

Study on Anomaly Detection Algorithm of QAR Data Based on Attribute Support of Rough Set

Hui Yang¹, Chunjing Xiao² and Yongwei Qiao³

^{1,2}College of Computer Science and Technology, Civil Aviation University of china,
Tianjin, china, 300300

³Engineering and Technology Training Center, Civil Aviation University of china,
Tianjin, china, 300300

¹yanghui_z@sina.com, ²chunjingxiao@163.com, ³yongweiqiao@163.com

Abstract

According to the characteristics of the large amount of QAR data, such as many parameters, time constraints, strong randomness and the problems of discrete data, together with attribute reduction and rules collecting during QAR anomaly detection, the paper proposed a anomaly detection algorithm of QAR data based on attribute support of rough set. Firstly, we discrete QAR data after converting the time sequence data into non-time sequence information system based on attribute support. Second, we carry on the attribute reduction using difference matrix based on statistics. At last, we get the fault rules with the decision tables based on fault occurrence statistics of the decision attributes. The method not only primely retains the time characteristics of QAR data, but also strengthens the relation between condition attributes and decision attributes. The efficiency of attribute reduction and rules extraction is high and the amount of calculation is small. The experimental results show that the method is feasible, reliable and availability.

Keywords: Regression, data structure, prediction, simulation

1. Introduction

QAR (Quick Access Recorder) is a kind of quick access recorder, recording a large amount of flight time sequence data. QAR monitoring is an important basis for flight inspection, safety assessment, safety incident investigation and aircraft maintenance. Anomaly detection is an important aspect of data mining. It can judge whether there is hidden fault or fault trend, provide accurate fault information for maintenance, avoid aircraft accidents and achieve the fault early warning purposes through the detection of abnormal QAR data [1-5].

Discretion, attribute reduction and rules acquisition is the main process of anomaly detection based on rough set. QAR data has the characteristics of large amount of data, such as many parameters, time constraints, strong interference, strong randomness and uncertainty, so that its data become single and lost some hidden information, and the contacts between condition attributes and decision attributes is poor after using the traditional discrete interpolation, while attribute reduction using the difference matrix and the general QAR rules acquisition strategy will cause a large amount of calculation and extraction rules, which is difficult for us to give valuable decision support. To solve these problems, this paper puts forward QAR abnormal data detection algorithm based on the important support degree of attributes of rough set. The method converts QAR data into non-time sequence information system, discretizes the data based on attribute support, conducts attribute reduction using difference matrix based on statistics and gets the fault rules using the decision tables based on fault occurrence statistics of the decision attributes. The method not only primely retains the

time characteristics of QAR data, but also strengthens the relation between condition attributes and decision attributes. The efficiency of attribute reduction and rules extraction is high and the amount of calculation is small. The experimental results show that the method is feasible, reliable and availability.

2. Related Work

The algorithms of anomaly detection are based on statistical, depth, density, distance and high dimensional data. The effect is not ideal for high dimensional data mainly because high dimensional data attributes are often interrelated. The QAR data is composed of discrete, digital and mixed discrete quantity, which has the strong relations between attributes and the strong implicitness of fault. Anomaly detection of QAR data belongs to the problem of multiple attribute decision. Rough set theory is a tool to deal with incomplete and uncertain problems, which can effectively deal with all kinds of incomplete information according to our understanding of the obtained data, find the connotative knowledge and reveal the potential rules [6]. Study on the combination of rough set theory and multi-attribute decision in a gradually mature stage [7-8]. AlamS.S (2002) proposed the AHP method based on the Rough method and directly improved on the traditional multiple attribute decision method [9]. A.IDimitras (1999) predicted Greece's company mergers and acquisitions by the method of Rough and achieved good results [10]. Meng Zuqiang and Cai Zixing (2004) designed the personalized decision rules mining algorithm using rough set theory [11]. The literature [12] proposed the real-time temporal logic framework, making the event variable to represent temporal sequence. The method in literature [13] detected timing mode using dynamic programming, but not suitable for depth discussion on the realization of algorithm and the processing of timing. The 14th Literature provided the concept of timing information systems (TIS) and real-time timing information systems (RTTIS), two ways of converting the time sequence into information system, made it possible to deal with the time sequence using the rough set, but they didn't elaborate the obtaining problem of timing information system [14].

3. TIS Converted to IS

Timing information system (TIS) reflects good timing characteristics and clearly sees the relationship between the timing and the non-sequential attributes through decision table. The object processed by rough set is non-sequential information system, thus QAR flight data must be converted to non-timing timing information systems (IS). It can eliminate the timing properties and implicitness in the decision-making table time characteristics by increasing the number of attributes. The following step is the conversion algorithm.

Input: Timing Information System (TIS)

Output: traditional information systems (IS)

Step 1: The customer determines the number of back time slices Δ

Step 2: Create $\Delta * |C|$ new attributes, corresponding to a time slice of the historical data

Step 3: The current time slice data combines with historical data to form a new object data

Step 4: Keep the decision attribute values of current time slice and ignore the former Δ decision attribute.

4. QAR Data Attribute Discrete based on the Degree of Support

The requirements of decision attributes in decision table through QAR data mining based on rough set are highly and largely affecting the generation of decision table rules. To

establish a valuable decision table, the relationship between the conditional attributes and the decision attributes must be established. This paper proposes the improved discrete method by combining QAR data domain knowledge with the attribute support degree:

Definition 1: $S = (U, A, V, f)$ is defined as a decision table, $A = C \cup D, C \cap D = \Phi, C$ is the condition attribute set, D is the decision attribute set. If $U / C = \{X_1, X_2, \dots, X_n\}$, $U / D = \{Y_1, Y_2, \dots, Y_m\}$, Then the support degree of the decision attribute D on condition attribute C (also known as C on D) is defined by Formula(1)

$$K_c(D) = \frac{1}{|U|} \sum_{i=1}^m |C_{-Y_i}| = \frac{1}{|U|} \sum_{i=1}^m |pos_c(Y)_i|, Y_i \in U / D \quad (1)$$

Where $|@|$ represents the number of elements which are contained in the collection. Support decision attribute in decision table is a measure of the overall classification capability. This can be the feedback information, measuring the quality of discrete data.

The algorithm steps based on attribute support degree are as follows:

Step 1 Select QAR data attribute set which is related to the fault, and these properties $a \in C$, then calculate the range $\min(V_a), \max(V_a)$ of range V_a , and $d = \frac{\max(V_a) - \min(V_a)}{k_a + 1}$, K_a is the maximum number of points, k_a is the number of

dynamic segmentation points, the segmentation point set of a is defined by Formula(2)

$$C_a^{k_a} = \begin{cases} \emptyset & k_a = 0 \\ \{\min(V_a) + id, i = 1, \dots, k_a\}, & k_a \geq 1 \end{cases} \quad (2)$$

The initial segmentation point set space is $\Omega_a = \{C_a^{k_a} : 0 \leq k_a \leq K_a\}$. The absolute error degree between the new decision attribute support and the original decision attribute support is β . Given the initial value β and $k_a = K_a$.

Step 2: C is the set of containing the m - k original discrete attribute. If $C \neq \Phi$, calculate the original decision attribute support degree $K_c(D)$ by formula (1). Otherwise, calculate it after respectively divide U by segmentation method for each attribute $a_i (i = 1, \dots, m)$, and take the highest value as the original decision attribute support. It is meaning to take $K_c(D) = K_{\{a_j\}}(D) = \max_{1 \leq i \leq m} K_{\{a_i\}}(D)$ and set $C = C \cup \{a_j\}$.

Step 3: Repeat for $p = 1, 2, \dots, k$

① Initialize segmentation points k_p of attributes $a_{i_p}, k_p = k_a$. For continuous attributes a_{i_p} and the initial segmentation points k_p , we can have a division of property k_{i_p} on the U by interval segmentation.

② set $C' = C \cup \{a_{i_p}\}$, then calculate the decision attribute support degree $K_{c'}(D)$ by formula(1).

③ Determine whether $|K_{c'}(D) - K_c(D)| \leq \beta$ is established. If it is established, then make $p = p + 1$ and turn ①, otherwise turn ④.

④make $k_p = k_p - 1$, for the new split points k_p , divide U by interval segmentation method and recalculate $K_{c.}(D)$, then turn ③.

Step 4: Encoded the discrete attribute values by 0,1,2, ...etc .

The attribute-based support discrete method maintains relationships between condition attributes and decision attributes, and it also ensures the validity of the decision table classification well. But because of the complexity and large number of QAR data, it does different treatment of numerical quantity, discrete data and hybrid amount.

Numeric quantity: for each of the 4 side frame consists of a frame and the average value as the frame value to solve the numerical Mauritius point, ensure the reliability of anomaly data mining, reduce the amount of data and save the computational time.

Discrete data: QAR data has fewer states and be relatively simple and the majority of it is Boolean type. In order to maintain consistency with the numerical amount, each is composed of four sub-frame one and makes OR operation between data to ensure the abnormal data state.

Discrete mixed amount: compared to discrete data, this kind of state is more complex and select a status which has most number of data blocks near the point as the point state. If the number is same, we select the nearest point state value as the point state value.

In addition, the discrete method based on attribute support is an optimization process, maximums the discrete degree by adjusting the size of k_a and keeps the original attribute support degree in the range of allowable error. In the calculation process, k_a can be continuously adjusted and it can be effective by averaging multiple computing values.

5. QAR Attribute Reduction of Decision Table based on Attribute Frequency with Time Characteristic

If the information table is very large, the calculation of attribute reduction based on difference matrix is huge. In order to reduce the computational complexity, we count the attribute frequency of difference matrix and give QAR data attribute reduction method, keeping the time characteristic of QAR data based on the principle of locality and the least used elimination algorithm recently. The main process is as follow:

- (1) For information table $S = \langle U, A \rangle$, set $\Delta t = 0$
- (2) Calculate the difference matrix M of decision table S and assign the core attributes to the reduction attribute set **Reduct**, namely **Reduct** = S_{ij} .
- (3) Calculate the attribute frequency in Δt and choose the highest frequency denoted by a_i , namely **{Reduct = Reduct \cup a_i }** and delete items in the identification matrix which contain conditional attribute a_i .
- (4) Determine whether $M = \emptyset$ is established.
- (5) If established, the output is attribute reduction set **Reduct**. Otherwise, if the attributes in Δt have been selected, we can choose attributes in $\Delta t = \Delta t + 1$ and turn to (3).

6. Rules Acquisition Strategy based on Fault Statistics of Decision Attribute

The general rule algorithms are as far as possible to obtain a minimal set of rules. The rules must be a lot because of the large QAR data quantity. According to the characteristics of QAR data, it was discretized into normal and abnormal states. The most important thing is to dig out the important degree for fault occurred. It simplifies the extraction of rules for the starting point of aircraft safety requirements greatly.

The following is new rule acquisition steps:

Input: decision-made table after attribute reduction $S = \langle U, A - reduct \rangle$

Output: the rules set consisting of the most likely fault occurred rules from all possible rule sets.

Step 1: use the general method to get the whole rules set.

Step 2: each rule is classified and based on whether the fault occurred. Rules which belong to a class of occurred fault are represented as disjunctive or conjunctive normal form.

Step 3: each rule of regular expression is converted into the main disjunctive normal form or principal conjunctive normal form. Then the method counts the major items and the small items.

Step 4: the related attributes of large items or small items have high occurrence frequency as an important attribute. Namely, the degree of importance of these attributes is very high when fault occurs. If the attribute is abnormal, the possibility of fault is very high. Then QAR data is divided into normal and abnormal states based on the normal range of given data.

Step 5: Output the rule set of the highest likelihood of fault occurred.

7. Experiment and Experimental Results Rules

7.1. Experimental Data

This paper selects Airborne QAR data of Boeing 737 aircraft from August 2008 to August 2009. The operating system is Windows XP and the experimental platform is developed using the C # programming language.

7.2. Data Discrete

Discrete experiment selects several key attribute of flap fault, such as trailing edge flap position(TEFlapPos), leading edge flap position(L.E.FlapPos), flap handle(FlapHandle), stick(CntrlCol), steering wheel(CntrlWhl), pitch(Ptich) and Slope(Roll).To illustrate the comprehensiveness of data processing algorithms on numerical, discrete data and discrete hybrid data, we have four discrete attributes, TE-FlpByPass, FlapExt-1, FlapTrans-1 and LE-Slat-1 added. They are numbered by C1, C2, ... , C11. Where "0" in C8 represents "normal","0"and "1" in C9, C10 respectively represents "FALSE" and "TRUE" ,and it represents "Extend" in C11. Fault data and the normal data are shown in Table 1.

Table 1. Examples Aircraft Flap Fault Diagnosis

U	related attributes of Flap fault											Fault
	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	D
1	0.8	0.7	0.6	0.9	0.7	0	0	0	1	0	0	1
2	0.9	0.7	0.6	0.9	0.6	0	0	0	1	0	0	1
3	0.6	0.6	0.3	0.9	0.3	0.3	0	0	1	0	0	0
4	1	0.3	0.4	0.3	1	0	1	0	1	0	0	0
5	0.8	0.7	0.3	0.9	0.3	0.2	0	0	1	0	0	0
6	1	0.3	0.8	0.3	1	0	1	0	1	0	0	0
7	0.9	0.3	0.9	0	0.9	0	1	0	1	0	0	0
8	0.8	0.2	0.7	0	0.8	0	1	0	1	0	0	0
9	0.6	0.2	0.6	0	0.6	0	0.9	0	1	0	0	0
10	0.1	0.4	0.4	0.3	0.6	0	0	0	1	0	0	0
11	0.9	0.4	0.6	0.4	0.9	0	0.8	0	1	0	0	0
12	0.6	0.4	0.7	0	0.7	0	0.5	0	1	0	0.1	1
13	0.7	0.6	0.7	0	0.9	0	0.5	0	1	0	0.1	0
14	0.9	0.5	0.9	0	0.8	0	0.8	0	1	0	0.2	0
15	0.8	0.5	0.9	0	0.8	0	0.8	0	1	0	0.1	0
16	0	0.5	0	0.1	0.1	0.8	0	0	1	0	0.9	1
17	0	0.8	0	0.1	0.1	0.8	0	0	1	0	0.9	1
18	0	0.5	0	0.1	0.1	0.8	0	0	1	0	0.9	0
19	0	0.8	0	0.1	0.1	1	0	0	1	0	1	1
20	0	0.8	0	0.2	0.2	1	0	0	1	0	1	1

Discrete results are shown in Table 2 when the initial quantization value $\beta = 0.01$ and $k_a = 5$.

Table 2. Discrete Attribute Table based on Improved Method

condition attributes	Quantization values of condition attributes					
	1	2	3	4	5	6
C1	[0,0.5]	(0.5,1]				
C2	[0.2,0.8]					
C3	[0,0.15]	(0.15,0.3]	(0.3,0.45]	(0.45,0.6]	(0.6,0.75]	(0.75,0.9]
C4	[0.1,0.9]					
C5	[0.1,0.55]	(0.55,1]				
C6	[0,0.17]	(0.17,0.34]	(0.34,0.51]	(0.51,0.68]	(0.68,0.85]	(0.85,1]
C7	[0,1]					
C8	0					
C9	1					
C10	0					
C11	0.1	0.2	0.9	1		

Flying height is the numeric attribute, containing the most of obvious time characteristics. From the view of data trends, we can see the discrete effects don't influence the data trend. Discrete and original data is very close and it can ensure the discrete effect, the time characteristics and the data trend at the same time.

7.3. Attribute Reduction of QAR Data Containing Time Characteristics in Decision Tables

Two attributes related to engine air parking is extracted from QAR data in the CSV file and the expert knowledge shows that when the engine air parking occurs, the Air/ground=air and $N2 < 25$. The threshold of N2 can be properly adjusted. Air/ground is a state data, and it only has two types: air and ground. The discrete state of N1 and N2 are {0,1,2,3,4,} and decision attribute $D = \{yes, no\}$. Sequential decision table of engine shutdown is shown in table 3. It was converted to a non-timing decision table shown in Table 4.

Table 3. Engine Shutdown Timing Decision Table

t	state	N1	N2	shutdown
1	S1	0	0	No
2	S2	1	1	No
3	S3	2	1	No
4	S4	2	2	No
5	S5	3	2	No
6	S6	3	3	Yes
7	S7	4	4	Yes
8	S8	3	3	Yes
9	S9	3	2	yes

Table 4. Engine Shutdown Non-sequential Information Decision Table

	N1 (t)	N2 (t)	N1 (n1) (t-1)	N2 (n2) (t-1)	fault (D)
S1	1	1	0	0	No
S2	2	1	1	1	No
S3	2	2	2	1	No
S4	3	2	2	2	No
S5	3	3	3	2	Yes
S6	4	4	3	3	Yes
S7	3	3	4	4	No
S8	3	2	3	3	No

We select most of the frequent attribute in the recent time slice, then delete items which contain it in the distinguish matrix and calculate the frequency of attribute until the distinguish matrix is empty. So we can get new difference Matrix (take N1 or N2) and frequency tables such as Table 5 and Table 6.

Table 5. The New Difference Matrix (Take N2)

	S1	S2	S3	S4	S5	S6	S7	S8
S1	--							
S2	N1n1n2	--						
S3	--	--	--					
S4	--	--	N1n2	--				
S5	--	--	--	--	--			
S6	--	--	--	--	--	--		
S7	--	--	--	--	n1n2	--	--	
S8	--	--	N1n1n2	n1n2	--	--	--	--

Table 6. Ultimate Attribute Occurrences Frequency Table

attribute	n1	n2
Number of occurrences (frequency)	2	2

Finally we can obtain reduction attribute (N1, N2, n1), and (N1, N2, n2). The number of attributes is reduced to 1/4 from the reduction results. The reduction efficiency is very effective and this method keeps the time characteristic of QAR data. The method can conduct the reduction to attributes of poor correlation and low-time characteristics, leave abnormal data block of fault rather than some fault point, and greatly increase the possibility of abnormal data mining in extracting rules.

7.4. Rules Acquisition

We can create a new decision table of engine failure information based on the attribute reduction result (N1, N2, n1). The new decision table is as shown in Table 7.

Table 7. New Decision Table of Engine Shutdown Fault Information

	Current N1	Current N2	Last N1(n1)	Fault (D)
S1	1	1	0	No
S2	2	1	1	No
S3	2	2	2	No
S4	3	2	2	No
S5	3	3	3	Yes
S6	4	4	3	Yes
S7	3	3	4	No
S8	3	2	3	No

(1) We can get the whole rule set according to the general method and list some of them.

Rule 1: the current N1 (1) → whether the fault occurred: NO(0)

Rule 2: the current N1 (2) ∧ the current N2 (1) → whether the fault occurred NO (0)

Rule 3: the current N2 (3) ∧ the current n1 (3) → whether the fault occurred YES (1)

Rule 4: the current N1 (4) ∧ the current N2 (4) → whether the fault occurred YES (1)

Rule 5: the current N2 (2) ∧ the current n1 (3) → whether the fault occurred NO (0)

(2) These rules are classified by decision attribute whether the fault occurred. The fault rules 3 and 4 are expressed by the principal disjunctive normal form.

Rule 3: the current N2 (3) \wedge the current n1 (3) \rightarrow YES (1)

Rule 4: the current N1 (4) \wedge the current N2 (4) \rightarrow YES (1)

(3) these rules are simplified qualitative representation based on the normal and abnormal range. For example, the normal range of N1 and N2 is [0,2] and the abnormal range of N1 and N2 is [3,4], and count the appear number of attributes: Current N2: 2 times; Current N1: 1 times; Last n1: 1 times.

Rule 3: the current N2 (abnormal) \wedge the current n1 (abnormal) \rightarrow YES (1)

Rule 4: the current N1 (abnormal) \wedge the current N2 (abnormal) \rightarrow YES (1)

(4) We can again extract rules according to the above rules, the number of statistical attributes and the range of normal or abnormal. It can be concluded that the most important abnormal attribute is current N2 when parking fault occurs. Therefore, we should first investigate the engine-related parts of N2 when conduct engine maintenance.

7.4. Model Validation and Evaluation

This paper selects three months QAR fault data of Boeing 737 of engine fault, flap fault, bump, air-conditioning fault provided by Ameco, and builds a large number of fault decision information table and minings a lot of rules based on this detection method with QAR abnormal data. The statistical fault decision information table of QAR data is shown in Table 8. The data in the table is the average results of 20 times. The run time is not static, but the basic data is in this range. The algorithm is stability and reliable.

Table 8. Statistical Information of QAR Fault Decision Table after Reduction

Initial input	Objects (U): 76, Condition attributes: 45, Decision attribute: 1						
Calculation Name (ms)	initialing file	difference matrix and attribute frequency	sort on difference matrix	attribute reduction	Attribute reduction set	The number of rules	Total time
data	15	16	31	0	[0]	347	78

We can get the important attribute related to flap fault according to the frequency based on rule acquisition method in this paper. The result is shown in Table 9. The occurrence number of L.E.Flappos, Pitch and Roll are 37, 35 and 38, which are consistent with the important attributes given by expert when the flap fault occurred. It demonstrates that the detection algorithm is feasible and effective.

Table 9. Related Attribute Frequency of Flap Fault

Attribute name	T.E.Flappos	L.E.Flappos	FlapHandle	CntrlCol	CntrlWhl	Pitch	Roll
Attribute frequency	12	37	2	11	13	35	38

8. Conclusion

The paper proposed an anomaly detection algorithm of QAR data based on attribute support of rough set according to the characteristics of QAR data. The experimental results show that this method not only retains the time characteristics of the QAR data, but also

enhances the relationship between conditional attributes and decision attributes, improves the efficiency of attribute reduction and rule extraction and it is feasible, reliable and effective.

Acknowledgements

This work was supported by National Natural Science Foundation of China Civil Aviation (No.61179063), and Fundamental research funds for the Central Universities (ZXH2012P009).

References

- [1] Q. Li Yong, "Research on the aircraft fault prediction and fault diagnosis system based on flight data [D]", Nanjing: Nanjing University of Aeronautics and Astronautics, (2007).
- [2] NASA AMES research center, APMS OVERVIEW [EB/OL]. <http://amps.arc.nasa.gov>. 2007-12-04
- [3] Scientific Monitoring Inc. I-Trend Introduction [EB/OL]. <http://www.scientificmonitoring.com>. 2007-10-10
- [4] H. Fang, "The key technology research on real-time monitoring of Aircraft flight safety [D]: [PhD thesis]", Nanjing: Nanjing University of Aeronautics and Astronautics, (2008).
- [5] T. Wang, "Research on the flight safety model based on QAR data [D]: [Master thesis]", Tianjin: Civil Aviation University of China, (2008).
- [6] L. Chao, "Analysis of time-series data based on rough set theory [D]", Changsha: Central South University, (2005).
- [7] C. Na, "Research on data mining method based on rough set theory [D]", Changchun: Changchun University of Technology, (2006).
- [8] M. Wu, "Research on attribute reduction of data mining algorithms based on rough set [D]", Xi'an: Xi'an University of Electronic Science and Technology, (2006).
- [9] Alam S. S. Ghosh Shrabonti. Rank by AHP:A Rough Approach, ISCF2002, pp. 185-190.
- [10] A. I. Dimitras, "Business failure prediction using rough sets [J]", European Journal of Operational Research, vol. 114, (1999), pp. 263-280.
- [11] M. strong and ZX. Cai, "Discovery on personalized decision rule: a method based on Rough Set", Control and Decision, vol. 9, (2004), pp. 994-998.
- [12] S. Osrtoff, "Temporal Logical for Real-Time System", Research Studies Perss, LTD, John Wiley & Sons Inc. (1989), pp. 21-50.
- [13] JD Berndt and J. Clifford, "Finding Patterns in Time Series: A Dynamic Programming Approach", In: U.M. Fayyad, G. Piatetsky — Shapiro, P. Smyth, eatl, Eds, Advances in Knowledge Discovery and Data Mining. Massachusetts: MIT Press, (1996), 18. 66 -74
- [14] TB Anders, "Mining Time Series Using Rough Set-A Case Study", In: Komorowski, J, Zkytkow, J. Eds. Proceeding of the First European Symposium, PKDD'97.Trondheim,Norway:Springer Lecture Notes in Artificial Intelligence, (1997), pp. 351-358.

Authors



Hui Yang, She was born in 1957 and a Professor of Civil Aviation University of China, the main research fields of her is artificial intelligence, data mining and fault diagnosis.



Chunjing Xiao, She was born in 1978, now is a lecturer at Civil Aviation University of china, the main research fields are artificial intelligence and data mining.



Yongwei Qiao, She was born in 1976, now is a lecturer at Civil Aviation University of china, the main research fields are aircraft engineering and aircraft fault diagnosis.

