# Airline

# SC1015 Mini-Project

Predicting Flight Delays

✈ **Members:**

Tan Leong Jun Joseph (U2321339H)
Tan Zhe Kai (U2322419A)

# TABLE OF CONTENTS

# 01 Problem Formulation

**Flight Delay:** | >15 MINS

**On-time:** | ~80%

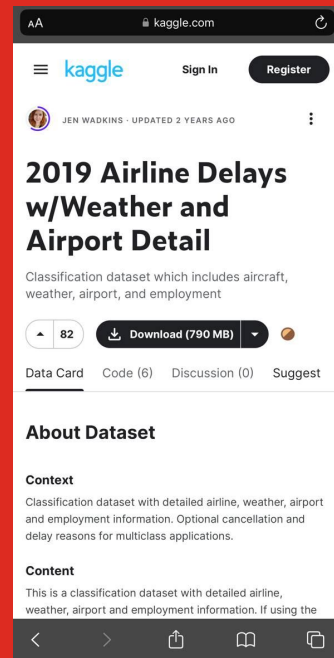**Average length:** | ~50 MINS

# Airline ✈

# 🗂 Dataset

- 8 data files

- Overlapping variables

- All data files connected

- Extensive data preparation and cleansing techniques needed

# Data Preparation

**Handling NaN values**

**Feature Engineering**

**Recheck for NaN values**

**Encode categorical variables**

**Drop irrelevant columns & duplicates**
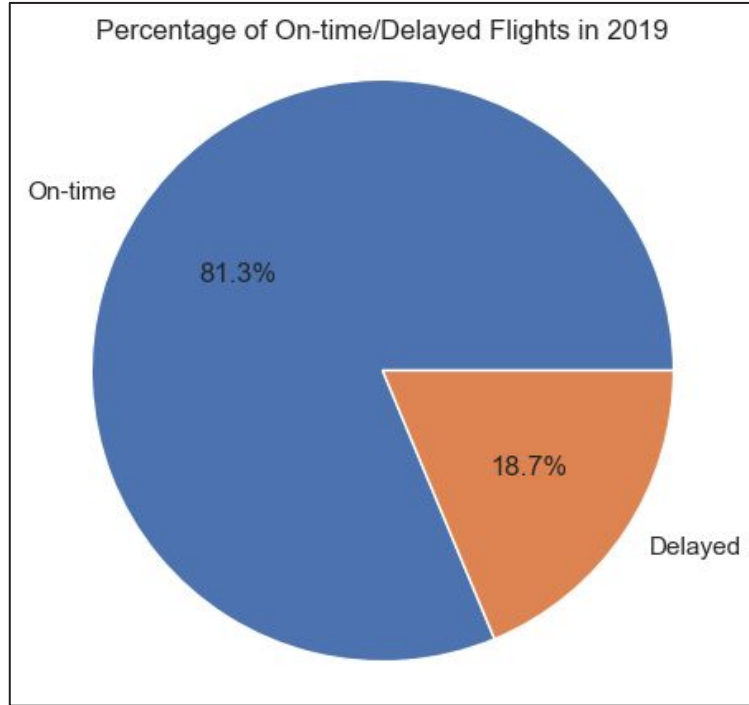
**Changing Data Types**

**Merging into one**

# 02
# Exploratory Data Analysis

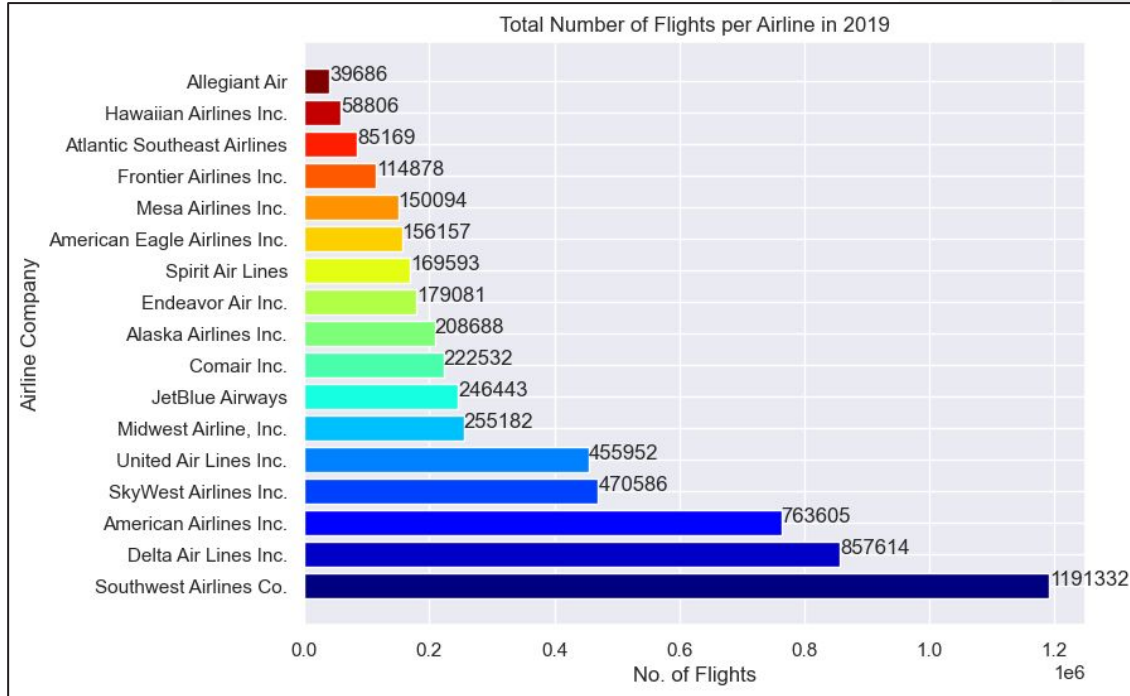Percentage of On-time/Delayed Flights in 2019

On-time 81.3%

Delayed 18.7%

# On-time vs Delayed

| ON-TIME:  4.5 million flights

| DELAYED: 1 million flights

Increasing trend in the number of flights per month

# Which month has the most flight delays?



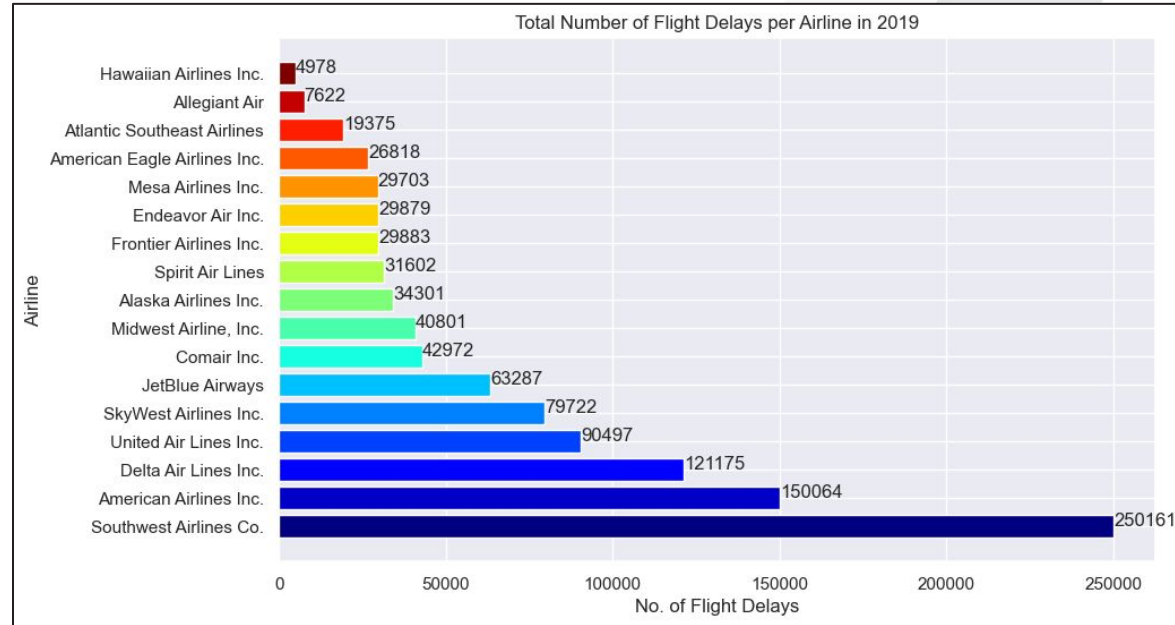Number of Flight Delays per Month in 2019

1. June, July & August (Summer)

2. December (Winter)

Possible reasons:
- Adverse weather conditions
- Vacation period
- High number of concurrent flights

# Which airline has the most flight delays?

Total Number of Flight Delays per Airline in 2019

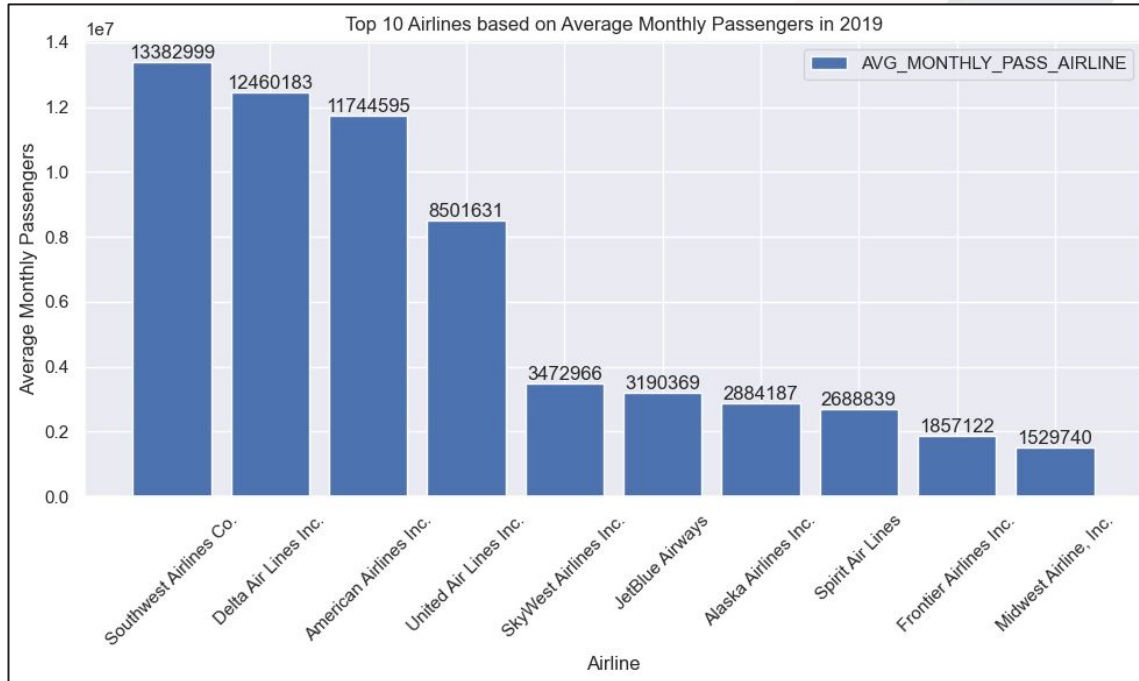| Airline | No. of Flight Delays |
|---|---|
| Hawaiian Airlines Inc. | 4978 |
| Allegiant Air | 7622 |
| Atlantic Southeast Airlines | 19375 |
| American Eagle Airlines Inc. | 26818 |
| Mesa Airlines Inc. | 29703 |
| Endeavor Air Inc. | 29879 |
| Frontier Airlines Inc. | 29883 |
| Spirit Air Lines | 31602 |
| Alaska Airlines Inc. | 34301 |
| Midwest Airline, Inc. | 40801 |
| Comair Inc. | 42972 |
| JetBlue Airways | 63287 |
| SkyWest Airlines Inc. | 79722 |
| United Air Lines Inc. | 90497 |
| Delta Air Lines Inc. | 121175 |
| American Airlines Inc. | 150064 |
| Southwest Airlines Co. | 250161 |

Most number of flight delays:

Southwest Airlines Co.

Possible reasons:
- Budget airline
- Popular

# Most Popular Airlines



Top 10 Airlines based on Average Monthly Passengers in 2019

1. Southwest Airlines Co.
2. Delta Air Lines Inc.
3. American Airlines Inc.

- Major American airlines

# 03 Models

**1st**

| Decision Tree

**2nd**

| Random Forest

**3rd**

| AdaBoost

# 4th
| XGBoost

# 5th
| Naive Bayes

# Comparison Summary

# Metrics used for models

🎯 | **Accuracy**  % of correct classifications

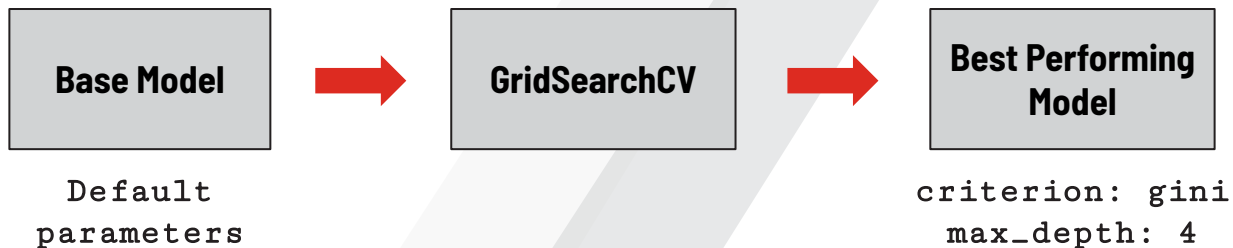🧭 | **Precision**  % of true positive/(true pos + false pos)

🔲 | **Recall**  % of true positive/(true pos + false neg)

🔗 | **F1 Score**  Mean of precision and recall

# Decision Tree Classifier



| Base Model | → | GridSearchCV | → | Best Performing Model |
|---|---|---|---|---|

```
Default
parameters
```

```
criterion: gini
max_depth: 4
```

Hyperparameter Tuning using
GridSearchCV:
● To obtain the best parameters
  to generate the model with
  the highest accuracy

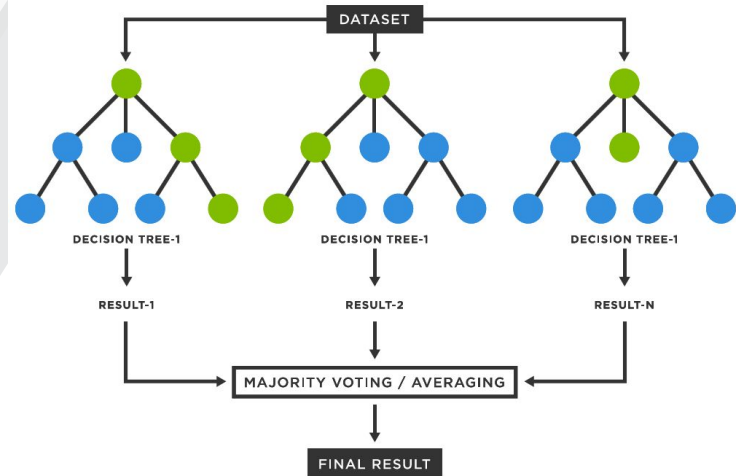| | |
|---|---|
| Accuracy | 81.3% |
| Precision | 67.0% |
| Recall | 0% |
| F1-Score | 0% |

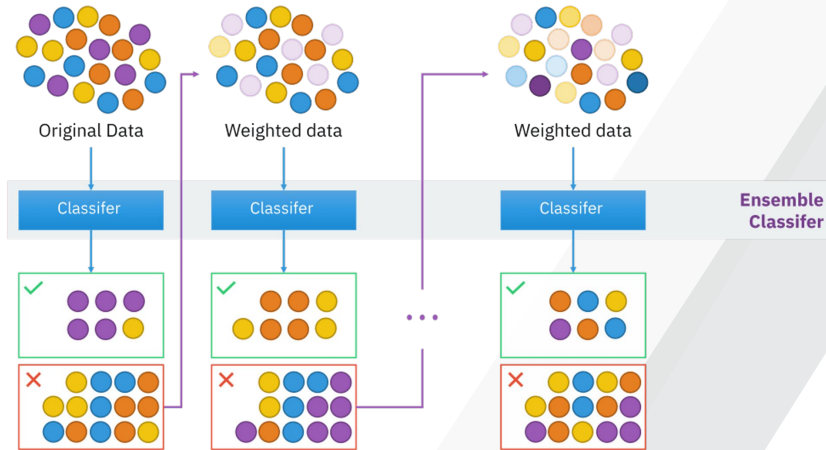# Random Forest Classifier

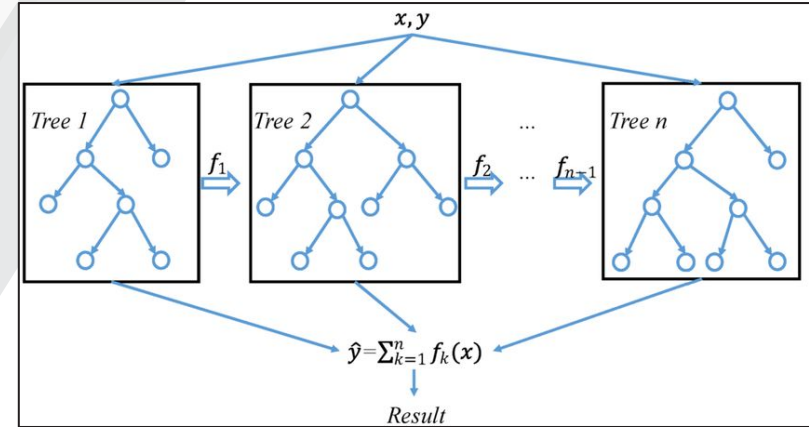Base Model → K-folds Cross Validation



- Combination of decision trees
- Prevent overfitting
- Result = avg no. of models, evaluated k times

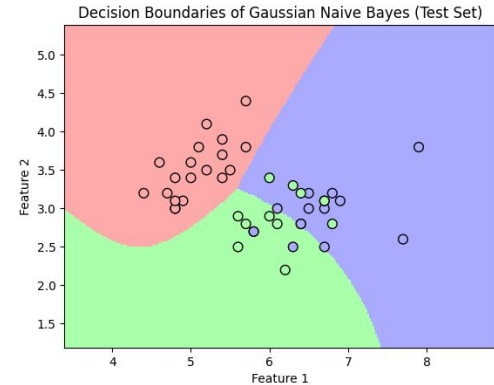# AdaBoost Classifier
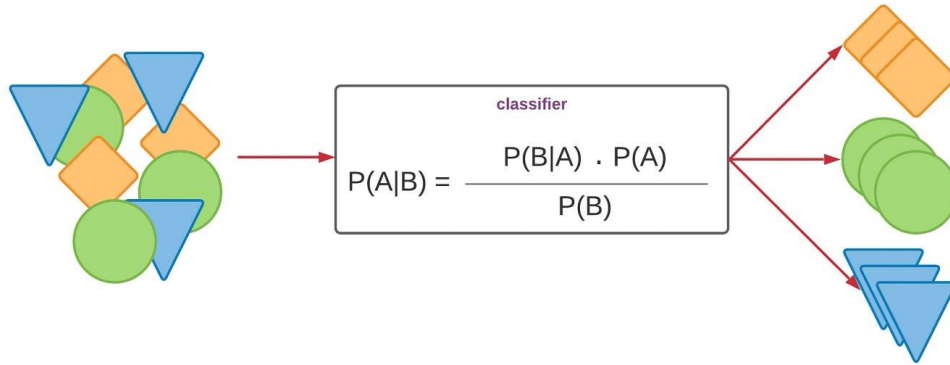
# XGBoost Classifier

# Naive Bayes Classifier



$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

classifier

Decision Boundaries of Gaussian Naive Bayes (Test Set)

- Supervised learning

- Probabilistic approach + Gaussian distribution

# Summary of Classification Models

| Model | Train Accuracy | Test Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Decision Tree | 81.3% | 81.3% | 67.0% | 0.0 | 0.0 |
| Random Forest | 99.3% | 82.0% | 55.5% | 20.1% | 29.5% |
| AdaBoost | 81.3% | 81.3% | 54.4% | 0.0% | 1.0% |
| XGBoost | 82.4% | 82.3% | 66.1% | 10.7% | 18.4% |
| Naive Bayes | 81.3% | 81.3% | 0.0% | 0.0% | 0.0% |

# Feature Importances on the Best Model

| Model | Train Accuracy | Test Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Random Forest | 99.3% | 82.0% | 55.5% | 20.1% | 29.5% |

**Top 5 Most Important Features**

| Feature | Importance |
|---|---|
| PLANE_AGE | 0.096223 |
| PREVIOUS_AIRPORT | 0.092637 |
| CONCURRENT_FLIGHTS | 0.082243 |
| DEP_TIME_BLK | 0.079219 |
| AWND | 0.077038 |

1. Age of Plane
2. Previous Airport
3. Concurrent Flights
4. Departure Time Block
5. Max Wind Speed

# Data Driven Insights & Recommendations



## | Travellers

Choose less popular seasons

## | Airlines

Prioritize newer aircraft models

## | Airport

Interact with each other

# Thank You !

✈ **References:**

- Slotnick, D. (2020, February 19). Why Your Airplane (Still) Might Be Delayed. The New York Times. Retrieved from https://www.nytimes.com/2020/02/19/business/air-travel-delays-airlines.html
- A general architecture of XGBoost. (n.d.). In ResearchGate. Retrieved from https://www.researchgate.net/figure/A-general-architecture-of-XGBoost_fig3_335483097
- Naive Bayes. (n.d.). In scikit-learn: Machine Learning in Python. Retrieved from https://scikit-learn.org/stable/modules/naive_bayes.html
- Dancerworld60. (n.d.). Demystifying Naïve Bayes—Simple Yet Powerful for Text Classification. Medium. Retrieved from https://medium.com/@dancerworld60/demystifying-na%C3%AFve-bayes-simple-yet-powerful-for-text-classification-ad92b14a5c7

Airline