

# **PROJECT REPORT**

## **Developing a Flight Delay Prediction Model using Machine Learning**

**PNT2022TMID35530**

<b>AKASHRAM J</b>	<b>2019115011</b>
<b>GAYATHRI P</b>	<b>2019115033</b>
<b>THARANYAA R</b>	<b>2019115113</b>
<b>VIJAY S</b>	<b>2019115120</b>

# **CHAPTER 1**

## **INTRODUCTION**

Delay is one of the best-known performance indicators in any transportation system. Civil aviation officials in particular understand delay as the time at which a flight is delayed or rescheduled. Delay can therefore be expressed as the difference between the scheduled flight time and the actual departure or arrival time. National regulators have a number of indicators relating to acceptable levels of flight delays. Flight delays are a major problem associated with air transportation systems.

Analysts and data scientists are immersed in this vast amount of data generated by sensors and IoT, enhancing their computational and data management skills to extract useful information from each data. In this context, the process of understanding domains, managing data, and applying models is called data science, a trend for solving challenging big data-related problems. In this project, extensive data analysis was performed to extract the key attributes/factors responsible for flight delays. In addition, there are other factors that can affect flight delays, such as: These factors, such as climate, natural disasters, pandemics, or technical problems with aircraft, vary from place to place and are not considered in this project as such problems rarely occur.

### **1.1 PROJECT OVERVIEW**

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. Using this algorithm model is built to predict the flight delay.

### **1.2 PURPOSE**

Flight delays result in significant financial and other losses to airlines, airports and passengers. Predictions are important in the decision-making process of all parties in the aviation industry. Therefore, predicting potential delays based on flight characteristics bridges an important information asymmetry between airlines and passengers. The main use cases for this algorithm are to forecast of possible delays on a given day for a given airport and airline.

## **CHAPTER 2**

### **LITERATURE SURVEY**

#### **2.1 EXISTING PROBLEM**

Many existing flight delay prediction methods are based on small samples and/or are complex to interpret with little or no opportunity for machine learning deployment. The proposed model gains insight into factors causing flight delays, cancellations and the relationship between departure and arrival delay using exploratory data analysis.

#### **2.2 REFERENCES**

[1] Jiang, Yushan, et al. "Applying machine learning to aviation big data for flight delay prediction." 2020 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech). IEEE, 2020.

[2] Liu, Fan, et al. "Generalized flight delay prediction method using gradient boosting decision tree." 2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring). IEEE, 2020.

[3] Ding, Yi. "Predicting flight delay based on multiple linear regression." IOP Conference Series: Earth and Environmental Science. Vol. 81. No. 1. IOP Publishing, 2017.

[4] Thiagarajan, Balasubramanian, et al. "A machine learning approach for prediction of on-time performance of flights." 2017 IEEE/AIAA 36th Digital Avionics Systems Conference (DASC). IEEE, 2017.

[5] Kim, Young Jin, et al. "A deep learning approach to flight delay prediction." 2016 IEEE/AIAA 35th Digital Avionics Systems Conference (DASC). IEEE, 2016.

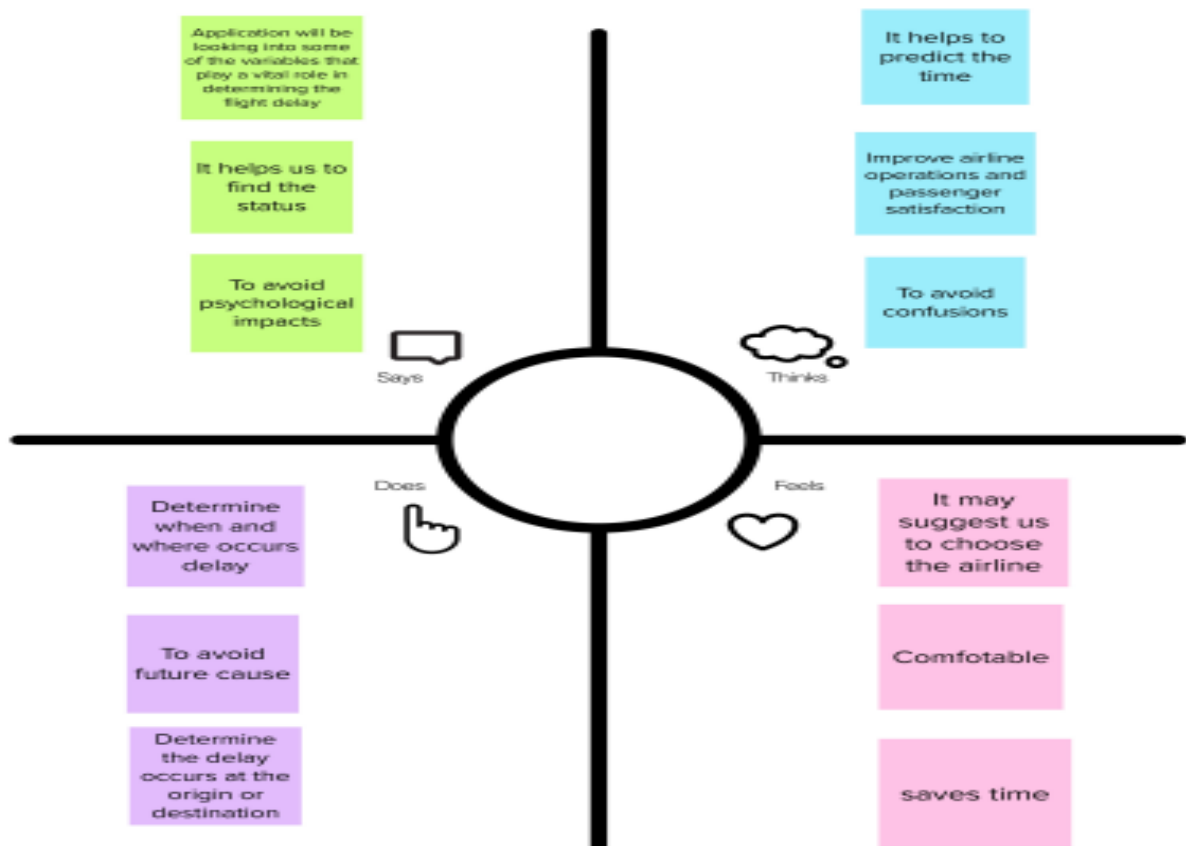
#### **2.3 PROBLEM STATEMENT DEFINITION**

Flight delays in air transportation are a major concern that has adverse effects on the economy, the passengers, and the aviation industry. This matter critically requires an accurate estimation for future flight delays that can be implemented to improve airport operations and customer satisfaction. Having said that, a massive volume of data and an extreme number of parameters have restricted the way to build an accurate model. Many existing flight delay prediction methods are based on small samples and/or are complex to interpret with little or no opportunity for machine learning deployment.

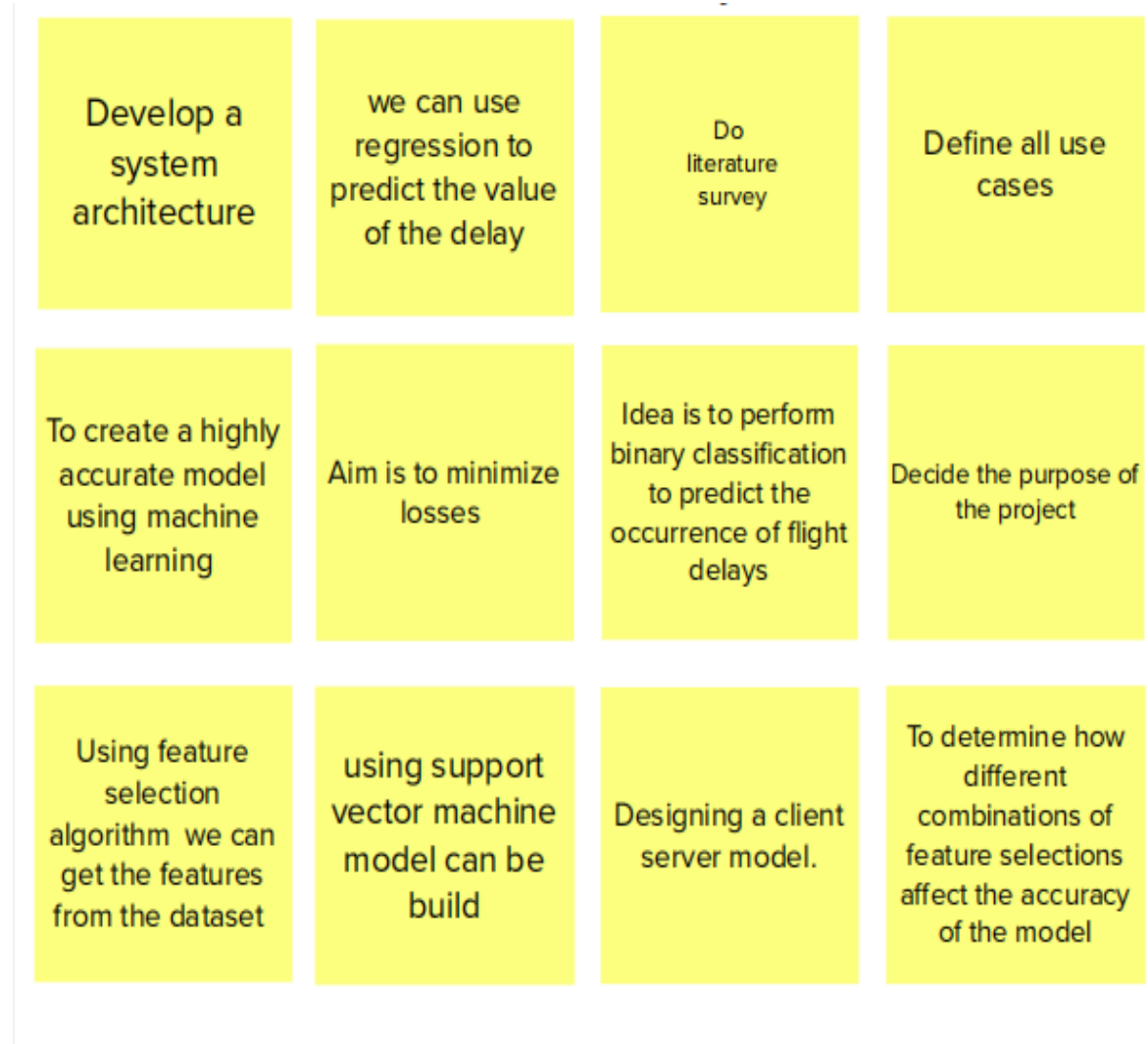
## CHAPTER 3

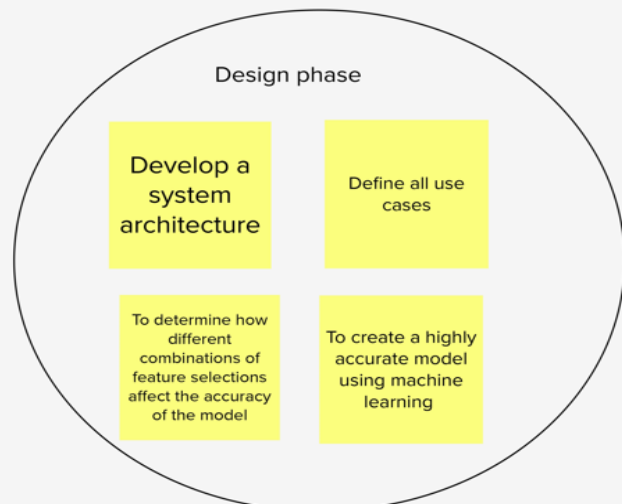
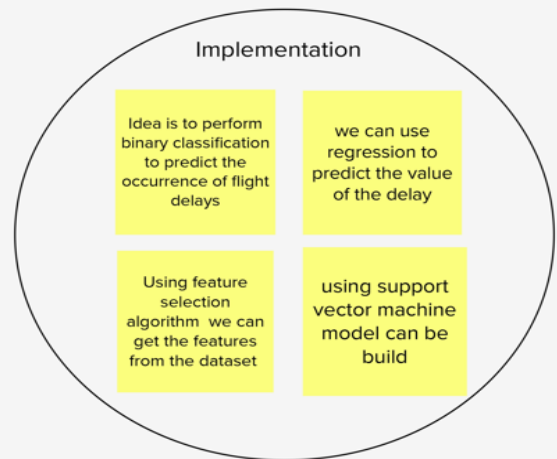
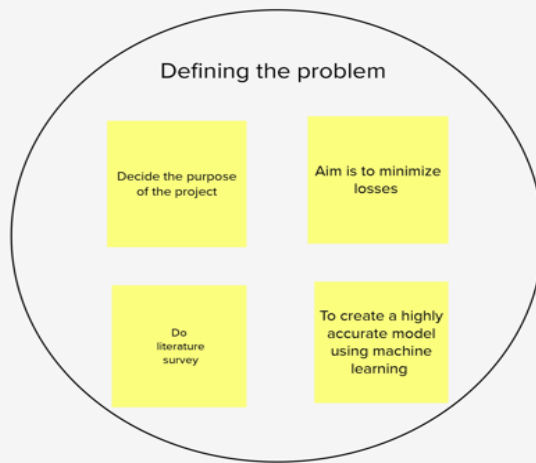
### IDEATION AND PROPOSED SOLUTION

#### 3.1 EMPATHY MAP CANVAS



### 3.2 IDEATION AND BRAINSTORMING









### 3.3 PROPOSED SOLUTION

S.No.	PARAMETER	DESCRIPTION
1	Problem Statement (Problem to be solved)	<p>Flight delays in air transportation are a major concern that has adverse effects on the economy, the passengers, and the aviation industry. This matter critically requires an accurate estimation for future flight delays that can be implemented to improve airport operations and customer satisfaction. Having said that, a massive volume of data and an extreme number of parameters have restricted the way to build an accurate model. Many existing flight delay prediction methods are based on small samples and/or are complex to interpret with little or no opportunity for machine learning deployment.</p>

2	Idea / Solution description	<p>The proposed model gains insight into factors causing flight delays, cancellations and the relationship between departure and arrival delay using exploratory data analysis. In addition, Random Forest (RF) algorithm is used to train and test the big dataset to help the model development.</p> <p>A web application has also been developed to implement the model and the testing results are presented with the limitation discussed</p>
3	Novelty / Uniqueness	<p>Many existing flight delay prediction methods are based on small samples and/or are complex to interpret with little or no opportunity for machine learning deployment. The proposed model gains insight into factors causing flight delays, cancellations and the relationship between departure and arrival delay using exploratory data analysis.</p>
4	Social Impact / Customer Satisfaction	<p>An accurate estimation of flight delay is critical for airlines because the results can be applied to increase customer satisfaction and incomes of airline agencies.</p> <p>Predicting flight delays can improve airline operations and passenger satisfaction, which will result in a positive impact on the economy</p>

5	Business Model (Revenue Model)	A web application has been developed to provide the end-users an interface to help predict flight delays. In future, we can implement the subscription plan for the prediction process and also if our model predicts well, we can sell it airlines, so they can prior inform the passenger.
---	--------------------------------	--

6	Scalability of the Solution	The proposed combined method of delay analysis and its prediction can also be further explored in other studies and also can extend the application in more comfortable with the end user. In the situation of airline, they can develop this system and make the passenger feels good and inform prior.
---	-----------------------------	--

### 3.4 PROBLEM SOLUTION FIT

Define CS, fit into CC	<b>1. CUSTOMER SEGMENT(S)</b> Passenger as well as airlines <b>CS</b>	<b>6. CUSTOMER CONSTRAINTS</b> No proper access to application. <b>CC</b>	<b>5. AVAILABLE SOLUTIONS</b> The solution is to deliver timely message and keep passengers informed about the status. <b>AS</b>	Explore AS, differentia
	<b>2. JOBS-TO-BE-DONE / PROBLEMS</b> Displaying the delaying time to the passengers is not done on time. <b>—</b>	<b>9. PROBLEM ROOT CAUSE</b> Flights being delayed due to natural occurrences like weather delay, extreme heat and severe storms and operational shortcomings, which is an expensive affair for the airlines, creating problems in scheduling and operations for the end-users thus causing bad reputation and customer dissatisfaction. <b>RC</b>	<b>7. BEHAVIOUR</b> Agree to a new connection Call the airline <b>BE</b>	

<b>3. TRIGGERS</b> Waste of time due to flight delay. <b>TR</b>	<b>10. YOUR SOLUTION</b> We implemented flight delay prediction through the proposed approaches that were based on machine learning algorithms. <b>SL</b>	<b>8. CHANNELS of BEHAVIOUR</b> <b>8.1 ONLINE</b> Check For Reimbursements <b>8.2 OFFLINE</b> Agree to A New Connection, call airline <b>CH</b>
--	--	---

<b>4. EMOTIONS: BEFORE / AFTER</b> BEFORE: Many people consider these delays to be a waste of time, and no customers enjoy waiting for long periods due to delays. AFTER: People can plan accordingly once they known the delay. <b>EM</b>		
--	--	--

## CHAPTER 4

### REQUIREMENT ANALYSIS

#### 4.1 FUNCTIONAL REQUIREMENTS

<b>FR No.</b>	<b>Functional requirements (epic)</b>	<b>Sub Requirement (Story / Sub-Task)</b>
FR-1	Details	Getting input like current year, month, date, selecting the airline and airport details from user.
FR-2	Data processing	Given data is fed to the model, using the algorithm it predicts
FR-3	Output	Displaying the result as delayed or not delayed.

#### 4.2 NON-FUNCTIONAL REQUIREMENTS

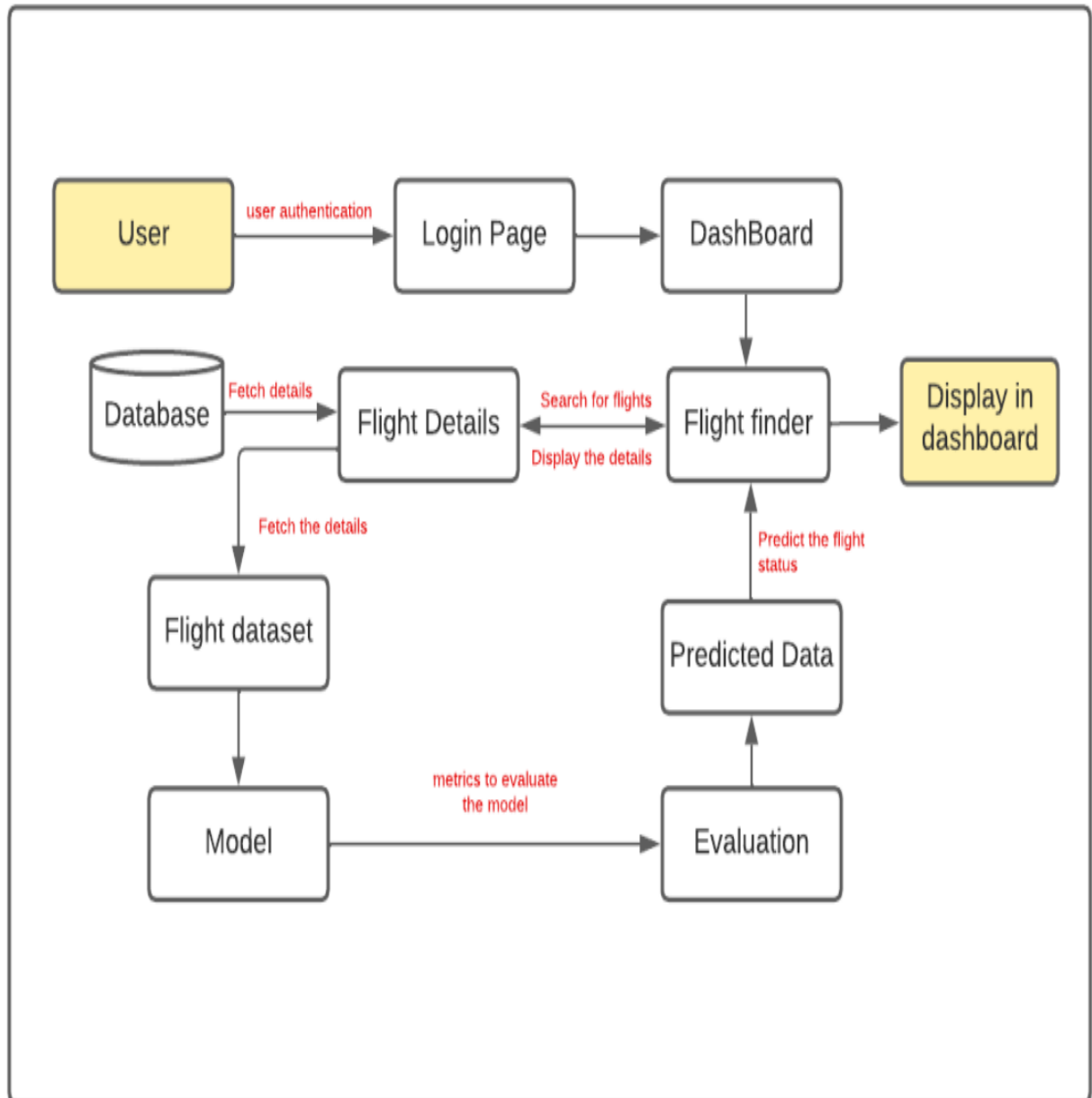
<b>FR No.</b>	<b>Non-Functional Requirement</b>	<b>Description</b>
NFR-1	<b>Usability</b>	User interface is very effective to use when compared with others.
NFR-2	<b>Security</b>	The data collected from the user will be stored securely in the cloud

NFR-3	<b>Reliability</b>	The user can trust the results from the application and they can check their flight status
NFR-4	<b>Performance</b>	Accurate prediction can be achieved.
NFR-5	<b>Availability</b>	Available if the network bandwidth of the user is of good range
NFR-6	<b>Scalability</b>	This application can be accessed from any place.

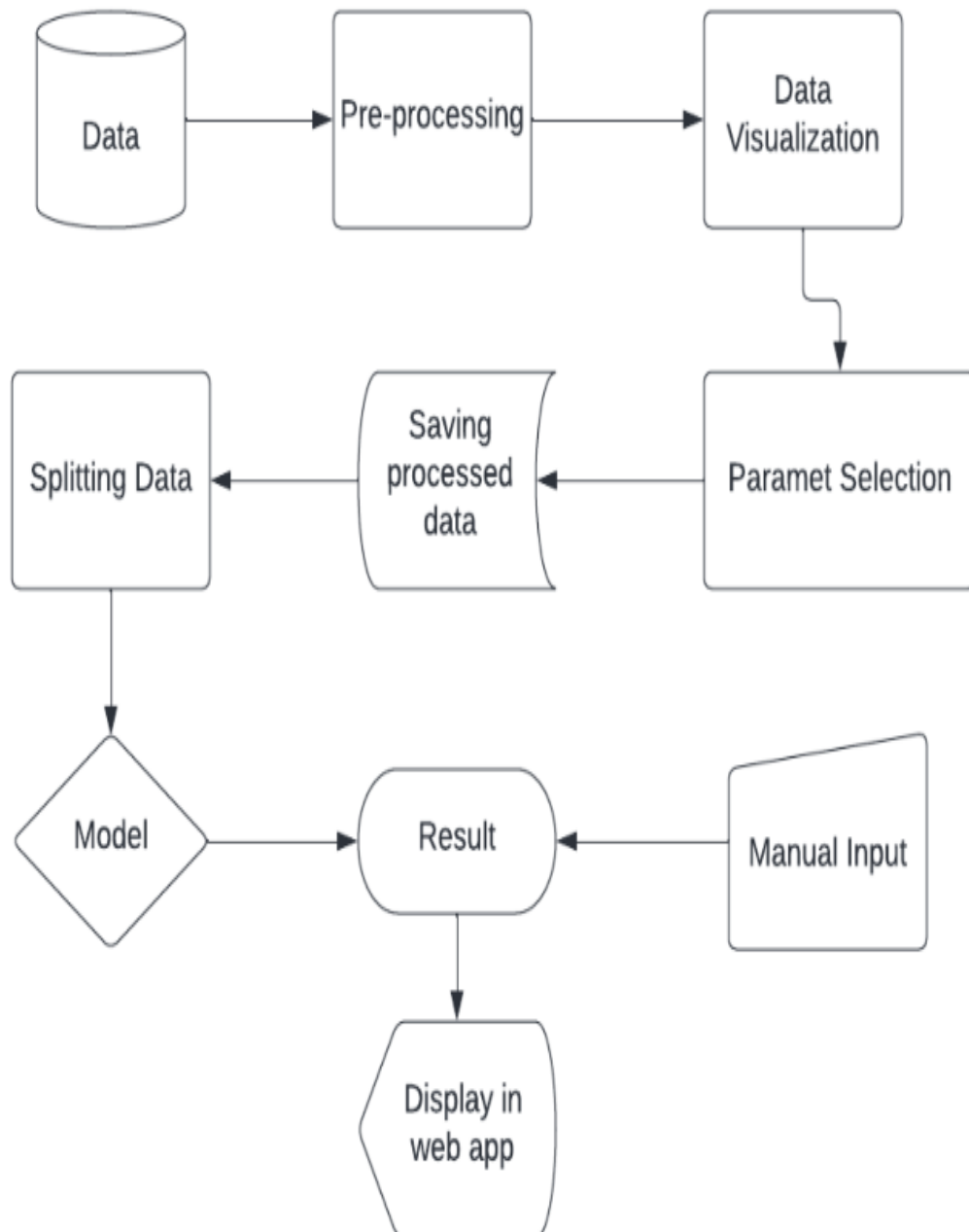
## CHAPTER 5

### PROJECT DESIGN

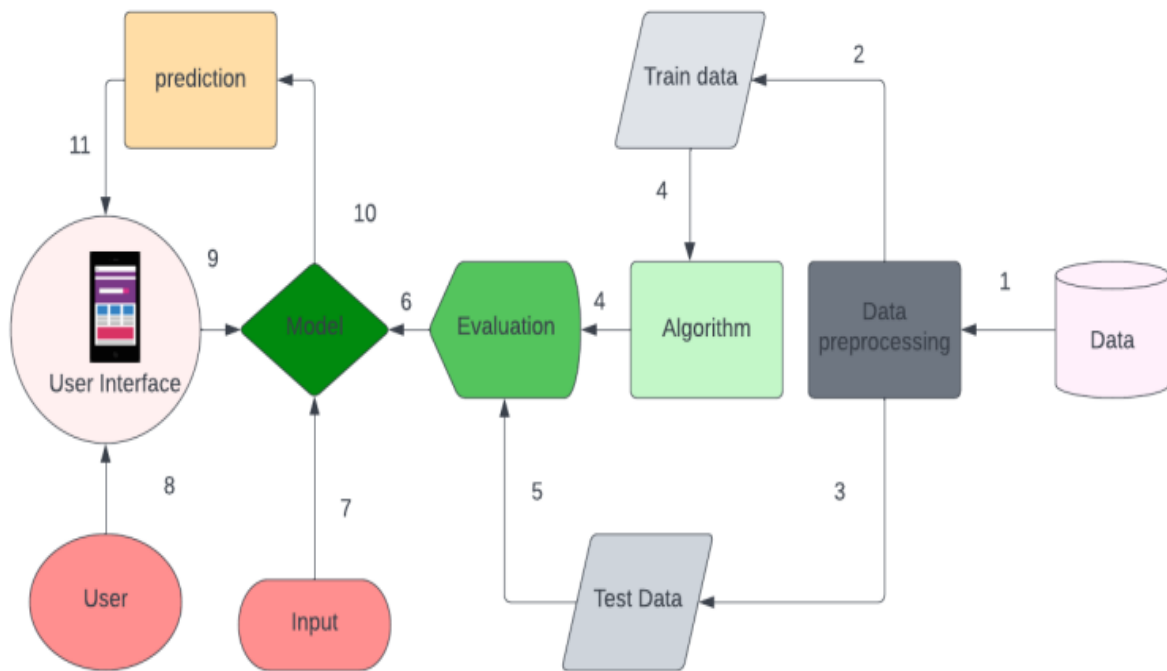
#### 5.1 DATA FLOW DIAGRAM



## 5.2 SOLUTION AND TECHNICAL ARCHITECTURE







### 5.3 USER STORIES

USER TYPE	USER STORY NUMBER	USER STORY/TASK	ACCEPTANCE CRITERIA	PRIORITY	RELEASE
Customer	USN-1	I can use this web app for flight delay prediction	I am getting the result	Medium	Sprint1
	USN-2	As a tourist person, I can able to get the accurate result.	As a user I can able to access the dashboard.	High	Sprint2

Customer (Web user)	USN-3	As a user, I can use the web application virtually anywhere	I can use the application in any device with a browser.	Medium	Sprint3
Administrative Management	USN-4	As an administrative I would provide all the IT support	Allows growth and success of the website	High	Sprint-3

## CHAPTER 6

### PROJECT PLANNING AND SCHEDULING

#### 6.1 SPRINT PLANNING AND ESTIMATION

<b>Sprint</b>	<b>Functional Requirement (Epic)</b>	<b>User Story Number</b>	<b>User Story / Task</b>	<b>Story Points</b>	<b>Priority</b>	<b>Team Members</b>
Sprint-1	Data Engineering	USN-1	Data Collection, Data Pre-processing and Feature Extraction	4	High	Tharanyaa R
Sprint-2	Machine Learning Prediction Model	USN-2	Building a Machine Model for Flight Delay Prediction.	4	High	Vijay S
Sprint-3	Flask Web Page	USN-3	Building Home Page.	4	Medium	Gayathri P
Sprint-4	Integration.	USN-4	Integrating the flask pages with the ML Model and IBM Cloud Deployment	4	Medium	Akashram J

#### 6.2 SPRINT DELIVERY SCHEDULE

<b>Sprint</b>	<b>Total Story Points</b>	<b>Duration</b>	<b>Sprint Start Date</b>	<b>Sprint End Date (Planned)</b>	<b>Story Points Completed (as on Planned End Date)</b>	<b>Sprint Release Date (Actual)</b>
Sprint-1	20	6 Days	24 Oct 2022	29 Oct 2022	20	29 Oct 2022
Sprint-2	20	6 Days	31 Oct 2022	05 Nov 2022	20	05 Nov 2022
Sprint-3	20	6 Days	07 Nov 2022	12 Nov 2022	20	12 Nov 2022

Sprint-4	20	6 Days	14 Nov 2022	19 Nov 2022	20	19 Nov 2022
----------	----	--------	-------------	-------------	----	-------------

## 6.3 REPORTS FROM JIRA

The screenshot displays the Jira 'FDP board' for the 'Flight delay prediction' software project. The interface includes a top navigation bar with options like 'Your work', 'Projects', 'Filters', 'Dashboards', 'People', 'Apps', and a 'Create' button. A search bar is also present. The left sidebar shows the project's navigation menu with 'Board' selected under the 'PLANNING' section. The main area shows the 'FDP board' with three columns: 'TO DO 1 ISSUE', 'IN PROGRESS 1 ISSUE', and 'DONE 1 ISSUE'. Each column contains one issue card with a title, a checkbox, and a Jira ID (FDP-1, FDP-2, FDP-3). The 'TO DO' column has the issue 'I faced package conflict'. The 'IN PROGRESS' column has 'Develop a Website'. The 'DONE' column has 'Deploy model in Cloud'. Each issue card has a '+ Create issue' button below it. The board is grouped by 'None'.

## **CHAPTER 7**

### **CODING AND SOLUTIONING**

We completed four sprints—Sprint 1, Sprint 2, Sprint 3 and Sprint 4—during the project development phase

#### **7.1 Sprint 1**

The dataset has been downloaded. The features are analysed and visualized and data has been cleaned and pre-processed. The independent and dependent variables are then identified and the dataset is split into train and test sets.

#### **7.2 Sprint 2**

Several machine learning algorithms have been applied for classification like logistic regression, K means, naïve bayes and random forest classifier and it is found that logistic regression gives the highest accuracy, so it is used for deployment.

#### **7.3 Sprint 3**

We had done building HTML files, written Python code, and running the application during Sprint 2.

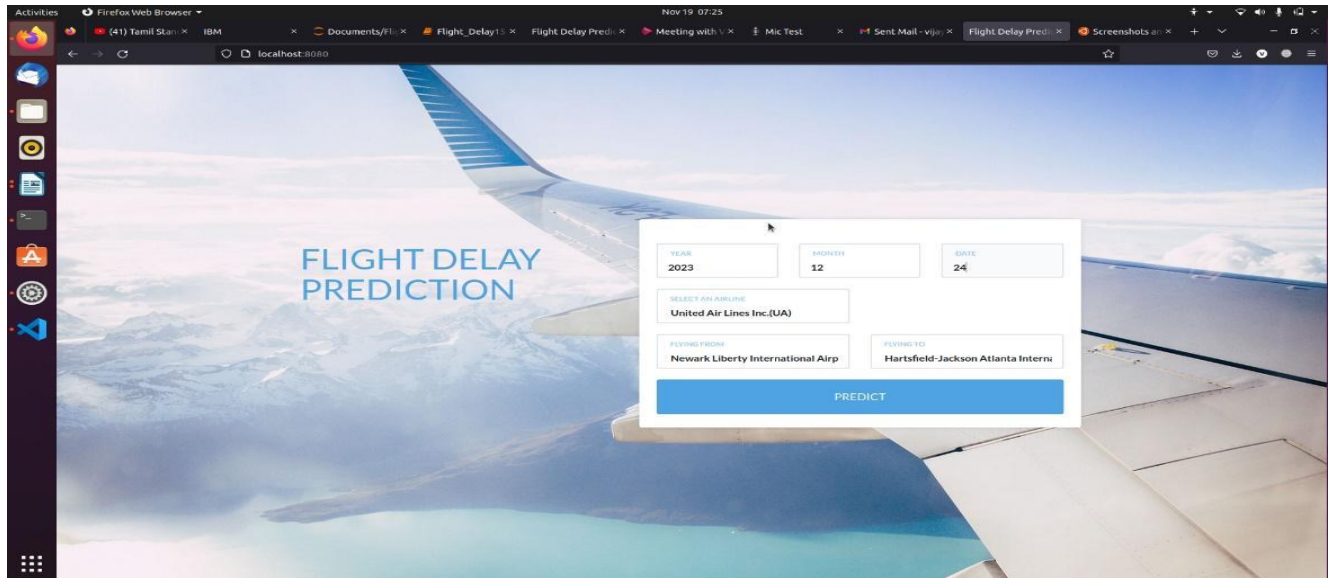
#### **7.4 Sprint 4**

We trained the ML model on IBM and integrated the flask. Registered on IBM cloud and activated Watson machine learning, cloud storage and Watson studio then trained the ML model on IBM using API KEY during sprint 4.

## CHAPTER 8

### TESTING

#### 8.1 TEST CASES



#### 8.2 USER ACCEPTANCE TESTING

Resolution	Severit y 1	Severit y 2	Severit y 3	Severit y 4	Subtot al
By Design	1	0	1	0	2
Duplicate	0	1	0	0	1
External	0	0	2	0	2
Fixed	4	1	0	0	5
Not Reproduced	0	0	1	1	2
Skipped	0	0	0	1	1
Won't Fix	1	0	0	0	1
Totals	6	2	4	2	14

Section	Total Cases	Not Tested	Fail	Pas s
Client Application	9	0	1	8
Security	2	0	0	2
Exception Reporting	4	0	1	3
Performance	4	0	0	4

## CHAPTER 9

### RESULTS

#### 9.1 PERFORMANCE METRICS

Model: Logistic Regression performance values

There is no big variation in the training and testing accuracy. Therefore, the Logistic Regression model is not overfit or underfit.

```
In [24]: print("Train set Accuracy: ", metrics.accuracy_score(y_train, LR.predict(X_train)))  
         print("Test set Accuracy: ", metrics.accuracy_score(y_test, LR.predict(X_test)))  
  
Train set Accuracy:  0.72045166414639  
Test set Accuracy:  0.7202230029020925
```

Model: Naive Bayes performance values

There is no big variation in the training and testing accuracy

```
In [13]: print("Train set Accuracy: ", metrics.accuracy_score(y_train, gnb.predict(X_train)))  
         print("Test set Accuracy: ", metrics.accuracy_score(y_test, y_pred))  
  
Train set Accuracy:  0.7207321224393608  
Test set Accuracy:  0.7196448209848886
```

Model: K means performance values

There is no big variation in the training and testing accuracy

```
In [14]: k = 6  
         neigh6 = KNeighborsClassifier(n_neighbors=k).fit(X_train, y_train)  
         yhat6 = neigh6.predict(X_test)  
         print("Train set Accuracy: ", metrics.accuracy_score(y_train, neigh6.predict(X_train)))  
         print("Test set Accuracy: ", metrics.accuracy_score(y_test, yhat6))  
  
Train set Accuracy:  0.7792924895752188  
Test set Accuracy:  0.7262257522529403
```

Model: Random Forest performance values

There is slight variation in the training and testing accuracy

```
In [21]: print("Train set Accuracy: ", metrics.accuracy_score(y_train, clf.predict(X_train)))  
         print("Test set Accuracy: ", metrics.accuracy_score(y_test, clf.predict(X_test)))  
  
Train set Accuracy:  0.8363318656342538  
Test set Accuracy:  0.7019318889988636
```

On comparing the four models built, based on the performance metrics it is clear that random forest gives the highest performance. Hence, that model is chosen for deployment.



## **CHAPTER 10**

### **10.1 ADVANTAGES**

- The application is fast and offers great accuracy in predicting the flight delay.
- Less maintenance is required.
- It is user friendly.
- It helps in reducing the tension of the passengers in knowing how long they will have to wait and lets passengers plan their schedule accordingly, thus in a way saving their time

### **10.2 DISADVANTAGES**

- It requires an internet connection for the website to work.

## **CHAPTER 11**

### **CONCLUSION**

From this study, we have developed a web application model that shows the flight delay prediction. In particular, by applying random forest algorithm to the prediction model, a reliable delay status of a single day could be acquired. Once the model was built it was integrated along with the Flask framework so that the users can enter their flight details and see if the flight would be on time or get delayed. Then this model is trained and deployed in the IBM Cloud.

As a result, anticipating delays can enhance airline operations and passenger satisfaction, which will benefit the economy and bring a positive impact.

## **CHAPTER 12**

### **FUTURE SCOPE**

The next steps are to apply other algorithms to the prediction and analyse the task of flight delays. It may yield important patterns and accuracy in flight delay data.

Web application can further be improved in which notification is sent via message or mail and allowing administrators to verify the identity of the user.

A section where the users can give their feedback can also be implemented.

## **CHAPTER 13**

### **APPENDIX**

#### **GITHUB LINK**

<https://github.com/IBM-EPBL/IBM-Project-7979-1658904775.git>

#### **PROJECT DEMO LINK**

[https://drive.google.com/file/d/1YxQk\\_UI7crIYlI4yUGX3uFkZ6InghEUu/view?usp=share\\_link](https://drive.google.com/file/d/1YxQk_UI7crIYlI4yUGX3uFkZ6InghEUu/view?usp=share_link)