

User behaviors attributes of database anomaly detection model

SaifAldeen Salim Ahmed

University : University of Technology in Iraq

College : Computer Science

department : Information System

ABSTRACT

This paper includes the description of designing a data-base anomaly detection system, which is capable of being more precise in depicting the behaviors of individuals and improving data-base abnormal detecting correctness. In designing the system, the Apriori approach is used first and depends on the k-means clustering and the Apriori methods. It is capable of more efficiently exploiting users' behaviors, and the data-base abnormal more efficient detecting. The relevant studies show that Apriori method according to time efficiency and precision of detection is more optimal than the sole utilization according to association rules mining approaches Apriori method.

Keywords: Anomaly detection structure, Apriori, users' behaviors.

Date of Submission: 02-11-2017

Date of acceptance: 08-12-2017

I. INTRODUCTION

Due to the improvement of intruder detecting and data mining approaches, the data-base audit registry that's considered to be a form of passive security element has been altered (1). It's a significant way of ensuring the security of data-base which data-base exception is efficiently found by users' behaviors profiles of the data-base mining (2). The above mentioned users' behaviors profiles is the repetitive accessing to specific resources which includes the data-base, data-tables etc. This repetitive operating of resources is properties of the users' behaviors patterns.

Nowadays, there exist basically 3 elements of deficiency exist in the user behaviors profiles mining. Initially, because of the different kinds of audit registries in the DBMS, it's quite hard selecting which audit logs maybe efficiently utilized for mining the behaviors profile of the users.

Secondly, due to the fact that the existing algorithms are not capable of properly describing the behaviors profile of the user, causing higher false positive rate of data-base anomalous detecting. Thirdly, even though part of the existing approaches are capable of achieving the mining the behaviors profile of users, when meaning the huge user data-base accessing registries they are noticeably inefficient (3). This research basically performs the next 2 main experiments for solving the above mentioned issues: First: A data-base structure is modeled for anomaly detecting. Second:

Apriori approach based on this structure is suggested which is a more sufficient detecting of abnormal data-base approach.

II. RELATED WORKS

There are numerous researches that propose using of data mining approaches in registry file analyzing procedure or the detection of security risks generally. One of the 1st method that utilize data mining approaches in intruder detecting was suggested by Lee and Stolfo (8). Who used 2 methods for detection, the association rules method and the redundant episodes approach. They illustrated that via the analysis of audit data it's possible discovering intrusion patterns.

Frei and Rennhard (9) utilized another method for searching for anomaly in registry files. They generated the Histogram Matrix, a registry file visualizing approach which aids security administrators find the anomalies. This method operates on each textual registry file. It depends on the idea that the brain is sufficient in the detection of patterns while observing images, thus, the registry file is viewed in a way which it's easy to observe changes from regular behaviors.

Fu and others (10) suggested an approach for anomaly detecting in unstructured system registries which doesn't need any application specific knowledge. In addition, they added an approach for the extraction of registry keys from free text messages.

Makanju, et.al(11) suggested a hybrid registryalarmdetectingmodel, with the use of each of anomaly and signature-based detectingapproaches.

III. BACKGROUND ON ANOMALY DETECTION

An anomaly (or outlier) is an observation that looks inconsistent with the rest(majority) of the dataset, and hence arouses the suspicion that it can be generated by a different mechanism. The objective of anomaly detection is to mine unusual and information of interest from a large amount of data. Detecting anomaly is extensively studied in different aspects, like the statistics, data mining, machine learning and information theory, and its applications have been greatly expanded to multiple areas like detecting of fraud, network intrusion, health monitoring, environmental monitoring and performance analysis.

A straightforward solution for anomaly detection is to construct a pattern of normal observations, and then one can use the pattern to identify anomalies. When applying the pattern on test data, the observations whose properties follow the regular pattern are labeled as normal, and those that deviate noticeably from the regular pattern are labeled as anomalies. According to the availability of labeled training data, anomaly detecting approaches may be divided into three main classes, which are the supervised, semi-supervised and unsupervised approaches.

In the first category all the training samples are required to be paired with a label or desired output, i.e., normal or abnormal observations, for the characterization of all anomalies or non-anomalies. Semi-supervised learning techniques make use of unlabeled records as well as labeled ones for training. Typically semi-supervised techniques are trained with a large amount of unlabeled records and a small amount of labeled records. It should be noted that pre-labeled data is not constantly available nor easily obtained in several of real-life utilizations, in addition, new kinds of observations (normal or abnormal) could occur that aren't included in the labeled training data. Unsupervised methods are often more appealing for anomaly detection, since they require no labeled data, rather they apply certain criteria to identify anomalous observations. An example type of such techniques are distance-based approaches, e.g., classifying records according to the average distance between each one of the data records to its mapping the nearest neighbor observations. If the measured distance for a specific record significantly exceeds nearest neighbor distance of all objects then the data record is considered as an anomaly, otherwise it is considered to be normal.

IV. DATABASE ANOMALOUS DETECTION MODEL

This paper includes the description of the design of a data-base anomaly detecting structure, as illustrated in Figure(1).

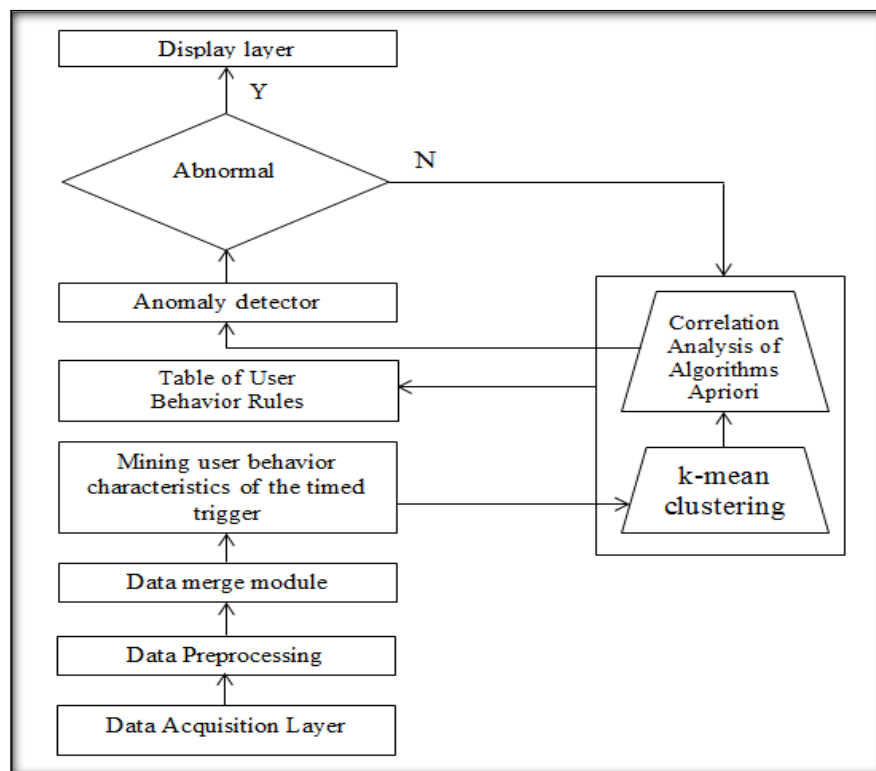


Figure (1): data-base anomaly detecting structure

The structure is basically made up of the following 5 modules, which are: data acquisition layer, data pre-processing, data merging, user behaviors mining, anomaly detecting.

The basic task: benefitting from registry miner and other data-base registry collecting tools or data-base system's own registry analyzing tools, like Oracle 11g Profiler for completing these of the existing testing data and the auditing of the training data.

4-1 Data Acquisition Layer Module

	C01_PRICE	C01_HARD	C01_PREMIUM	C01_ID	C01_TREND	C01_SCREEN	C01_SPEED	C01_RAM	C01_MULTI	C01_CD	C01_ADS
1	506	33 no		506	139	4	1,775	170 no	14	yes	
2	507	33 no		507	139	8	2,490	340 no	15	yes	
3	508	33 no		508	139	16	3,599	340 no	17	yes	
4	509	50 no		509	139	8	2,690	340 no	14	yes	
5	510	66 no		510	139	8	3,195	540 no	15	yes	
6	511	66 no		511	139	8	3,695	452 no	14	yes	
7	512	50 no		512	139	8	2,645	250 no	15	yes	
8	513	66 no		513	139	16	3,090	452 no	15	yes	
9	514	66 no		514	139	2	1,890	107 no	15	yes	
10	515	33 no		515	139	4	1,999	170 no	14	yes	
11	516	50 no		516	139	8	2,935	250 no	17	yes	
12	517	25 no		517	139	4	1,990	214 no	14	yes	
13	518	50 no		518	139	4	2,290	214 no	14	yes	
14	519	66 no		519	139	4	2,390	214 no	14	yes	
15	520	50 no		520	139	4	2,025	170 no	14	yes	
16	521	33 no		521	139	4	2,095	214 no	14	yes	
17	522	33 no		522	139	8	2,590	340 no	14	yes	
18	523	25 no		523	139	4	1,499	170 no	14	yes	
19	524	33 no		524	139	8	2,325	250 no	15	yes	
20	525	66 no		525	139	4	2,099	120 no	14	yes	
21	526	66 no		526	139	16	2,999	245 no	15	yes	
22	527	66 no		527	139	8	2,790	340 no	15	yes	
23	528	33 no		528	139	2	1,590	107 no	14	yes	
24	529	50 no		529	139	4	2,499	170 no	14	yes	
25	530	50 no		530	139	8	2,575	250 no	15	yes	
26	531	25 no		531	139	8	2,390	340 no	15	yes	
27	532	33 no		532	139	2	1,495	107 no	14	yes	

Fig. (2): the type of gathered auditing registry data detecting

4-2 Data Pre-processing Module

The basic task: post collecting data (audit training data, the current data inspection), for removing the genuine noise, like the inconsistent auditing information, no particular operating persons, no operating item information and other valuable data, in addition to stemming some information isn't in association to the needed data from relevant communications links protocols when data-base servers are reconnected via clients.

4-3 Data Merging Module

The basic task: initially, to statistic the number of functional data-base items, like the data table, data view, etc. used by this person. Then, dealing with the task item and user via the numeric identifying, and after that storing in the established sessions registry tables. Moreover, every

one of the connected operation records is considered a transaction T, each one of the transactions T generates the data-base D, which will be utilized to mine user behaviors properties. Every one of the transactions belonged to the data-base D is made up of the following fields: the session connection ID, Data-base process user, process object, process path, etc. with an overall of 12 fields.

4-4 User Behaviors Mining Module

The basic task: Firstly, the procedure on the data-base run via the users and their data-base items must be clustered with the use of the k-means approach, this way finishing the Preliminary Characterizing of the behaviors of the user. In addition making the 1st preparation for more mining of user behaviors properties. Secondly, improving the sufficiency of mining property rules

concerning the user behaviors with the use of the Apriori algorithm.

4-5 Anomaly Detecting Module

The basic task: comparing the training stage of the rules generated by a regular user (obtained from old-

rules table) with the testing stage of mining association rules (obtained from new rules), in the case where the data-base exception happened then the irregular details have to be timely registered and shown.

The following is the anomaly detecting algorithm:

Input: The current audit data to be detected

Output: The information of detecting anomalies

Method:

Step1: acquire the begin time of detection: start time

Step2: While (each rule in old-rules) {

Step3: While (each rule in new-rules) {

Step4: If ((the antecedent of rules in old-rules
 == the antecedent of rules in new-rules)

&& (the consequent of rules in old-rules
 == the consequent of rules in new-rules)

&& (the difference of support and confidence
 between the two rules greater than 35%))

Step5: {make intrusion information record
 in the Intrusion table} }

Step6: the end time of detection: endtime

Step7: record the time spent testing(starttime-endtime)

Step8: If (interrupt) {

Step9: break } }

V. THE RESEARCH OF THE MODEL OF MAIN ECHNOLOGIES

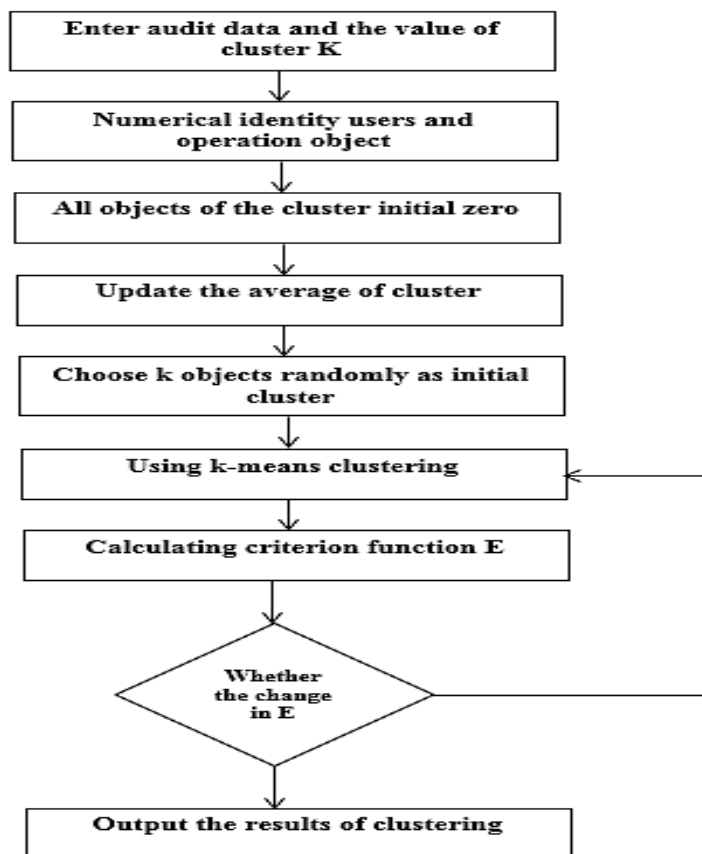


Figure (4): the procedure of clustering

5-1 Preliminary Characterization of User Behaviors Properties

According to the K-means

This research includes the description of using the k-means clustering method for dealing with the data-base users and operating object data, for the sake of characterizing the users' initial behaviors profiles. The main task E of k-means method is described in equation (1).

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - m_i|^2 \quad (1)$$

The clustering procedure of data-base items is depicted in Fig. 4.

5-2 Apriori association rule mining algorithm

Theorem (2): $\forall c \in C_k, R(c)$ could be produced by the two items of $R(x), R(y) (x \neq y)$ in R_{k-1} , and $R(c) = R(x) \cap R(y)$.

2) The main idea behind this Algorithm

Data storing format which is a transaction identifier is corresponding to several types is called the standardized data formatting in the data-base of transactions. While, data storing format which is a transaction identifier is corresponding to several transaction identifiers relevant to the element is called the vertical data formatting.

The data which has standardized data formatting is transformed to vertical formatting via the scan process of the data-base one time. In the same time, every one of the items in the itemset and transaction identifiers which corresponds to the elements are stored separate from one another with the use of two dynamical list storages. The supporting counting of C_k may be reached by $L_{k-1} \cap L_1$ with no need to repeat

This approach has been suggested for the Boolean associating rules mining redundant element group of the main method, in the year of 1994 by Agrawal et.al.(6). Even though the use of Apriori method own nature may raise the effectiveness to a specific degree. For the sake of mining the data-base registry data more efficiently, this study presents the Apriori method.

1) Relevant theories and conclusions

Theorem (1) : Supposing redundant k itemset is capable of generating (k+1) Itemset, then the number of itemset of the redundant k itemset is definitely $>= k$.

data-base scanning for the sake of obtaining the supporting counting of C_k . Connecting conditions of Apriori method are improved with the use of the conclusions which have been proven in this study.

Redundant itemset k are obtained via connecting conditions of $L_{k-1} \cap L_1$, and the final item of L_{k-1} based on the index is subject to comparison with every one of the items of L_1 when $L_{k-1} \cap L_1$, avoiding duplicating comparing of connecting conditions of $L_{k-1} \cap L_1$, and with no consideration if (k-1) is sub-set of candidate itemset C_k is in L_{k-1} .

The repetitions of redundant item set k are improved via Theorem 1, and the supporting counts of nominated item set k are handled with Theorem 2.

3) The following is the description of the Apriori algorithm:

```

Input:   D: Transaction database
         Min_support : minimum supportcount
Output:  L: Frequent item sets in D

Method:
Step1: C1 = {The set of all items}
Step2: for each Candidate itemsets c ∈ C1 {R(C) = σ ;}
Step3: for each transaction t ∈ D {
Step4: Ct = Sub-Set (t, C1);
        // the table of generated each item ID in item set
Step5: for each c ∈ Ct {R(c) = R(c) ∪ t_TID ;}
Step6: Count_Support (C1);
Step7: L1 = {c ∈ C1 | c_delete = 0}; R1 = {R(c) | c ∈ L1};
Step8: Up_Sort (L1);
Step9: for (k=2; |Lk-1| > k-1; k++) {
Step10: Up_Sort (Lk-1); Ck = apriori_gen (Lk-1);
Step11: L = ∪ kLk;
Step12: function Count_Support(C: candidate_k-itemset)
        for each Candidate itemsets c ∈ C; {
            if (c_support = R(c) < min_support) {
                c_delete = 1; } }
Step13: function apriori_gen(Lk-1: Frequent (k-1) itemsets)
        { for each lk-1 ∈ Lk-1 {
            if (c_support = |R(c)| < min_support) { c_delete = 1; }
        }
    }
    }
    
```

VI. EXPERIMENTAL RESULTS AND DISCUSSIONS

Experimental Results:

The precision of the Anomaly Detecting of the Model the comparing of detecting precision between user behaviors diggers paired via anomaly

detections suggested in the model and the conventional structure of data-base anomalies detecting according to the Apriori algorithm. In the experiments, the use of trend data-base events for testing the detecting precision. The outcome is illustrated in Figure (5).

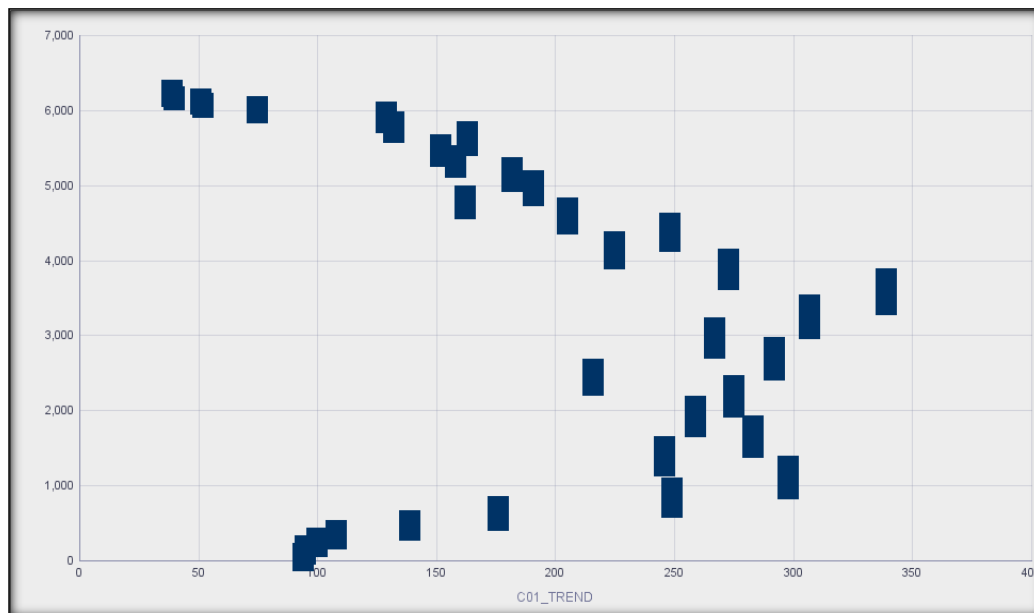


Fig. (5): Experimental Outputs
 Detection Accuracy & Database events

The difference between the behaviors of two users (user1 & user2) by using Apriori Algorithm as Shown in Figure (6), the explore Data behaviors of user shown in figure (7).

The results of the experimentation proved the fact that mining

properties of users behaviors post the auditing registry clustering may be excavating user behaviors rules with more efficiency, therefore, the precision of data-base anomaly detecting is also developed with more efficiency.

Itemsets: 1,000 out of 122,027

ID	Items	Support(%)	Item Count
52	110	76.9231	1
1197	88, 110	69.2308	2
1486	100, 105	69.2308	2
22	75	69.2308	1
46	100	69.2308	1
49	105	69.2308	1
10930	75, 105, 100	61.5385	3
14421	88, 105, 100	61.5385	3
17098	100, 110, 105	61.5385	3
666	70, 88	61.5385	2
822	75, 88	61.5385	2
833	75, 100	61.5385	2
836	75, 105	61.5385	2

Itemset Details:

ID: 22

Item List

75

Support (%)	69.2308
Item Count	1

Figure (6): Apriori user 1

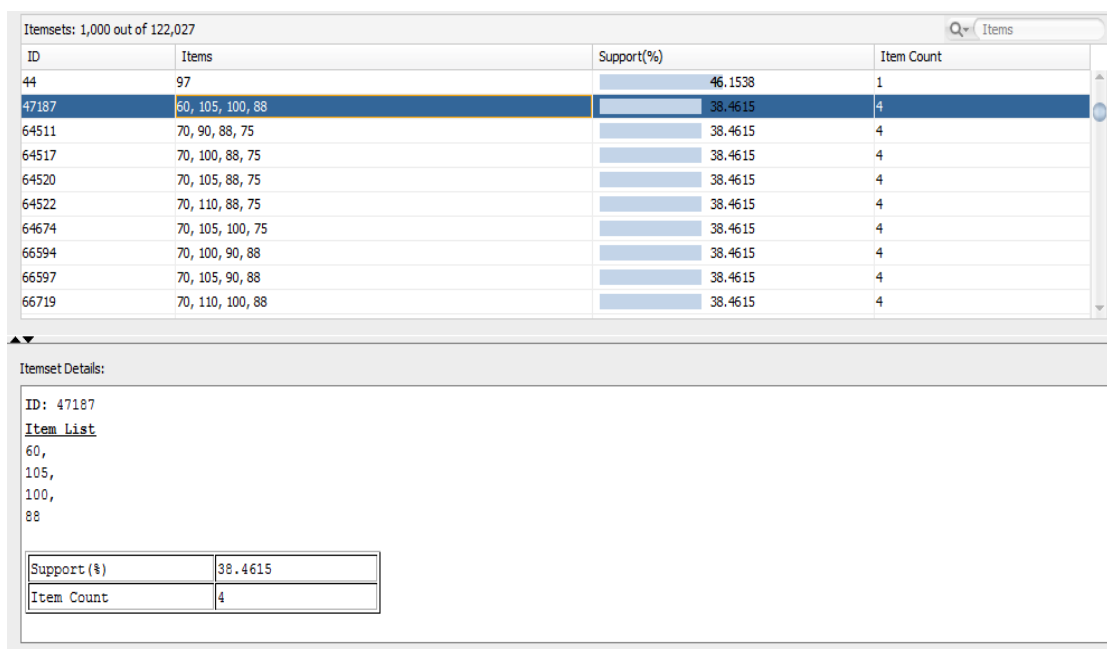


Figure (6): Apriori user2

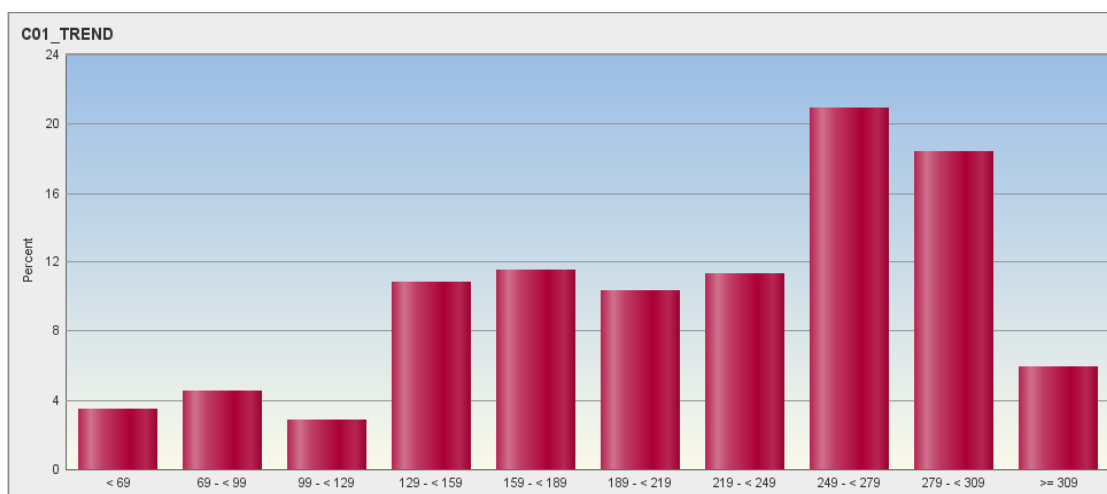


Figure (7): Explore Data

VII. CONCLUSION

This study explains a structure of data-base anomaly detecting for the sake of efficiently improving the precision of data-base anomaly detecting. Based on this structure, a sufficient Apriori method is suggested. This algorithm may be able to get rid of

some insignificant rules to a specific degree and depicting data-base users' behaviors profiles more clearly. At last, sequence patterns mining used in user that access the data-base registries will be a significant research topic in the future.

SaifAldeen Salim Ahmed "User behaviors attributes of database anomaly detection model." International Journal of Engineering Research and Applications (IJERA) , vol. 7, no. 12, 2017, pp. 29-35.