

# Prediction of Flight Arrival Delay Time with Machine Learning Approaches on the U.S. Bureau of Transportation Statistics

Jiarui LI

University of Nottingham, Ningbo China  
scyj16@nottingham.edu.cn

Ran JI

University of Nottingham, Ningbo China  
scyrj1@nottingham.edu.cn

Yik Lun Yelan LAU

University of Nottingham, Ningbo China  
scyy118@nottingham.edu.cn

**Abstract**—According to the data from the Bureau of Transportation Statistics (BTS), the number of passengers and flights has been increasing year by year. However, flight delay has become a pervasive problem in the United States in recent years due to various factors, including human factors such as security regulations, as well as natural factors such as bad weather. Flight delay not only affects the profits of airlines but also affects the satisfaction of passengers. Therefore, a model that can predict the arrival time of airplanes needs to be developed. Machine learning methods have been widely applied to prediction problems. In this paper, a variety of machine learning methods, including linear regression, decision tree (DT), random forest (RF), gradient boosting (GB), and gaussian regression models were trained on the U.S. Department of Transportation's (DOT) BTS dataset. The results show that the (model name) performs best and can be used to predict the arrival time of the U.S. flights in advance, which is beneficial for airlines and passengers to make timely decisions.

**Index Terms**—Big data, Air flight, Airport, Delay, Machine learning, Spark, Regression

## I. INTRODUCTION

Big data analysis has been successfully applied in many fields, including but not limited to biology [1], healthcare [2], geography [3], traffic and transportation [4], and the Internet [5]. Using big data analysis can help with discovering new information and making better informed decisions [6], [7]. While travelling by air is becoming a more popular option, illustrated by the growth in passenger count between 2016 and 2019 in the Bureau of Transportation Statistics (BTS) data [8], the application of big data in the aviation field is still quite limited [4].

One serious problem airlines and travellers share alike with aviation is flight delays. In 2019, only 79% of flights in the United States (US) arrived on time. Delays cause crucial loss of time, money and resources for airlines and passengers [4], [8]–[10], and are estimated to cost \$30 billion for all parties involved in 2019 in the US alone [11]. Other than losses directly caused by delays, there are also indirect costs such as the inevitable decrease in demand due to unreliable service, which can amount to an even larger amount of money [4].

Flight delays are generally classified as five main reasons, air carrier, extreme weather, national aviation system, late-arriving aircraft and security [12]–[15]. Therefore, the effective prediction of flight delays enables airlines and traffic

management departments to make timely adjustments, which will have a positive effect on passengers and various aviation departments and reducing losses. Machine learning has been a hot-field method, a good quantum of research uses machine-learning method to predict flights.

In this work, a model for predicting arrival delays based on flight departure information is proposed. The dataset used to train and test our model includes the flight data of the relevant airports in the states of the United States obtained from the DOT Bureau of Transportation Statistics database. It mainly contains (data column). In the prediction process, the preprocessed flight departure data are used to predict the possible arrival difference and delay information of the flight through the model. The machine learning algorithms implemented in this work include Linear Regression, Decision Tree (DT), Random Forest (RF), Gradient Boost (GB) and Gaussian Regression.

The rest of this report is organized as follows: Section II shows the related work in flight delay prediction. Section III illustrates the problem statement of our work. Section IV illustrates the dataset and preprocessing step in this work. In the next part of Section IV, the procedure of data preprocessing is discussed. In the third part of Section IV, models are described. Section V discusses the results of the model. Section VI provides the conclusion of our work.

## II. LITERATURE REVIEW

As mentioned in the introduction, flight delays are a common problem encountered by airlines at present. When the arrival or departure time of a flight exceeds the estimated or scheduled time, the flight will be delayed.

Attempts to predict delays have been made by exploring patterns in air traffic. Rebello and Balakrishnan et al. [15] created variables indicative of NAS status and used systematic dependencies between airports to predict future network-related delays. Klein et al. [16] focus on prediction from weather conditions and propose a delay prediction model based on the Weather Impacted Traffic Index (WITI), which can be used to evaluate the severity and impact of weather. Belcastro et al. [17] proposed a method based on big data. A MapReduce program using parallel algorithms is executed on the cloud platform by analyzing and mining flight information

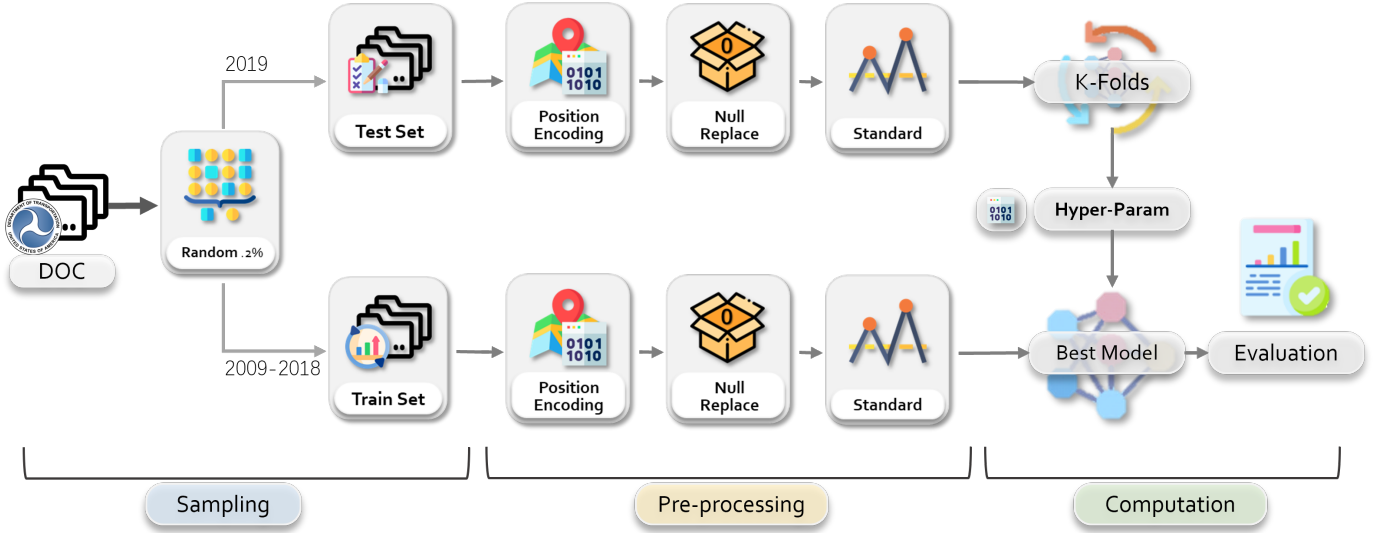


Fig. 1. The data process and predict procedure.

and weather conditions, and it is used to predict flight delays caused by weather.

The latest research has also focused on using machine learning methods to predict delays and cancellations. All studies have concluded that there is a close relationship between arrival delays and departure delays [18]. Therefore, studies often use departure delays as inputs to predict arrival delays, although some studies predict arrival delays by using information of weather condition as inputs. Some studies categorize flight delay prediction as a classification problem, while others categorize delay prediction as a regression problem, predicting delays in minutes.

Jiang et al. [4] used Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), Multilayer Perceptron (MLP) and data visualization methods to solve the two-category problem of flight delays and the classification task of five-level delays. Alonso et al. [18] proposed a unimodal model established through machine learning methods to classify the length of flight delays for arrivals at Porto Airport based on given airport departure flight information. They transformed the regression problem of delay time into a classification problem of time intervals. The model achieved a coefficient result of 0.7 for its prediction outcomes.

In the research of [13], a two-stage model is proposed, which firstly solves the classification problem of delay, and then predicts the actual delay time by regression methods. Xu et al. [19] used Bayesian network to analyze flight delays in two lines of the National Airspace System according to weather conditions. Choi et al. [12] also used supervised machine learning methods such as DT, RF, and AdaBoost to predict flight arrival delays based on weather conditions. In the studies of [20], [21], neural networks are used to predict flight delays for John F. Kennedy (JFK) and Esenboga International Airport respectively. Chakrabarty et al. [22] trained a Gradient Boosting model using flight details covering the top five

busiest airports in the United States and the Gradient Boosting Tree (GBT) model is applied to predict flight arrival delays. In reference [23], researchers used logistic regression and random forest models to predict flight delay status and duration based on flight data recorded by the BTS. In addition, they considered some non-flight attributes such as weather and peak tourism data. The study results showed that the model's prediction accuracy was between 80% and 85%.

### III. PROBLEM STATEMENT

Flight delays can result in significant financial losses for airlines, as well as inconvenience for passengers and issues with air traffic control systems. As a consequence, it is necessary to predict flight delay information in advance, in order to allow other passengers and air traffic control departments to be informed and take action proactively. Within the United States, machine learning regression models are utilized to predict flight delay information based on flight departure data. This technology can serve as a decision-making tool for airports and air traffic control departments.

### IV. METHODOLOGY

#### A. Data Collection

The information regarding domestic delayed flights in the United States between 2009 and 2019 is sourced from the U.S. DOT BTS, which is responsible for monitoring the punctuality of domestic flights that are operated by major air carriers. It mainly includes delay data for domestic routes in the United States, which also provides other related information, such as airline date, original airport, destination airport, scheduled and actual departure time, taxi-out time, taxi-in time, wheels-off time, wheels-on time, scheduled and actual arrival time, flight time and delay time by various reasons, etc.

There are 68,979,001 samples in this total dataset. The randomly chosen flight data from 2009 to 2018 is used as the

Flight Count in Each Year from 2009 to 2018 (Trainset)



Fig. 2. Flight count of each year in training set

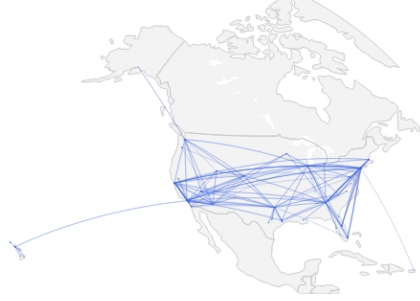


Fig. 3. Airline in training set demonstrates in the map.

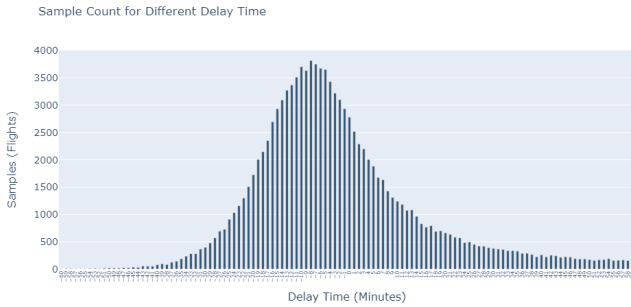


Fig. 4. Flight count of different delay level.

training set, and the randomly chosen data from 2019 is used as the test set. The train and test data set contains 121,513 and 14,861 flights, including flight data of 358 airports. The following are all the attributes and descriptions in the dataset:

TABLE I  
SELECTED ATTRIBUTES AND THEIR TYPES (PART 1).

Attributes	Type	Comment
Flight Date(mmdd)	String	The date of the flight.
Carrier	String	Unique Carrier Code of different carriers.
Carrier Number	Integer	An identification number assigned by US DOT to identify a unique airline (carrier).
Original Location	String	The identifier of original airport.
Destination Location	String	The identifier of of arrival airport.
Scheduled Departure Time	Float	Time (min) of the flight scheduled departure time.
Actual Departure Time	Float	Time (min) of the flight actual departure time.
Taxi-out Time	Integer	Time (hhmm) of the flight taxing-out time.
Taxi-in Time	Integer	Time (hhmm) of the plane putting away wheels.
Wheels-off Time	Integer	Time (hhmm) of the plane putting away wheels.
Wheels-on Time	Integer	Time (hhmm) of the plane putting down wheels.
Delay Departure Time	Integer	Delay time (hhmm) of the flight departure.
Scheduled Arrival Time	Integer	Time (min) of the flight scheduled arrival time.
Actual Arrival Time	Integer	Time (min) of the flight actual arrival time.

## B. Data Pre-processing

To enable machine learning models to recognize the data, four data pre-processing methods are applied. Firstly, instead of encoding airport identifiers as numbers as in previous research [3], the airports are transformed to their longitude and latitude information to enable the models recognize the location of the airport better. Secondly, the operation carrier identifier is encoded into unique integer identifiers using the string indexer embedded by Spark. Next, null values in features such as Weather Delay Time, Carrier Delay Time, NAS Delay Time, Security Delay Time, and Late Aircraft Delay Time are replaced with 0 to identify flights that did not experience delays under those conditions. Finally, all features are standardized by removing the mean and scaling to unit variance.

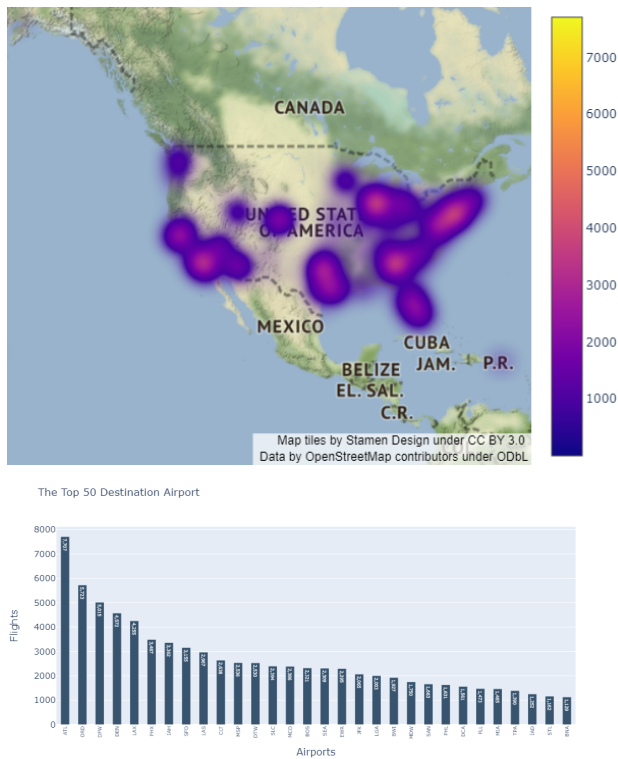
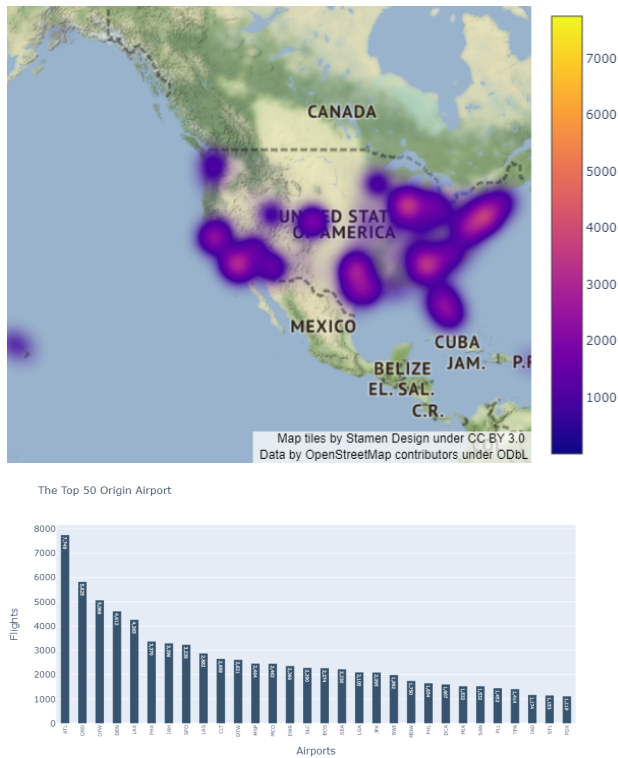


TABLE II  
SELECTED ATTRIBUTES AND THEIR TYPES (PART 2).

Attributes	Type	Comment
Carrier Delay Time	Integer	Time (min) of the flight delay by carrier factors. Carrier Delay is within the control of the air carrier, such as aircraft cleaning, aircraft damage, awaiting the arrival of connecting passengers or crew, baggage, bird strike, etc.
Cancelled	Integer(1 or 0)	The flight was cancelled or not.
Cancelled Code	Char (A,B,C,D)	The reason for cancelling the flight.
Weather Delay Time	Integer	Time (min) of the flight delay by weather factors.
NAS Delay Time	Integer	Time (min) of the flight delay by national air system factors, such as non-extreme weather conditions, airport operations, heavy traffic volume, etc.
Security Delay Time	Integer	Time (min) of the flight delay by security factors, such as evacuation of a terminal or concourse, re-boarding of aircraft because of security breach, etc.
Late Aircraft Delay	Integer	Time (min) of the flight delay by the late arrival of the previous flight that utilized the time of the same plane that will be departing.
Distance	Float	Distance of the airline.

After feature selection, We extracted the following information as input information for the dataset because these factors are considered to be effective factors affecting flight delays in our work and related works:

### C. Prediction Models

To solve the predicting flight arrival time and judging delay information problems, as mentioned in the literature review, the following popular algorithms are used for training and tests:

- Linear Regression [24]
- Decision Tree (DT) [8], [25]
- Random Forest (RF) [9], [26]
- Gradient Boosting (GB) [27]
- Gaussian Regression [23], [28]

To improve results, the following methods are applied:

- 1) *K-Fold Cross-validation* [29]: Cross-validation is a measure for evaluating model performance. K-fold cross-validation is a technique that randomly divides the original sample into K sub-samples. In our work, we split it into 3 folds. Then, a single sub-sample is used as validation data for testing the model, and the remaining 2 sub-samples are used as training data. These processes

TABLE III  
SELECTED ATTRIBUTES AND THEIR TYPES.

Attributes	Type	Comment
Flight Date(mmdd)	Integer	The flight date of this flight.
Carrier ID	Integer	An identification number assigned by US DOT to identify a unique airline (carrier).
Original Location	(Float, Float)	The coordinator of the original airport including latitude and longitude.
Destination Location	(Float, Float)	The coordinator of the destination airport including latitude and longitude
Scheduled Departure Time (min)	Float	Time (min) of the flight scheduled departure time.
Departure Delay Time (min)	Float	Delay time (hhmm) of the flight departure. Type in Integer.
Distance	Float	Distance of the airline.
Carrier Delay Time (min)	Integer	Time (min) of the flight delay by carrier factors. Carrier Delay is within the control of the air carrier, such as aircraft cleaning, aircraft damage, awaiting the arrival of connecting passengers or crew, baggage, bird strike, etc.
Weather Delay Time (min)	Integer	Time (min) of the flight delay by weather factors.
NAS Delay Time (min)	Integer	Time (min) of the flight delay by national air system factors, such as non-extreme weather conditions, airport operations, heavy traffic volume, etc.
Security Delay Time (min)	Integer	Time (min) of the flight delay by security factors, such as evacuation of a terminal or concourse, re-boarding of aircraft because of security breach, etc.
Late Aircraft Delay Time (min)	Integer	Time (min) of the flight delay by the late arrival of the previous flight that utilized the time of the same plane that will be departing.

are repeated 3 times. As the original dataset is large, as mentioned in the previous section, we performed random selection to handle it. Additionally, the k-fold cross-validation procedure is utilized to mitigate any bias caused by the random selection of data [30].

- 2) *Hyper-Parameter Tuning* [31]: To ensure a good fit of the machine learning model to the data, hyperparameters must be carefully configured prior to training. In our approach, several sets of hyperparameters are tested using K-Fold validation, and the set of parameters that yields the best performance is selected as the final hyperparameter configuration. For the Decision Tree and Random Forest models, we test different values for the maximum depth of the tree, including 5, 6, 7, and 8. For the Linear Regression and Gaussian Regression models, we test different convergence tolerances, including  $1e-4$ ,  $1e-5$ ,  $1e-6$ , and  $1e-7$ .

## V. RESULTS AND DISCUSSION

TABLE IV  
RESULTS ON TEST SET OF DIFFERENT MODELS IN THIS WORK

Regressor	RMSE	$R^2$ Score
Decision Tree	30.2331	0.6314
Random Forest	29.1984	0.6562
Gradient Boost	27.4730	0.6956
Linear Regression	14.5720	0.9144
Gaussian Regression	14.5720	0.9144

### A. Prediction Results

To analyze the performance of the models, root mean squared error (RMSE) and  $R^2$  score are used for evaluation regression results. As mentioned in previous section, training data is in 2009-2018, 10 years, and test data is in 2019. To evaluate the models, K-Fold Cross Validation (K=3) are

TABLE V  
RESULTS ON TRAIN SET OF DIFFERENT MODELS IN THIS WORK

Regressor	RMSE	$R^2$ Score
Decision Tree	17.9888	0.7935
Random Forest	17.4475	0.8058
Gradient Boost	18.9884	0.7700
Linear Regression	9.4382	0.9432
Gaussian Regression	9.4382	0.9432

Method Root Mean Square Error Comparison

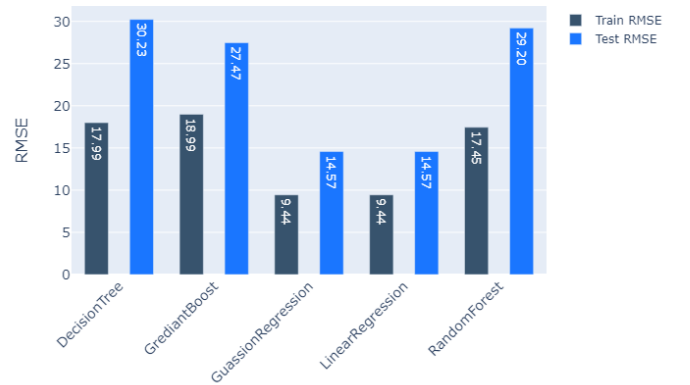


Fig. 7. Comparisons of different models on RMSE

performed, and the regression results of which are shown in Table V and Table IV.

The RMSE and  $R^2$  scores results of Linear Regression [24], Decision Tree [8], Random Forest [9], and Gaussian Regression [23] are shown in Fig 7, Table IV and Table V. Linear Regression and Gaussian Regression are shown outperforms than other two methods, with the RMSE 14.5720 and  $R^2$  score 0.9144.



## B. Findings

Based on the results, the Decision Tree based models performs unsatisfactorily, because of the discrete property of these models. This also can be seen on the performance difference among Decision Tree, Random Forest, and Gradient Boost. With more trees adapted, the model can fit the data more smoothly and achieve better result. Compared to Decision Tree based models, linear regression can predict the time difference of airplane arrival to some extent, because it can fit the data pattern smoothly. Additionally, the time and distance information of the flight departure can be used as inputs to predict flight delay. From the literature review, weather information such as temperature, humidity, and air pressure at departure time could also be used as input features for prediction, but this dataset lacks such data. This can be considered as future work. As mentioned in the introduction, flight delay is a prevalent problem, and the predictive model we propose can help airlines and passengers make timely decisions, thus reducing losses caused by flight delays to some extent.

## VI. CONCLUSION

In our work, we proposed that departure information of flights can be used to predict flight delay time. To handle the large-scale dataset efficiently, we applied Spark for data processing and distribution among the processing members, thus reducing the overall processing time. The proposed Linear Regression performs better than other our evaluated methods. The results can prove that flight delay time can be predicted to reduce loss of airlines and passengers for economy.

## VII. ADDITIONAL WORK

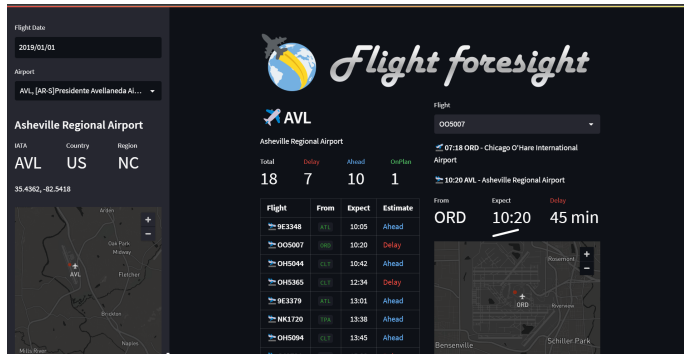


Fig. 8. The screenshot of the flight airline delay prediction visualization board.

This project is aimed to tackle the influence caused by airline arrival difference to the airline control and destination scheduling. To enable our project can contribute to the industry and daily life, a visualization client named flight foresight is comprised. It can provide the predict flight delay information for select date and destination airport. And users can check the status of any flight they want.

## REFERENCES

- [1] V. Marx, "The big challenges of big data," in *Nature*, no. 498, 2013, p. 255–260.
- [2] W. Raghupathi and V. Raghupathi, "Big data analytics in healthcare: promise and potential," *Health information science and systems*, vol. 2, pp. 1–10, 2014.
- [3] R. Kitchin, "The real-time city? big data and smart urbanism," *GeoJournal*, vol. 79, pp. 1–14, 2014.
- [4] Y. Jiang, Y. Liu, D. Liu, and H. Song, "Applying machine learning to aviation big data for flight delay prediction," in *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/CBDCCom/CyberSciTech)*, 2020, pp. 665–672.
- [5] C. Snijders, U. Matzat, and U.-D. Reips, "big data": big gaps of knowledge in the field of internet science," *International journal of internet science*, vol. 7, no. 1, pp. 1–5, 2012.
- [6] Z. Lv, H. Song, P. Basanta-Val, A. Steed, and M. Jo, "Next-generation big data analytics: State of the art, challenges, and future research topics," *IEEE Transactions on Industrial Informatics*, vol. 13, no. 4, pp. 1891–1899, 2017.
- [7] G. Dartmann, H. Song, and A. Schmeink, *Big data analytics for cyber-physical systems: machine learning for the internet of things*. Elsevier, 2019.
- [8] R. Boggavarapu, P. Agarwal, and R. K. D.H., "Aviation delay estimation using deep learning," in *2019 4th International Conference on Information Systems and Computer Networks (ISCON)*, 2019, pp. 689–693.
- [9] K. Sharma, R. L. Eliganti, B. S. K. Meghana, and G. Gayatri, "Error calculation of flight delay prediction using various machine learning approaches," in *2022 IEEE International Conference on Current Development in Engineering and Technology (CCET)*, 2022, pp. 1–5.
- [10] B. Thiagarajan, L. Srinivasan, A. V. Sharma, D. Sreekanthan, and V. Vijayaraghavan, "A machine learning approach for prediction of on-time performance of flights," in *2017 IEEE/AIAA 36th Digital Avionics Systems Conference (DASC)*. IEEE, 2017, pp. 1–6.
- [11] F. A. Administration, "Cost of delay estimates 2019," [https://www.faa.gov/data\\_research/aviation\\_data\\_statistics/media/cost\\_delay\\_estimates.pdf](https://www.faa.gov/data_research/aviation_data_statistics/media/cost_delay_estimates.pdf), 07 2020.
- [12] S. Choi, Y. J. Kim, S. Briceno, and D. Mavris, "Prediction of weather-induced airline delays based on machine learning algorithms," in *2016 IEEE/AIAA 35th Digital Avionics Systems Conference (DASC)*, 2016, pp. 1–6.
- [13] B. Thiagarajan, L. Srinivasan, A. V. Sharma, D. Sreekanthan, and V. Vijayaraghavan, "A machine learning approach for prediction of on-time performance of flights," in *2017 IEEE/AIAA 36th Digital Avionics Systems Conference (DASC)*, 2017, pp. 1–6.
- [14] S. Addu, P. R. Ambati, S. R. Kondakalla, H. Kunchakuri, and M. Thotempudi, "Predicting delay in flights using machine learning," in *2022 International Conference on Applied Artificial Intelligence and Computing (ICAIC)*, 2022, pp. 374–379.
- [15] J. J. Rebollo and H. Balakrishnan, "Characterization and prediction of air traffic delays," *Transportation research part C: Emerging technologies*, vol. 44, pp. 231–241, 2014.
- [16] A. Klein, C. Craun, and R. S. Lee, "Airport delay prediction using weather-impacted traffic index (witi) model," in *29th Digital Avionics Systems Conference*. IEEE, 2010, pp. 2–B.
- [17] L. Belcastro, F. Marozzo, D. Talia, and P. Trunfio, "Using scalable data mining for predicting flight delays," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 8, no. 1, pp. 1–20, 2016.
- [18] H. Alonso and A. Loureiro, "Predicting flight departure delay at porto airport: A preliminary study," in *2015 7th International Joint Conference on Computational Intelligence (IJCCI)*, vol. 3, 2015, pp. 93–98.
- [19] N. Xu, G. Donohue, K. B. Laskey, and C.-H. Chen, "Estimation of delay propagation in the national aviation system using bayesian networks," in *6th USA/Europe Air Traffic Management Research and Development Seminar*. FAA and Eurocontrol Baltimore, 2005.
- [20] S. Khanmohammadi, S. Tutun, and Y. Kucuk, "A new multilevel input layer artificial neural network for predicting flight delays at jfk airport," *Procedia Computer Science*, vol. 95, pp. 237–244, 2016.
- [21] E. Demir and V. B. Demir, "Predicting flight delays with artificial neural networks: Case study of an airport," in *2017 25th Signal Processing and Communications Applications Conference (SIU)*, 2017, pp. 1–4.

- [22] N. Chakrabarty, "A data mining approach to flight arrival delay prediction for american airlines," in *2019 9th Annual Information Technology, Electromechanical Engineering and Microelectronics Conference (IEMECON)*. IEEE, 2019, pp. 102–107.
- [23] V. Natarajan, S. Meenakshisundaram, G. Balasubramanian, and S. Sinha, "A novel approach: Airline delay prediction using machine learning," in *2018 International Conference on Computational Science and Computational Intelligence (CSCI)*, 2018, pp. 1081–1086.
- [24] Y. Ding, "Predicting flight delay based on multiple linear regression," in *IOP conference series: Earth and environmental science*, vol. 81, no. 1. IOP Publishing, 2017, p. 012198.
- [25] I. H. Witten and E. Frank, "Data mining: practical machine learning tools and techniques with java implementations," *Acm Sigmod Record*, vol. 31, no. 1, pp. 76–77, 2002.
- [26] K. Sharma and T. S. Kiranmai, "Prediction of cardiovascular diseases using genetic algorithm and deep learning techniques," *INTERNATIONAL JOURNAL OF EMERGING TRENDS IN ENGINEERING AND DEVELOPMENT*, vol. 3, no. 10, p. 26808, 2021.
- [27] S. Manna, S. Biswas, R. Kundu, S. Rakshit, P. Gupta, and S. Barman, "A statistical approach to predict flight delay using gradient boosted decision tree," in *2017 International conference on computational intelligence in data science (ICCIDS)*. IEEE, 2017, pp. 1–5.
- [28] Z. Chen, D. Guo, and Y. Lin, "A deep gaussian process-based flight trajectory prediction approach and its application on conflict detection," *Algorithms*, vol. 13, no. 11, p. 293, 2020.
- [29] K. Ito and R. Nakano, "Optimizing support vector regression hyperparameters based on cross-validation," in *Proceedings of the International Joint Conference on Neural Networks, 2003.*, vol. 3, 2003, pp. 2077–2082 vol.3.
- [30] H. Moayedi, A. Osouli, H. Nguyen, and A. S. A. Rashid, "A novel harris hawks' optimization and k-fold cross-validation predicting slope stability," *Engineering with Computers*, vol. 37, pp. 369–379, 2021.
- [31] P. Probst, A.-L. Boulesteix, and B. Bischl, "Tunability: Importance of hyperparameters of machine learning algorithms," *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1934–1965, 2019.