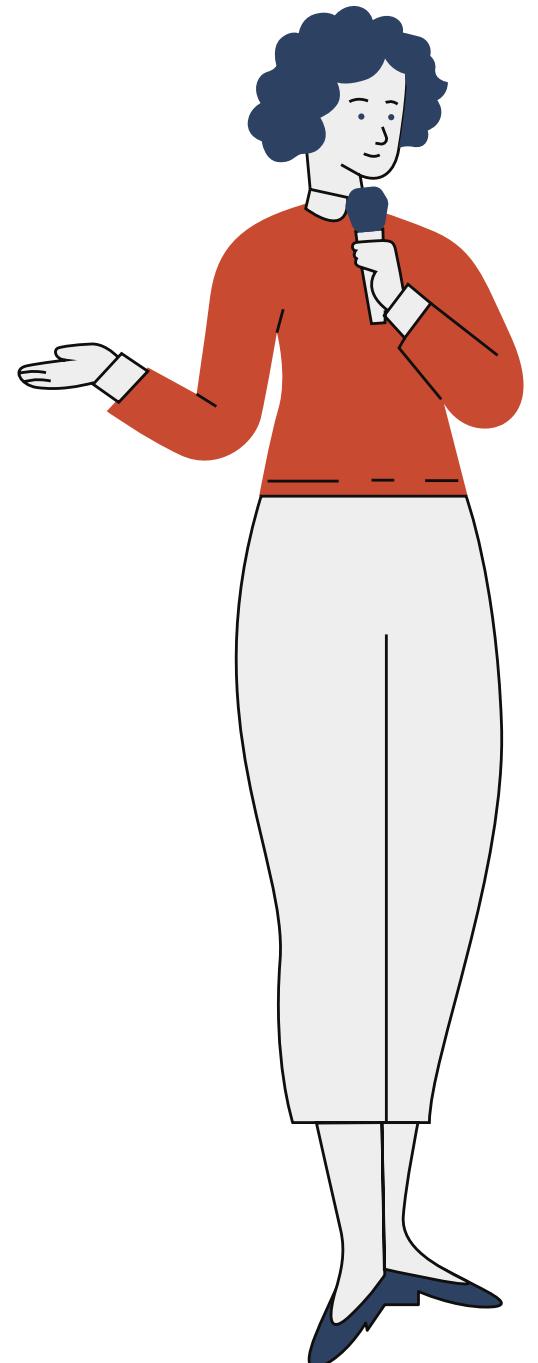
- Random Forest Classifier resulted in **better confidence** in the decision tree as seen from the gini coefficient
- It helped us identify the variable with the **most impact** on aviation incident



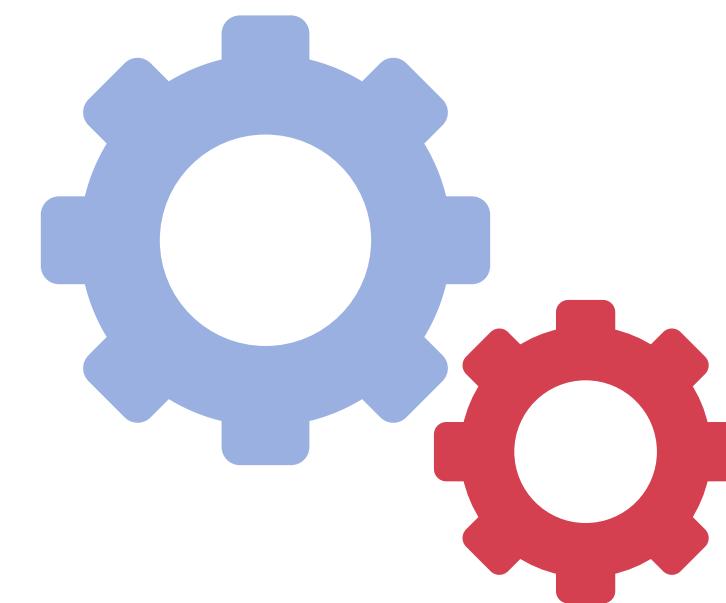
## WEATHER CONDITION

- WEATHER CONDITION HAS THE GREATEST IMPACT ON AVIATION INCIDENT
- IT IS HARD TO PREDICT AND REQUIRE CONSISTENT MONITORING
- USE MACHINE LEARNING MODELS TO BETTER PREDICT WEATHER CONDITIONS
- WITH BETTER WEATHER PREDICTION AIRLINES CAN MAKE MORE INFORMED DECISIONS AND PILOTS CAN BE MORE WARY OF MORE SUITABLE FLIGHT PATHS

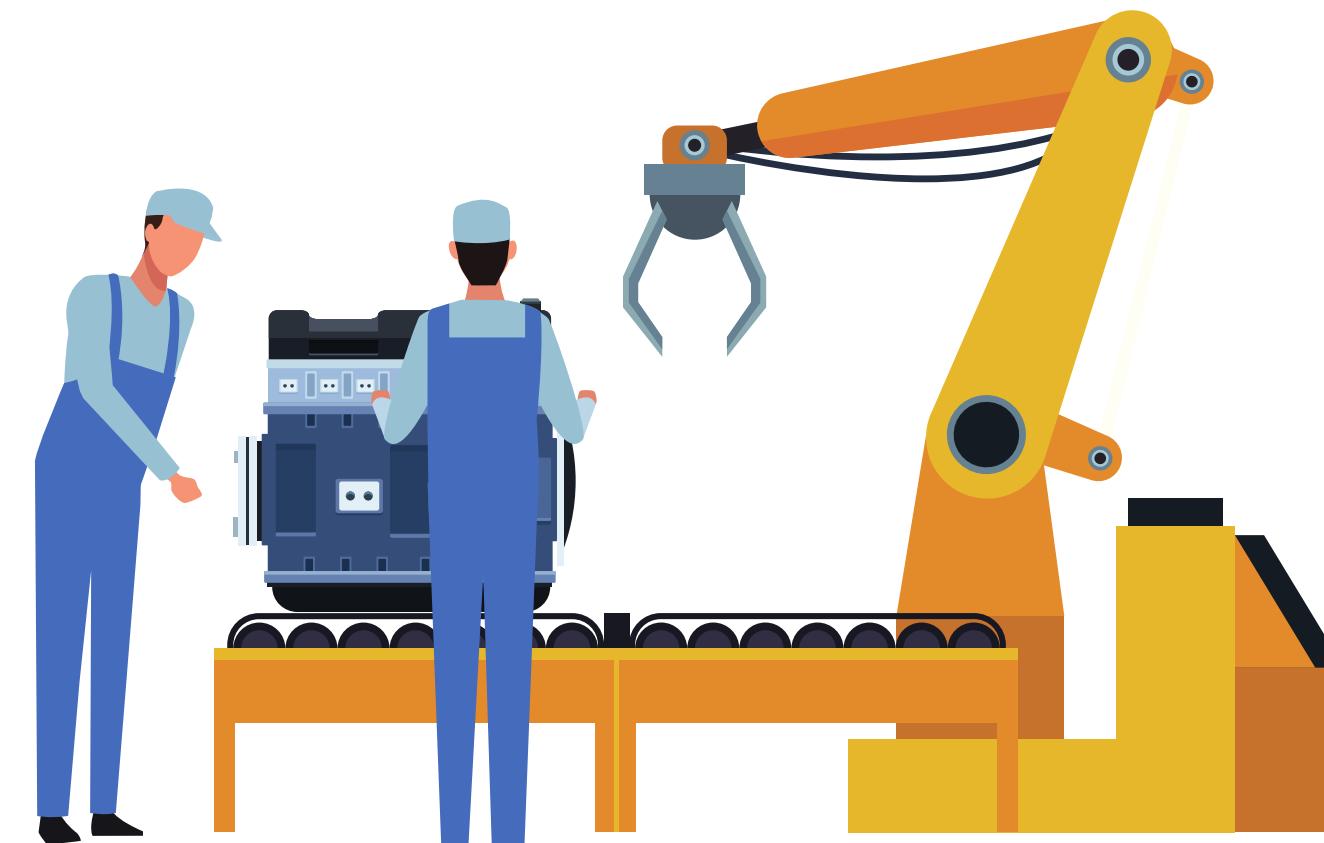
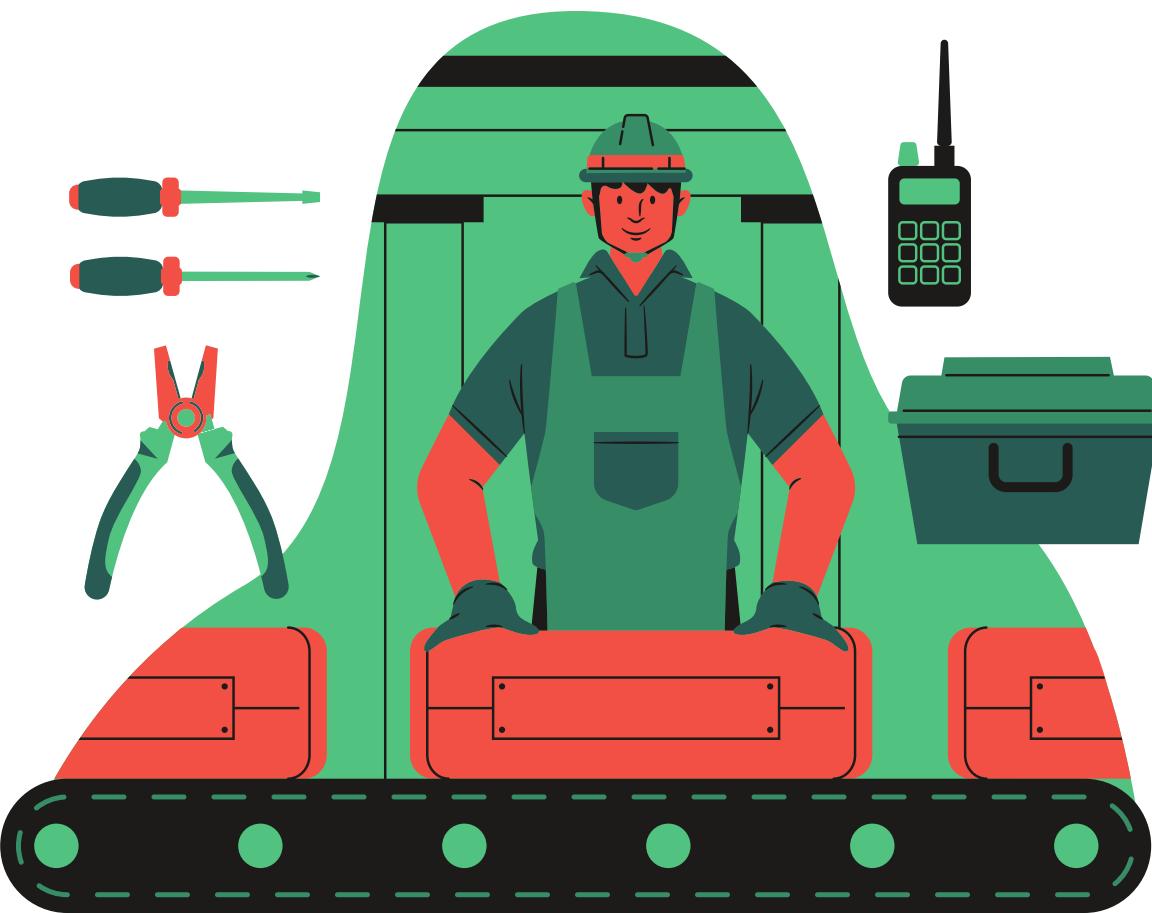
## COUNT OF ENGINE IS A MAJOR FACTOR CONTRIBUTING TO AVIATION INCIDENTS



- Manufacturers of planes should consider improving the design of their planes i.e to include more engines
- Although it may be less economical, it results in better ethical practice and lower chance of aviation accidents



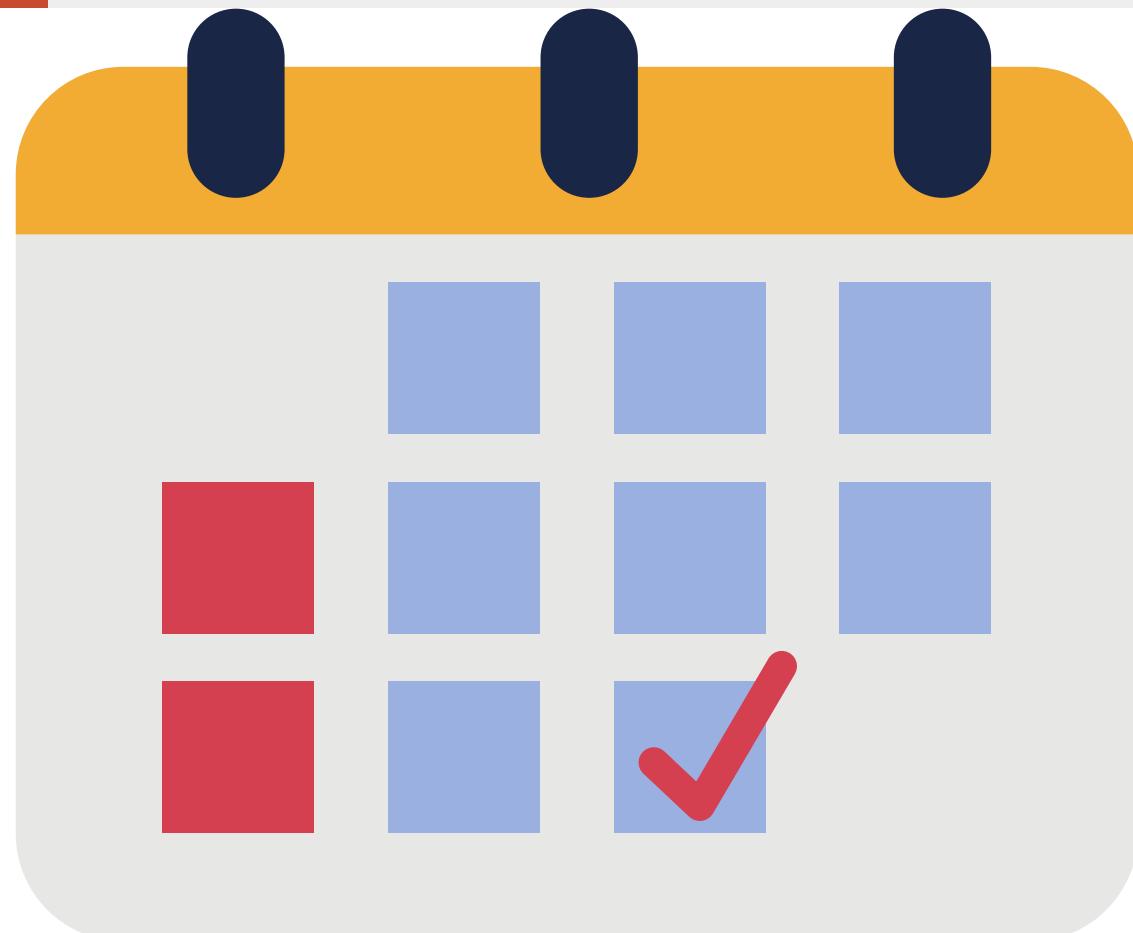
# CONCLUSION AND INSIGHTS DRIVEN



## Company

- WHILE THE MANUFACTURER DOES NOT HAVE A SIGNIFICANT IMPACT ON AVIATION INCIDENTS BASED ON OUR DATA ANALYSIS , THEY STILL HAVE AN IMPACT ON AVIATION INCIDENTS
- GOVERNMENT BODIES SHOULD MAKE AN EFFORT TO MAKE SURE MANUFACTURER UPHOLD A CERTAIN STANDARD FOR THE PARTS OF THE PLANE THEY PRODUCE

# CONCLUSION AND INSIGHTS DRIVEN



## MONTH

- WE REALISE THE VARIABLE MONTH HAD THE LEAST IMPACT ON AVIATION INCIDENTS AS COMPARED TO THE OTHER VARIABLES
- MOST LIKELY BECAUSE THE CERTAIN MONTH HAS MORE FLIGHTS BUT THERE IS LESS AVIATION INCIDENTS

**Thank you  
for listening!**

