# Forecasting Airport Surface Traffic Congestion based on Decision Tree

Zhaoyue Zhang[a,*], An Zhang[a], Cong Sun[a], and Shanmei Li[b]

*[a]School of Aeronautics, Northwestern Polytechnical University, Xi'an, 710072, China*
*[b]College of Air Traffic Management, Civil Aviation University of China, Tianjin, 300300, China*

**Abstract**

To improve the operational efficiency of airport surfaces, this paper studies the air traffic congestion prediction of airport surfaces, demonstrates the limitations of traffic congestion prediction, and proposes a prediction method for airport surface traffic congestion based on decision tree. Firstly, the definition and measurement methods of traffic congestion in airport surfaces are promoted. Then, the key factors affecting traffic congestion are extracted, and a prediction model of traffic congestion is established. Finally, we verify the validity of the model based on actual operation data from Atlanta. The results show that the accuracy of the prediction is 70%.

## 1. Introduction

With the rapid development of the air transport industry, air traffic demand has risen sharply. The traffic supply and demand contradictions are prominent, and air traffic congestion is becoming increasingly serious, especially in the case of airport surface traffic congestion. Traffic congestion often causes huge economic losses, environmental pollution (mainly air pollution and airport noise pollution), and increased controller workload, which seriously affects the safety and efficiency of air traffic operations.

To solve the problem of traffic congestion of airport surfaces, it is necessary not only to further strengthen the infrastructure construction of the air traffic control department, but also to maximize the utilization of the existing surface resources. Therefore, it is particularly important to establish an evaluation system of "airport surface traffic congestion". The airport traffic congestion evaluation system has a great impact on the actual traffic operation of the whole airport: it helps airport traffic managers operate aircraft efficiently and safely. The quantitative relationship between traffic congestion and average delay time, traffic saturation, queue length, and so on can be analyzed comprehensively. At the same time, it also helps the relevant aircraft release departments and controllers make correct and wise release decisions.

At present, the problem of air traffic congestion has been studied by many scholars. In 2009, Tao et al. analyzed the relationship between demand, capacity, and flight delays of airport surface traffic and divided the air traffic states by studying the delays of all departing aircraft [1]. In 2012, Wang Lei of Xi'an University of Technology analyzed and developed a short-term prediction model of air traffic flow based on a combination of linear regression and support vector machine. A short-term prediction system, which provides data support and decision-making basis for the relevant control departments to solve airport surface traffic congestion, was also developed [2]. In 2013, Li Shanmei of the Civil Aviation University of China carried out further research by analyzing congestion behavior. The congestion indicator system and air traffic congestion spread prediction were established [3]. The enhanced air traffic management system (Enhanced Traffic Management System, ETMS) in the United States holds that air traffic congestion occurs when air traffic demand is greater than air traffic capacity. Air traffic states are recognized based on the above rule [4]. In 2001, Chatter and Sridhar obtained that there is a nonlinear relationship between the dynamic density index and the controller's workload, and the neural

---

\* Corresponding author.
*E-mail address*: zy_zhang@cauc.edu.cn

network method was used to model the dynamic density [5]. In 2002, Wang and Tene et al. studied the queuing delay time of aircraft arriving and departing the airport and regarded it as a measure of airport congestion. It was pointed out that the different queuing delay time of aircraft was caused by different airport traffic demands and airport traffic capacity, and the propagation phenomenon of queuing delay time was studied [6]. In 2005, the MRTIE organization, represented by Wanke, also compared airport demand with capacity to determine airport surface traffic states [7].

Although there are some advantages of the above studies, the definition and measurement method of airport surface traffic congestion are insufficient. The definition and measurement method of traffic congestion is the basis of air traffic congestion identification and prediction. At present, there is no unified standard for the definition of traffic congestion. In the past, research on air traffic congestion mainly focused on comparing the air traffic demand and capacity. When the demand is greater than the capacity, air traffic congestion occurs. For the complex air traffic system, this definition method is obviously too simple, lacks relevant theoretical analysis and specific quantitative analysis indicators, and does not reflect the essential characteristics of traffic congestion. Moreover, it cannot reflect the dynamic process of the emergence and development of traffic congestion.

Thus, this paper explores the influencing factors of airport surface traffic congestion, establishes the airport surface traffic congestion state prediction model based on decision tree algorithm, and finally gives an example analysis. A decision tree of air traffic congestion for the Atlanta airport is established. The validity of the prediction model is verified.

## 2. Analysis of Traffic Congestion in Airport Surface

### 2.1. Definition of Airport Surface Traffic Congestion and Its Influencing Factors

Based on previous studies, we promote the definition of airport traffic congestion as an imbalance between airport traffic demand and capacity. Demand is the number of aircraft expected to pass through a taxiway or runway in a given time or space. Capacity refers to the maximum number of aircraft that can be accommodated on a taxiway or runway in a specific time or space. When the aircraft is running smoothly, the traffic demand of the airport is often less than the capacity. When the traffic demand of the airport is greater than its capacity, air traffic congestion of the airport surface occurs. Thus, the traffic congestion at the airport surface is determined by the traffic capacity of the airport and the traffic demand of the airport. When some factors change the air traffic capacity and demand of airport surface, the air traffic state of the airport surface will be changed. These factors are usually divided into human factors and environmental factors [8-9], including time period, special circumstances, airport infrastructure quality, holidays, controllers' workload, weather conditions, and traffic volume.

### 2.2. Measurement of Air Traffic Congestion at Airport Surface

Measures of airport surface traffic congestion status include taxi time, taxi speed, and taxi delay. Here, we use taxi time to measure the congestion state and define the traffic congestion index ( $CI$ ) as the ratio of the aircraft taxi delay time to the actual average taxi time. The formula is as follows:

$$CI = \frac{t - t_0}{t_0} \tag{1}$$

Where $t$ is the actual taxi time of a certain period and $t_0$ is the average taxi time.

We classify the traffic congestion state of the airport surface and divide it into four states: smooth, slight congestion, moderate congestion, and severe congestion. By consulting with the first line controllers and references, we propose the classification criteria of airport surface traffic congestion state.

The traffic congestion state of the airport surface is classified according to the following rules:

(1) $CI \leq 0.10$, the traffic state is smooth;
(2) $0.10 < CI \leq 0.30$, the traffic state is slight congestion;
(3) $0.30 < CI \leq 0.50$, the traffic state is moderate congestion;
(4) $0.50 < CI$ , the traffic state is severe congestion.

*2.3. Attribute Extraction of Airport Surface Traffic Congestion*

We explore air traffic congestion properties of airport surfaces by accessing relevant information, as shown in Table 1.

Table 1. Crowding attribute selection

| Attribute name | Attribute value |
|---|---|
| Weather conditions (rain, fog, snow, hail···) | 0-No impact, 1-Impact |
| Time slot | 0-Peak, 1-Low peak |
| Holiday and vacations | 0-Y, 1-N |
| Traffic volume | 0-Large, 1-Small |
| Airport infrastructure integrity | 0-Complete, 1-Incomplete |
| Ability of controllers | 0-Outstanding, 1-Average, 2-Poor |
| Aircraft size | 0-Large, 1-Small |
| Air force restrictions | 0-Yes, 1-No |
| Traffic volume in the previous period | 0-Large, 1-Small |
| Working day | 0-Yes, 1-No |

Based on the integrity of the airport infrastructure and the performance of airport taxiway and runway, the air traffic data of the airport surface is not well obtained. The influence of human factors cannot be quantified well. Thus, we extract influencing factors from environmental factors, such as weather conditions, time periods, holidays, traffic volume, and traffic volume of the previous time period. The identification method of air traffic congestion states for airport surfaces is established by these factors.

(1) Weather conditions: According to the weather data we find, the weather conditions are not clearly indicated. There is only weather delay data, so the weather attribute is affected when the weather delay is bigger than zero, and the weather attribute has no effect on air traffic when there is no weather delay.

(2) Time period: according to the statistical data and the working experience of front-line controllers, we define 8:00-11:00 in the morning and 2:30-5:00 in the afternoon as the peak period, and the rest of the time is the low peak period. The air traffic congestion often occurs in the peak periods.

(3) Holidays: These include national holidays.

(4) Traffic volume: According to the experience and the actual situation of the airport surface operation, we count the air traffic demand every 15 minutes. When the demand is equal or greater than to 10, the value of traffic volume attribute is 0. The value of traffic volume attribute is 1 when the demand is less than 10.

(5) Whether working day: Monday to Friday are working days, and Saturday and Sunday are non-working days.

**3. Airport Surface Traffic Congestion Prediction based on Decision Tree Algorithm**

Decision tree [10-12] is a non-parametric supervised learning method. It can summarize decision rules from a series of data with features and labels, and present these rules with the structure of tree graph to solve the problem of classification and regression.

*3.1. Decision Tree C4.5 Algorithm*

The decision tree C4.5 algorithm is an improved version of the decision tree ID3 algorithm. The attribute with the highest information gain rate is adopted as the criterion for selecting the branch attribute while inheriting advantages of the ID3 algorithm. The basic principle of the decision tree C4.5 algorithm is expressed as follows:

$S$ is a collection set of $s$ samples. Suppose there are $M$ classes $(i = 1, 2, \cdots, m)$. The expected information required for the classification of a given sample is shown as follows:

$$I(s_1, s_2, \cdots, s_m) = -\sum_{i=1}^{m} p_i \log_2(p_i) \tag{2}$$

Where $p_i$ is the probability that any sample belongs to $C_i$ and $s_i$ is the number of samples belonging to class $i$. $p_i$ is equal to $s_i/s$.

Let attribute $A$ have $v$ subsets $s_1, \cdots, s_v$. The samples in $s_j$ have a value $a_j$ on $A$. If $A$ is selected as a test attribute, then these subsets correspond to branches that are grown by the node representing the set $S$. Let $s_{ij}$ be the number of samples of class $C_i$ in subset $s_j$. The entropy according to the subset divided by $A$ is calculated by the following formula:

$$E(A) = \sum_{j=1}^{v} \frac{s_{1j} + s_{2j} + \cdots + s_{mj}}{s} I(s_{1j}, s_{2j}, \cdots, s_{mj}) \tag{3}$$

Where $\dfrac{s_{1j} + s_{2j} + \cdots + s_{mj}}{s}$ is the weight of the $j^{\text{th}}$ subset and is equal to the number of samples in the subset (i.e., the value of $A$ is $a_j$) divided by the total number of samples in $s_j$. The smaller the entropy value, the higher the purity of the subset partition. For a given subset $s_j$, there are:

$$I(s_{1j}, s_{2j}, \cdots, s_{mj}) = -\sum_{j=1}^{m} p_{ij} \log_2(p_{ij}) \tag{4}$$

Where $p_{ij} = \dfrac{s_{ij}}{s_j}$ is the probability that the samples in $s_j$ belongs to class $C_i$.

The corresponding information gain value can be obtained from the expected information and the entropy value, and the information gain value obtained by branching A is obtained by the following formula:

$$Gain(A) = I(s_{1j}, s_{2j}, \cdots, s_{mj}) - E(A) \tag{5}$$

The information gain in $Gain(A)$ is the same as that in the ID3 algorithm, and the split information $SplitInfo(S, A)$ represents the breadth and uniformity of splitting sample set $S$ according to attribute $A$.

$$SplitInfo(S, A) = -\sum_{i=1}^{c} \frac{s_i}{s} \log_2 (\frac{s_i}{s}) \tag{6}$$

Thus, the information gain rate for an attribute can be calculated as follows:

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInfo(S, A)} \tag{7}$$

The C4.5 algorithm selects the attribute with the highest information gain ratio as the test attribute of a given set $S$ by calculating the information gain rate of each attribute. Each node is created and the attribute value is marked, and then branches will be created according to the attribute value.

### 3.2. Establishment of Congestion Prediction Model based on Decision Tree Algorithm

The decision tree is constructed based on the data. Before the establishment of the decision tree, we need to collect and count all kinds of required data, attributes, etc., carry on the preliminary processing to obtain the effective data, and then carry on the above process to establish the decision tree.

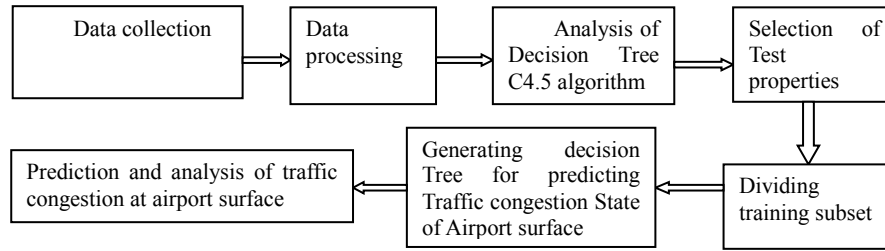To sum up, the process of building a decision tree is shown in Figure 1.

Figure 1. Decision tree algorithm flow

## 4. Example Analysis

Because Atlanta Airport is one of the largest and busiest airports in the world, it is also an academic research object for many scholars and experts in the civil aviation field. Therefore, we also conduct research on the air traffic states at the Atlanta airport surface using the congestion prediction model established above.

### 4.1. Data Analysis and Processing

The air traffic data of the Atlanta International Airport surface from January 1, 2015 to January 7, 2015 is obtained from the US Transportation Administration website. There is a total of 308 time periods in the daytime, and every period is 15 minutes. The traffic congestion index ( $CI$ ) of every time period is calculated according to Equation (1). Then, the air traffic states of the airport surface can be identified based on the value of $CI$ . There are four states: smooth, slight congestion, moderate congestion, and severe congestion. In additional, the attribute value of every time period is calculated. Some of the calculation are shown in Table 2.

Table 2. Partial calculation data

| Sample number | Weather | Time slot | Holiday or not | Traffic | Working day or not | Traffic volume in the previous period | CI | Crowding state |
|---|---|---|---|---|---|---|---|---|
| 1 | no influence | peak | yes | small | yes | small | 0.10 | smooth |
| 2 | no influence | peak | yes | large | yes | small | 0.17 | moderate congestion |
| 3 | no influence | peak | yes | small | yes | large | 0.10 | smooth |
| 4 | influential | peak | yes | large | no | small | 0.15 | moderate congestion |
| 5 | influential | low peak | no | large | yes | small | 0.16 | moderate congestion |

The $CI$ value is calculated to two decimal places. Because the actual taxiing time of some samples is less than the average taxiing time, the $CI$ value is negative, so the $CI$ value is directly represented by "0" when the $CI$ value is negative.

After further data processing statistics, the number of samples of different attribute values in each attribute is obtained, as shown in Table 3.

Table 3. Number of samples for different attribute values

| Smooth | 250 | Slight congestion | 45 | Moderate congestion | 10 | Severe congestion | 3 |
|---|---|---|---|---|---|---|---|
| no influence | 232 | | | influential | 57 | | |
| peak | 154 | | | low peak | 154 | | |
| holiday and vacations | 132 | | | non-holidays | 176 | | |
| traffic volume | 212 | | | small traffic volume | 96 | | |
| working day | 220 | | | non-working days | 88 | | |
| large traffic volume in the previous period | 211 | | | small traffic volume in the previous period | 97 | | |

### 4.2. Construction of Decision Tree

$S$ is the set of all samples, and the number of samples is $s = 308$. They are divided into four different categories, that is, smooth $C_1$, slight congestion $C_2$, moderate congestion $C_3$, and severe congestion $C_4$. Let $s_1$ be the number of samples of class $C_1$, $s_2$ be the number of samples of class $C_2$, $s_3$ be the number of samples of class $C_3$, and $s_4$ be the number of

samples of class $C_4$. From the above table, $s_1 = 250$, $s_2 = 45$, $s_3 = 10$, $s_4 = 3$. The expected information required to classify the samples is obtained as follows:

$$
\begin{aligned}
I(s_1, s_2, s_3, s_4) &= -\sum_{i=1}^{4} p_i \log_2(p_i) \\
&= -\frac{250}{308}\log_2\frac{250}{308} - \frac{45}{308}\log_2\frac{45}{308} - \frac{10}{308}\log_2\frac{10}{308} - \frac{3}{308}\log_2\frac{3}{308} \\
&= 0.875379
\end{aligned}
\tag{8}
$$

Next, the information gain rate of each attribute is calculated separately, and then the attribute with the largest information gain rate is selected as the split attribute of the decision tree.

Suppose that "whether holiday" is used as the split attribute, and the "whether holiday" attribute has two different attribute values: $\{a_1 = \text{yes}, a_2 = \text{no}\}$. The number of samples with $a_1 = \text{yes}$ is 132, and the number of samples with $a_2 = \text{no}$ is 176. $s_{ij}$ is the number of samples belonging to class $C_i$ in the subset $a_j$, e.g. $s_{11}$ indicates the number of samples whose traffic state is smooth in holidays, and $s_{12}$ is the number of samples whose traffic state is smooth in non-holidays.

When the attribute "whether holiday" is "yes", $s_{11} = 103$, $s_{21} = 21$, $s_{31} = 6$, $s_{41} = 2$. According to Equation (4), we can obtain the following:

$$
\begin{aligned}
I(s_{11}, s_{21}, s_{31}, s_{41}) &= -\sum_{i=1}^{4} p_{i1} \log_2(p_{i1}) \\
&= -\frac{103}{132}\log_2\frac{103}{132} - \frac{21}{132}\log_2\frac{21}{132} - \frac{6}{132}\log_2\frac{6}{132} - \frac{2}{132}\log_2\frac{2}{132} \\
&= 0.995469
\end{aligned}
\tag{9}
$$

Similarly, when "whether holiday" is "no", $s_{12} = 147$, $s_{22} = 24$, $s_{32} = 4$, $s_{42} = 1$, and $I(s_{12}, s_{22}, s_{32}, s_{42}) = 0.775392$. The entropy of every subset divided by the "whether holiday" attribute can be calculated as follows:

$$
\begin{aligned}
E(\text{Whether it is holiday}) &= \sum_{j=1}^{2} \frac{s_{1j} + s_{2j} + s_{3j} + s_{4j}}{s} I(s_{1j}, s_{2j}, s_{3j}, s_{4j}) \\
&= \frac{132}{308} I(s_{11}, s_{21}, s_{31}, s_{41}) + \frac{176}{308} I(s_{12}, s_{22}, s_{32}, s_{42}) \\
&= 0.869711
\end{aligned}
\tag{10}
$$

According to Equation (5), the information gain is:

$$
\begin{aligned}
Gain(\text{Whether it is holiday}) &= I(s_1, s_2, s_3, s_4) - E(\text{Whether it is holiday}) \\
&= 0.005668
\end{aligned}
\tag{11}
$$

Based on Equation 6), we can obtain the following:

$$
\begin{aligned}
SplitInfo(\text{Whether it is holiday}) &= \sum_{j=1}^{2} p_j \log(p_j) \\
&= -\frac{132}{308}\log_2\frac{132}{308} - \frac{176}{308}\log_2\frac{176}{308} \\
&= 0.985228
\end{aligned}
\tag{12}
$$

According to Equation (7), the information gain rate is calculated as follows:

$$GainRatio(\text{Whether it is holiday}) = \frac{Gain(\text{Whether it is holiday})}{SplitI(\text{Whether it is holiday})} = \frac{0.005668}{0.985228} = 0.005753 \qquad (13)$$

In the same way, the information gain rate of other attributes can be calculated:

$$GainRatio(\text{Time period}) = 0.079358$$

$$GainRatio(\text{Time volume}) = 0.032634$$

$$GainRatio(\text{Working day}) = 0.047388$$

$$GainRatio(\text{Traffic volume in the previous period}) = 0.022269$$

$$GainRatio(\text{Weather}) = 0.016909$$

Then, we sort the information gain rates as follows:

$GainRatio$(Time period)
$> GainRatio$(Working day)$GainRatio$(Time volume)
$> GainRatio$(Working day)$GainRatio$(Traffic volume in the previous period)
$> GainRatio$(Weather)
$> GainRatio$(Whether it is a holiday)

The decision tree of air traffic congestion prediction for the Atlanta International Airport surface is shown in Figure 2. It can be seen that the first split attribute is "time period", the second split attribute is "traffic volume", and so on.



Figure 2. Airport surface traffic congestion state prediction decision tree

Where:

$x_1$ —the traffic volume during the previous time period;
$x_2$ —the time period;
$x_3$ —the traffic volume;
$x_4$ —the weather;

$x_5$—whether it is a working day;

$x_6$—whether it is a holiday;

$y = 1$—the traffic state is smooth;

$y = 2$—the traffic state is slight congestion;

$y = 3$—the traffic state is moderate congestion;

$y = 4$—the traffic state is severe congestion.

According to the decision tree, the relationship between each attribute and the final traffic congestion state can be clearly obtained, for example:

$$x_2 - x_5 - x_4 - x_3 - x_6 - x_1 - y = 1$$

When the time period = peak, whether work day = yes, weather conditions = no impact, traffic volume = large, whether holidays = yes, and traffic volume in the previous time period = large, the traffic state of the airport surface is smooth. Thus, if the air traffic managers get the values of all the attributes, they can predict the air traffic state of the next time period based on the decision tree.

### 4.3. Verification of Decision Tree

We randomly select another ten sets of data samples to verify the established decision tree. The data samples and verification results are shown in Table 4.

Table 4. Verification results

| Sample number | Time | Time slot | Working day or not | Weather | Traffic | Holiday or not | Traffic volume in the previous period | Crowding state | Predictive state |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 2014/12/21 | peak | no | influential | large | no | small | severe congestion | smooth |
| 2 | 2014/12/21 | low peak | yes | no influence | small | no | large | smooth | smooth |
| 3 | 2014/12/22 | peak | yes | no influence | large | no | large | smooth | smooth |
| 4 | 2014/12/23 | peak | yes | influential | small | no | large | slight congestion | slight congestion |
| 5 | 2014/12/23 | low peak | yes | influential | small | no | small | smooth | smooth |
| 6 | 2014/12/24 | low peak | yes | influential | small | yes | large | smooth | smooth |
| 7 | 2014/12/25 | peak | yes | no influence | large | yes | large | slight congestion | smooth |
| 8 | 2014/12/25 | low peak | yes | no influence | large | yes | small | smooth | smooth |
| 9 | 2014/12/26 | low peak | yes | no influence | large | no | large | smooth | smooth |
| 10 | 2014/12/27 | peak | no | no influence | small | no | large | smooth | moderate congestion |

It can be obtained from the table that in the ten samples of the verification data, there are seven groups in which the predicted congestion state is the same as the actual traffic congestion state, so it can be concluded that the accuracy rate of the predicted traffic state prediction decision tree of the airport surface is 70%. Through verification, although the constructed decision tree has certain errors in the prediction of traffic congestion state, it is still feasible. In actual work, traffic conditions can also be predicted based on our decision model. When the traffic managers become aware that the future traffic state is severe congestion through our method, they can take appropriate measures in advance to relieve the degree of air traffic congestion and avoid the occurrence of severe congestion.

## 5. Conclusion

It is an important measure for the controller to ensure the safety of the aircraft, reduce the workload, and improve the operational efficiency by analyzing and predicting the traffic congestion of the airport surface and preparing the solution in advance.

This paper refers to the research on traffic congestion of airport surfaces at home and abroad, and it analyzes the influencing factors of traffic congestion in various airport surfaces. Some extreme factors and unquantifiable factors are excluded due to the limitation of conditions. The decision tree algorithm is selected from many knowledge expression methods to establish a decision tree of traffic state prediction. The decision tree is verified based on the actual operational data of the Atlanta airport surface. We establish a congestion state prediction decision tree for the Atlanta airport, and the prediction accuracy is 70%.

## Acknowledgements

## References

1. T. Xu, J. Ding, B. Gu, and J. Wang, "Flight Delay Warning based on Incremental Array Support Vector Machines," *Acta Aeronautica ET Astronautica Sinica*, Vol. 30, No. 7, pp. 1256-1262, 2009
2. B. M. He, "Study on Short-Term Prediction Model of Air Traffic Flow," *Journal of Wuhan University of Technology* (*Transportation Science and Engineering*), Vol. 12, No. 1, pp. 334-356, 2012
3. S. M. Li, "Research on Identification and Prediction Methods of Air Traffic Congestion," Tianjin University Press, 2013
4. Volpe National Transportation Systems Center, "Enhanced Traffic Management System (ETMS) Functional Description," U.S. Dept. of Transportation, Cambridge, MA, 2002
5. G. B. Chatterji and B. Sridhar, "Measures for Air Traffic Controller Workload Prediction," in *Proceedings of 1st AIAA Aircraft Technology*, *Integration and Operations Forum*, pp. 104-125, 2001
6. P. T. R. Wang, N. Tene, and L. Wojcik, "Relationship Between Airport Congestion and at-Gate Delay," in *Proceedings of 21st Digital Avionics Systems Conference*, pp. 67-88, 2002
7. C. R. Wanke, L. Song, S. Zobell, D. Greenbaum, and S. Mulgund, "Probabilistic Congestion Management," 6th USA/Europe Seminar of Air Traffic Management R&D, pp. 27-30, 2005
8. Z. Zhao, "Research on Airspace Capacity Assessment and Forecast," Nanjing University of Aeronautics and Astronautics Press, pp. 34-41, 2015
9. X. N. Dong, "Sector Capacity Evaluation and Complexity Analysis," Nanjing University of Aeronautics and Astronautics Press, pp. 65-80, 2017
10. Y. X. Sun, C. F. Shao, D. Zhao, and S. Ou, "Traffic Accident Severity Prediction Model based on C5.0 Decision Tree," *Journal of Chang'an University* (*Natural Science Edition*), Vol. 34, No. 5, pp. 123-132, 2014
11. X. W. Wang, C. Q. Yuan, and M. Huang, "A Motion Prediction Mechanism based on Fuzzy Decision Tree," *Computer Science*, Vol. 32, No. 9, pp. 1176-179, 2005
12. R. Li, Y. Liu, J. H. Li, X. P. Gu, D. X. Niu, and Y. Q. Liu, "Study on Daily Characteristic Load Prediction based on Improved Decision Tree Algorithm," *Proceedings of the CSEE*, Vol. 25, No. 23, pp. 36-41, 2005