CENTERIS - International Conference on ENTERprise Information Systems / ProjMAN - International Conference on Project MANagement / HCist - International Conference on Health and Social Care Information Systems and Technologies, CENTERIS/ProjMAN/HCist 2018

# Predictive Modelling: Flight Delays and Associated Factors, Hartsfield–Jackson Atlanta International Airport

Roberto Henriques, Inês Feiteira *

*aNova Information Management School, Universidade Nova de Lisboa Campus de Campolide, Lisboa 1070-312, Portugal*

## Abstract

Nowadays, a downside to traveling is the delays that are constantly being advertised to passengers resulting in a decrease in customer satisfaction and causing costs. Consequently, there is a need to anticipate and mitigate the existence of delays helping airlines and airports improving their performance or even take consumer-oriented measures that can undo or attenuate the effect that these delays have on their passengers. This study has as main objective to predict the occurrence of delays in arrivals at the international airport of Hartsfield-Jackson. A Knowledge Discovery Database (KDD) methodology was followed, and several Data Mining techniques were applied. Historical data of the flight and weather, information of the airplane and propagation of the delay were gathered to train the model. To overcome the problem of unbalanced datasets, we applied different sampling techniques. To predict delays in individual flights we used Decision Trees, Random Forest and Multilayer Perceptron. Finally, each model's performance was evaluated and compared. The best model proved to be the Multilayer Perceptron with 85% of accuracy.

*Keywords:* Data Mining; Predictive Analysis; Flight Delays; Hartsfield–Jackson International Airport; Atlanta International Airport

* Corresponding author. Tel.: +351 213 828 610
  E-mail address: m2015283@novaims.unl.pt

## 1. Introduction

The airline industry has grown over the years, approximately 5% per year over the last 30 years [1]. Its demand grew steadily with a global passenger air traffic growth of 6.5% in 2015, far above of the 10-year average annual growth of 5.5% [2].

However, a high demand does not only translate into success since it can diminish the capacity to respond and deal with it. As a result of the high demand is the congestion of the system caused by the disproportional growth between flights and airport capacity [3]. Consequently, this volume, together with several factors (mechanical problems, atmospheric conditions, traffic control issues, delay propagation, labor scarcity and Ground Delay Programs (GDP) among others), promotes two possible scenarios. In the best one, it is possible to exchange aircraft between flights or to request aircraft not currently used. In the worst scenario, the cancellation and or delay of the flight happens [4]. For the Bureau of Transportation Statistics (BTS) [5] a flight delay is defined as a flight that is late 15 or more minutes than the scheduled time.

A study, developed by Ball et al. [6], estimated the total impact of the costs of delays in the American economy by 32.9 billion dollars. This cost affects the country's Gross Domestic Product (GDP) indirectly, and the passengers, airlines and lost demand, directly.

The concept of traveling has been shaping over time. In the past was seen as a privilege; today, travel often represents a necessary evil – the result of air delays, increased security measures and degradation of services provided [6] – in achieving goals of greater value to the passenger, such as sightseeing or even work.

The analysis of air delays becomes important since a better knowledge of their existence and triggers can improve the performance of airlines and, consequently airports in their operations by the possibility of anticipation, construction of schedules, etc.

In this study we have focused in the analysis of the delays with emphasis on the arrivals since these are more related with the passenger's satisfaction and because one arrival delay may trigger a delay in a departure [7]. The Hartsfield-Jackson International Airport (ATL) in Atlanta is on the top of the passenger traffic results for the busiest airports in 2015 having its strategic location as a major "gateway" to entry into North America and is reported to be two hours away from 80% of the population of the USA [8].

To predict the delays in the arrivals at ATL several approaches were implemented such as the 1) use of different sampling techniques (SMOTE and Undersampling); 2) the inclusion and exclusion of outliers; 3) the application of different attribute selection methods (GainRatio and Correlation-based Feature Selection); 4) the inclusion and exclusion of two variables (delay in departure and real departure time of a flight), and finally; 5) by applying three different supervised machine learning algorithms (Decision Trees, Random Forest and Multilayer Perceptron).

The paper is organized as follows: Section 2 provides a brief explanation on the need for KDD, Data Mining, and Machine Learning when dealing with issues with large volumes of data and, state of the art about the study of air delays, both at an industrial community level as well as at the scientific community level. Section 3 describes the various stages defined for the methodology followed in this work. In section 4 the classifiers' performance are analyzed, the best approaches are selected for each type of advance of the prediction and compared to related works in the same area of research and with the same target type variables but also to websites that predict delays. Lastly, conclusion and space for future works are presented in section 5 which concludes the paper.

## 2. Background and Related Work

Since 2010, the amount of delays has been fluctuating, with 2014 having the highest percentage of delay. In 2015, the percentage of delayed flights was 19.53%, the third highest value since 2010 [9]. Over time it has gained great importance and, as in other areas, there is an exponential growth of data that overload us and that we cannot control. As result, a gap between the production of data and our understandability of it emerges [10]. To acquire knowledge from data the use of tools is required to allow the discovery of hidden information in databases being Knowledge Discovery in Databases (KDD) the provider of it consisting in the process of discovering new, valid, useful and perceptible patterns in the data [11]. Furthermore, Data Mining (DM) is considered one of the steps of KDD [11–14] and represents the application of algorithms to extract patterns from the data, translating it in information [11]. Consequently, Machine Learning (ML) is a type of approach to the discovery of knowledge in the data. Focus on the

construction of computational algorithms that can learn through data (from the past) to make predictions [10] and when associated with DM tools can help to perceive complex phenomena as well as solve various problems from many areas [15].

The importance of flight delays leads to many investigations, both industrial and scientific. An example of the first strand is Kaggle. Kaggle is a platform for analysis and predictive modeling competitions with monetary awards for people with interest in the area [16].

For the purposes of this study, the focus of analysis is in on the use of machine learning techniques to predict the delay in individual flights by a classification problem [15,17,18].

Y. J. Kim, Choi, Briceno, & Mavris [17], in addition to predicting the status of the day-to-day delay at an airport level used neural networks to be able to predict the class of the delay of an individual flight. The authors gathered data about the flight and the weather conditions. Furthermore, they have allied the status of the day delay of the airport computed in the first stage of their study. Their results showed that their model achieved 87.42% accuracy, better than the best predictions until then demonstrated, of 83.4% [18] and 81% [19].

Choi, Kim, Briceno, & Mavris [18] also presented a classification model in which the main objective was to predict flight delays of individual flights caused by climatic conditions. They used data collected from the On-time Performance dataset for the years of 2005 to 2015 using features like flight schedules and day. The authors also included weather variables obtained by the Integrated Surface database of the National Oceanic and Atmospheric Administration (NOAA), at the origin and destination locations. Random Forest was the best classifier considered, and they concluded that the results with real climatic conditions, with no forecast horizon, have better accuracy (80.36%). The authors also mentioned that accuracy is higher when sampling technique is not applied.

In the same line of thought, Belcastro, Marozzo, Talia, & Trunfio [15] applied a parallel version of the Random Forest algorithm where the main objective was to be able to predict, with a few days in advance, the delay in the arrival of an individual flight due to the weather. They used data about the flight information from the airline On-time Performance dataset comprising the years of 2009 till 2013, as well as weather variables at the origin and destination acquired from the Quality Controlled Local Climatological Data (QCLCD). Considering a threshold of 15 minutes, they achieved 74.20% of accuracy. Nonetheless, the authors stated that not considering the weather features, the model would achieve only an accuracy of 69.1%.

The work presented in this paper, implements a set of algorithms to predict the delay in an individual flight taking into account (1) flight information, similar to [17]; (2) climate at origin and destination, similar to [15] and [18], but also; (3) information of the aircraft as well as (4) possible congestion of the system. The goal is to understand if it is possible, to improve the performance of the models already presented in this type of approach.

## 3. Methodology

### 3.1. Selection

This step was used as a basis for the discovery of standards [11] where the data needed are targeted, selected and gathered. Data about U.S. domestic airline traffic with arrival at ATL and weather data for 2 months of each season of the year of 2015 (April, May, July, August, October, November, January and February) was collected from the Airline On-time Performance dataset of the Bureau of Transportation Statistics (BTS) and from the archive of airport weather observations of the Iowa Environmental Mesonet, respectively. The weather information was collected 1) at the origin, at the schedule time of departure; 2) at the destination, on the schedule time of departure in origin; and, 3) at the destination, at the schedule time of arrival.

Also, information about the airplanes used in each flight is collected from the U.S. Civil Aviation Register database of the Federal Aviation Administration (FAA). Additional information about the time zone between origin-destination pairs for calculations of weather variables and time flight duration of routes was collected from the Travel Math website. For information about possible congestion, we used the Office Personnel Management to include the federal holidays for the year of 2015.

The dataset has 34 independent variables with a total of 248 956 observations. Considers different types of attributes concerning information of flight, airport of origin, airplane details, delay propagation/ congestion and weather factors. The inclusion of two specific independent variables, departure delay and real departure time, in our

study will affect directly the performance of the model because of the type of information that represents. For that reason, and because of its inclusion will not allow a possibility of a prediction before an airplane departs, the study of the inclusion and exclusion is implemented to see what approach improves the model performance.

As a dependent variable, we used the arrival delay as binary, where if the actual time of arrival is equal or greater than 15 minutes from the scheduled time the variable assumes the value 1. Otherwise, it assumes the value of 0, according to the Department of Transportation (DOT) definition [5].

Since we faced a problem of unbalanced classes with 86% of flights without delay and 14% of delayed flights and since this problem can lead to a false classification accuracy [20,21] we applied a sampling technique to minimize it. The goal is to have a dataset with a similar distribution of classes [21]. Synthetic Minority Over-sampling Technique (SMOTE) creates synthetic examples of the minority class being oversampled. It introduces synthetic examples along the line segments joining any/all of the *k* minority class nearest neighbors [20,21]. Its main advantage over traditional oversampling techniques is that instead of duplicating examples it generates new examples based on the existents. It is also possible to deal with this problem by implementing an undersampling technique to prevent common classes to hide rare classes [22] eliminating examples from the majority class. It has an advantage of turning the process faster with fewer observations but as disadvantage presents the fact that it can neglect useful information [23].

For the data partitioning, as it is necessary to see how a predictive model can succeed in unseen data, the holdout method was applied with a division of 70% for training and 30% for testing.

### 3.2. Data Pre-processing

By pre-processing data, all unnecessary information is eliminated and corrected to assure coherence [11]. For that reason, missing values [24] that can represent an issue in data quality and compromise the interpretation of data, must be treated. Among several possibilities to deal with missing values, we remove observations, impute by median (numerical variables) or mode (categorical variables) and the remove variables when they had a high percentage of missing values [25].

Extreme values that lie near the limits of the data range or go against the trend of the remaining data [25] are denominated as outliers [26,27]. Although sometimes, these values are correct (e.g., the case of extreme values of temperatures) and do not represent an outlier, in other cases, they can be considered outliers and can cause less accurate models [25]. In those cases, we need to identify them and test if their presence can improve or not the model.

### 3.3. Data Transformation

This step consists in transforming existing data and in dimensionality reduction [11]. Witten, Frank, and Hall [10] describe it as "engineering the input data into a form suitable for the learning scheme chosen and engineering the output to make it more effective". First, to increase the model's performance, the different scales of variables are scaled by applying the min-max method, and variables are rescaled to a range between 0 and 1 [25].

Next, we applied a variable selection step, which is "the process of identifying and removing as much of the irrelevant and redundant information as possible" [28]. Variables were selected through single attribute evaluator (Gain Ratio Attribute Evaluator with Ranker search method [10]) that evaluates attributes by measuring their gain ratio concerning the class and ranking them. We also used the filter method (Correlation-based Feature Selection with BestFirst search method [10]) from the attribute subset evaluator that takes into consideration the variables correlation that when used could cause instability and inaccurate results for the model [25].

Both methods were applied to test their performance, except for the wrapper [29] and embedded [10] methods because of their trade-off between high performance and computational cost.

### 3.4. Data Mining

Part of the KDD process where algorithms are applied to find patterns in the data and to translate it into information depending on the objective of the work [11]. In this study three learning algorithms were applied to predict if a flight will be delayed or not, hence being a binary classification. The selection of these algorithms was made regarding the best algorithms identified in work presented at [15,17,18]. Those algorithms are Random Forest and Multilayer

Perceptron Neural Networks. For comparison purposes, we included a simpler algorithm that by introducing simplicity, understandability, and the ability in handling all types of variables could present an interesting performance (Decision Trees).

Decision Trees (DT) are known as a "collection of decision nodes, connected by branches, extending downward from the root node until terminating in leaf nodes" [25]. Have a divide and conquer philosophy and advantages like the easy interpretation and ability for processing a sort of data types [10,30]. The algorithm used was the C4.5 that came to handle the shortcomings of the ID3 algorithm [30,31].

Random Forests (RF) are well-known of ensemble methods (an ensemble of many individual decision trees) where the prediction is made through voting, i.e. each classifier has equal weight, and the winner is who has the majority of votes or weighted where each base classifier has different voting power [10,32–34].

Multilayer Perceptron Neural Networks (MLP) are an algorithm based on the human brain, where input signals (information) are received by the nodes (neurons). An additional input included as a bias is connected via interconnection weights and then processed through an activation function converting it into an output [35,36]. MLP neural networks use the backpropagation training algorithm which is characterized by its decreasing gradient through the networks that made them capable of minimizing the squared error between the network output and target value for the outputs [30].

### 3.5. Evaluation

After the application of all algorithms, it is necessary to evaluate their performance in a given dataset [33]. For that reason, a wide variety of measures can help in the choice of the "best" algorithm. The measures used in this paper are the Area Under the ROC Curve [10,37], the F-measure [10] and the Accuracy [37] to reverse the effect of unbalanced classes if there is still.

## 4. Results and Discussion

As already mentioned, to improve the prediction performance, different factors such as the used sampling technique (SMOTE and Undersampling), the treatment performed to outliers (inclusion and exclusion of it), the variable selection method (GainRatio for 10, 15 and 20 variables and Correlation-based Feature Selection) and the inclusion and exclusion of some extra variables (delay in departure and real departure time) were implemented. For each of these cases, three learning algorithms were applied: DT, RF, MLP.

As expected, it is possible to observe in Table 1 that most of the classifiers performed better on training data than on unseen data except for the MLP algorithm when the variables "Dep_Delay" and "Real_Dep_Time" are excluded. Probably explained by the possibility of the modeling still unbalanced despite the attempt to eliminate this problem.

If our interest is to predict if a flight is delayed at the destination after taking off, then the two variables have to be taken into consideration. For this purpose, the most suitable algorithm is the MLP with the SMOTE technique, excluding the outliers and with the selection of variables through GainRatio approach selecting ten variables.

Otherwise, if our interest is to predict if a flight will be late at the arrival beforehand, those two variables cannot be included. In that scenario, the algorithm presenting better results, although with lower ROC Curve value, is also the MLP using SMOTE technique, but including the outliers and selecting the attributes through Gain Ratio approach, retaining fifteen variables.

Furthermore, it is visible that the decision tree has a smaller performance in the area under the ROC Curve value than the RF and MLP when considering the two variables. When not considering those two variables the DT has a higher area under the ROC curve compared to RF but still a lower value compared to MLP.

By comparing the results of this study with the previous studies, we can distinguish their approaches and results and this work as seen in Table 2 below. Nonetheless, it is important to mention that these studies differ from the one presented in this paper by using different variables and sources, periods of time and types of approaches. For comparison, the results used are the ones that do not account for the use of the two variables because the other studies also do not consider them. In the case of the MLP, better results are accomplished by other authors [17]. Regarding the RF algorithm, our study attained better results in comparison to the other two studies [15,18].

Additionally, it is possible to compare our best algorithm (MLP) result with some websites prediction performances. As we can see in Table 2 KnowDelay website predictions has the highest accuracy followed by FlightCaster and DelayCast.

Table 1. Results of best approaches for each algorithm. ROC represents the area under the ROC curve, F refers to the F-Measure and ACC is the accuracy.

| | Including Dep_Delay and Real_Dep_Time | | | | | | Excluding Dep_Delay and Real_Dep_Time | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Training Set | | | Test Set | | | Training Set | | | Test Set | | |
| | ROC | F | ACC | ROC | F | ACC | ROC | F | ACC | ROC | F | ACC |
| DT | 0.95 | 0.88 | 94.62 | 0.84 | 0.82 | 79.77 | 0.81 | 0.86 | 87.48 | 0.53 | 0.80 | 83.06 |
| | With Outliers; SMOTE; GainRatio (10) | | | | | | With Outliers; SMOTE; GainRatio (10) | | | | | |
| RF | 1 | 0.99 | 99.14 | 0.85 | 0.82 | 79.62 | 1 | 1 | 100 | 0.52 | 0.79 | 84.26 |
| | With Outliers; Undersampling; GainRatio (10) | | | | | | Without Outliers; SMOTE; GainRatio (10) | | | | | |
| **MLP** | **0.88** | **0.93** | **94.01** | **0.89** | **0.86** | **84.06** | **0.73** | **0.75** | **78.06** | **0.56** | **0.79** | **85.63** |
| | **Without Outliers; SMOTE; GainRatio (10)** | | | | | | **With Outliers; SMOTE; GainRatio (15)** | | | | | |

Table 2. Related Work Comparison.

| Article/ Website | Information used | Objective | Type of Prediction | Algorithms Accuracy (%) | |
| --- | --- | --- | --- | --- | --- |
| | | | | RF | MLP |
| [17] | Delay status of the day; Flight data; Weather data | Predict Class of Delay of an individual flight | Multiclass | - | 87.42% |
| [15] | Flight data; Weather data | Predict arrival delays of individual flight due to weather conditions | Binary | 74.20% | - |
| [18] | Flight data; Weather data | Predict arrival delays of individual flight | Binary | 80.36% | - |
| FlightCaster [38] | Flight; Weather | - | - | 85% | |
| KnowDelay [39] | Airport Performance; Weather | - | - | 90% | |
| DelayCast [40] | Flight; Weather; Airline Company History; Number of Passengers | - | - | 80-90% | |
| Our Work | Flight data; Weather data; Airplane info; Delay Propagation information | Predict arrival delays of individual flight | Binary | 84.26% | 85.63% |

## 5. Conclusions and Future Work

After presenting this work, we can conclude that the better model for prediction is the MLP. It is also important to mention that the SMOTE technique presents better results than the Undersampling technique.

Contrarily to what was expected variables such as the month, scheduled and real timetables and distance, among others, turn out not to be so important. The variables that most contribute to the existence of delay for the best model was mainly the variable of delay in departure when included due to accessing this type of information improves the models' performance. Also, the variables of airplane and congestion denoted some importance in both models nonetheless, lower than the previous one. Variables of weather, in general, were also present in both models but with lower power. However, it is important to mention that all three phases of the weather variables observations turn out to be present in the best models.

For future work, we recommend the implementation of a similar study for the Portuguese case. The development of this study using another type of tool than the one used here (Weka) is also recommended to overcome the difficulty

of not reaching some results due to computational cost and time. We purpose tools such as Apache Hadoop or Python that can handle a large amount of data.

We also recommend the application of other types of attribute selection, parameters tuning, and algorithms not considered that could improve the models results. As suggestion, the extent of this classification problem to a regression one would also be interesting. Finally, as reference for future works, a prediction model that can overcome all the lacks here existent and fulfill the needs not addressed in this study is proposed.

## References

[1]      Belobaba P, Odoni A, Barnhart C. The global airline industry. 2009.

[2]      International Air Transport Association (IATA). Air Passenger Market Analysis. 2015.

[3]      U.S. Department of Transportation, Federal Aviation Administration, Corporation M. Airport Capacity Benchmark Report 2004. Washington, D.C.: 2004.

[4]      Jarrah AIZ, Yu G, Krishnamurthy N, Rakshit A. A Decision Support Framework for Airline Flight Cancellations and Delays. Transp Sci 1993;27:266–80. doi:10.1287/trsc.27.3.266.

[5]      Bureau of Transportation Statistics. Airline On-Time Performance and Causes of Flight Delays 2016. http://www.rita.dot.gov/bts/help_with_data/aviation/index.html#q8 (accessed October 31, 2016).

[6]      Ball M, Barnhart C, Dresner M, Neels K, Odoni A, Peterson E, et al. Total Delay Impact Study -- A comprehensive assessment of the cost and impacts of flight delay in the United States. 2010.

[7]      Tu Y, Ball MO, Jank WS. Estimating Flight Departure Delay Distributions—A Statistical Approach With Long-Term Trend and Short-Term Pattern. J Am Stat Assoc 2008;103:112–25. doi:10.1198/016214507000000257.

[8]      Airports Council International. Preliminary World Airport Traffic Ranking. 2016.

[9]      Bureau of Transportation Statistics. On-Time Performance - Flight Delays at a Glance 2017. https://www.transtats.bts.gov/HomeDrillChart.asp?URL_SelectMonth=4&URL_SelectYear=2017 (accessed July 17, 2017).

[10]     Witten IH, Frank E, Hall MA. Data mining: Practical Machine Learning Tools and Techniques, Third Edition. 3rd ed. Elsevier Inc.; 2011. doi:10.1002/1521-3773(20010316)40:6<9823::AID-ANIE9823>3.3.CO;2-C.

[11]     Fayyad U, Piatetsky-Shapiro G, Smyth P. From Data Mining to Knowledge Discovery in Databases. AI Mag 1996;17:37. doi:10.1609/aimag.v17i3.1230.

[12]     Han J, Kamber M. Data Mining: Concepts and Techniques. vol. 12. 2nd ed. Morgan Kaufmann Publishers, Inc.; 2011. doi:10.1007/978-3-642-19721-5.

[13]     Kononenko I, Kukar M. The Name of the Game. Mach. Learn. Data Min. Introd. to Princ. Algorithms, Horwood Publishing Limited; 2007, p. 2–4.

[14]     Wang XWX. Intelligent Quality Management Using Knowledge Discovery in Databases. 2009 Int Conf Comput Intell Softw Eng 2009:1–4. doi:10.1109/CISE.2009.5364999.

[15]     Belcastro L, Marozzo F, Talia D, Trunfio P. Using Scalable Data Mining for Predicting Flight Delays. ACM Trans Intell Syst Technol 2016;8:1–20. doi:10.1145/2888402.

[16]     GE Aviation. GE Flight Quest: Think you can change the future of flight? 2012. https://www.kaggle.com/c/flight (accessed July 20, 2017).

[17]     Kim YJ, Choi S, Briceno S, Mavris D. A Deep Learning Approach to Flight Delay Prediction. AIAA/IEEE Digit Avion Syst Conf - Proc 2016;December:1–6. doi:10.1109/DASC.2016.7778092.

[18]     Choi S, Kim YJ, Briceno S, Mavris D. Prediction of Weather-induced Airline Delays Based on Machine Learning Algorithms. AIAA/IEEE Digit Avion Syst Conf - Proc 2016;2016–Decem:1–6. doi:10.1109/DASC.2016.7777956.

[19]     Rebollo JJ, Balakrishnan H. Characterization and prediction of air traffic delays. Transp Res Part C Emerg Technol 2014;44:231–41. doi:10.1016/j.trc.2014.04.007.

[20]     Chawla N V., Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: Synthetic minority over-sampling technique. J Artif Intell Res 2002;16:321–57. doi:10.1613/jair.953.

[21]     Hoens TR, Chawla N V. Imbalanced Datasets: From Sampling to Classifiers. In: He H, Ma Y, editors. Imbalanced Learn. Found. Algorithms, Appl. First Ed. 1st Editio, 2013, p. 43–59. doi:10.1002/9781118646106.ch3.

[22]     Weiss GM. Mining with Rarity: A Unifying Framework. SIGKDD Explor 2004;6:7–19. doi:10.1145/1007730.1007734.

[23]     Liu X-Y, Wu J, Zhou Z-H. Exploratory Undersampling for Class Imbalance Learning. IEEE Trans Syst Man Cybern 2009;39:539–50. doi:10.1109/TSMCB.2008.2007853.

[24]     Allison PD. Missing Data. Quant. Appl. Soc. Sci., 2001, p. 72–89. doi:10.1136/bmj.38977.682025.2C.

[25]     Larose DT. Discovering Knowledge in Data: An Introduction to Data Mining. John Wiley & Sons, Inc.; 2005.

[26]     Hawkins DM. Introduction. Identif. Outliers, Springer, Dordrecht; 1980, p. 1–12. doi:10.1007/978-1-4614-0406-4.

[27]     Pyle D. Data Preparation for Data Mining. vol. 17. Morgan Kaufmann Publishers, Inc.; 1999. doi:10.1080/713827180.

[28]     Hall MA, Holmes G. Benchmarking Attribute Selection Techniques for Discrete Class Data Mining. vol. 15. Hamilton, New Zealand: 2002. doi:10.1109/TKDE.2003.1245283.

[29]     Kohavi R, John GH. Wrappers for feature subset selection. Artif Intell 1997;97:273–324.

[30]     Mitchell TM. Machine learning. McGraw-Hill Science/Engineering/Math; 1997.

[31]     Quinlan JR. C4.5: Programs for Machine Learning. San Mateo, California: Morgan Kaufmann Publishers, Inc.; 1993. doi:10.1007/BF00993309.

[32]     Kumar GR, Kongara VS, Ramachandra DG. An Efficient Ensemble Based Classification Techniques for Medical Diagnosis. Int J Latest Technol Eng Manag Appl Sci 2013;II:5–9.

[33]     Dean J. Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners. New Jersey: John Wiley & Sons, Inc; 2014.

[34]     Wang CW. New Ensemble Machine Learning Method for Classification and Prediction on Gene Expression Data. Annu Int Conf IEEE Eng Med Biol Soc 2006;1:3478–81. doi:10.1109/IEMBS.2006.259893.

[35]     Palit AK, Popovic D. Neural Networks Approach. Comput. Intell. Time Ser. Forecast. Theory Eng. Appl. (Advances Ind. Control. 1st ed., Springer-Verlag London; 2005, p. 372. doi:0.1007/1-84628-184-9.

[36]     Zhang G, Patuwo BE, Hu MY. Forecasting with artificial neural networks: The state of the art. Int J Forecast 1997;14:35–62. doi:10.1016/S0169-2070(97)00044-7.

[37]     Kononenko I, Kukar M. Measures for Performance Evaluation. Mach. Learn. Data Min. Introd. to Princ. Algorithms, Horwood Publishing Limited; 2007, p. 68–81.

[38]     Smartertravel. SMARTERTRAVEL. New Tool Predict Flight Delays 2009. https://www.smartertravel.com/2009/08/19/new-tool-predicts-flight-delays/ (accessed January 30, 2018).

[39]     Knowdelay. KNOWDELAY. HOW IT Work Tak a Look What Makes Our Site Unique 2018. http://www.knowdelay.com/how-it-works.html (accessed January 30, 2018).

[40]     Tourism Review. Tourism Review. DELAYCAST – Predict FLIGHT DELAYS 2008. https://www.tourism-review.com/delaycast-predicting-flight-delays-news856 (accessed January 30, 2018).