

1 Article

2 Applications of Data Mining Algorithms for Network 3 Security

4 Kai Chain

5 Department of Computer and Information Science, Republic of China Military Academy, Kaohsiung 83059,
6 Taiwan; chainkai@mail2000.com.tw

7 * Correspondence: chainkai@mail2000.com.tw; Tel.: +886-7-743-8179; Fax: +886-7-747-9515

8 **Abstract:** Typical modern information systems are required to process copious data. Conventional
9 manual approaches can no longer effectively analyze such massive amounts of data, and thus
10 humans resort to smart techniques and tools to complement human effort. Currently, network
11 security events occur frequently, and generate abundant log and alert files. Processing such vast
12 quantities of data particularly requires smart techniques. This study reviewed several crucial
13 developments of existent data mining algorithms, including those that compile alerts generated by
14 heterogeneous IDSs into scenarios and employ various HMMs to detect complex network attacks.
15 Moreover, sequential pattern mining algorithms were examined to develop multi-step intrusion
16 detection. These studies can focus on applying these algorithms in practical settings to effectively
17 reduce the occurrence of false alerts. This article researched the application of data mining
18 algorithms in network security. The academic community has recently generated numerous
19 studies on this topic.

20 **Keywords:** data mining; network security; association rules; DDoS

21

22 1. Introduction

23 In short, data mining can be defined as the computational process of discovering meaningful
24 patterns or rules through the semi-automatic or automatic exploration and analysis of copious data.
25 The most commonly used analytical techniques for general data mining algorithms are cluster
26 analysis, classification analysis, link analysis, and sequence analysis. Cluster analysis is the task of
27 grouping data into multiple categories or clusters so that data within the same cluster share a higher
28 similarity with each other than with data in other clusters. Clustering can be used to identify dense
29 and sparse areas in the data space, thereby revealing the global distribution of the data and the
30 relationships between data characteristics. As the most widely applied data mining technique,
31 classification analysis can be used to model the relationships between crucial data or to predict data
32 trends. It involves classifying new data under predefined categories (e.g., normal and invaded). Two
33 of the most common classification methods are decision trees and rules. In 1993, Agrawal et al. at the
34 IBM Almaden Research Center proposed a problem concerning the association rules between
35 itemsets in a customer transaction database and presented the Apriori algorithm [13]. Their study
36 prompted intensive research on the problem of discovering association rules; this further facilitated
37 the application of data mining analysis techniques in various fields, particularly classification
38 analysis. Sequence analysis is a longitudinal analysis technique that identifies data characteristics by
39 employing temporal variables.

40 Cluster analysis, link analysis, and sequence analysis are the results of theories and research
41 that involved analysis of small-scale databases. Data mining has now extracted or integrated large
42 quantities of data from several large-scale databases. For instance, Wal-Mart utilizes an effective
43 data warehouse system that enables the provision of clear and detailed sales information to clients
44 and product suppliers; this enables presenting data regarding what products are most required by
45 its clients and when the clients need those products. Most importantly, Wal-Mart applied data

46 mining analyses of customer behavior models to increase the effectiveness of inventory control and
47 customer relationship management. The application of data mining for commercial purposes has
48 presently reached a certain level of maturity. In recent years, data mining techniques are often used
49 in research on security [1-12]. One of the most typical examples of a data mining application in
50 network security is an intrusion detection system (IDS). In general, an IDS collects a large quantity of
51 alert information. However, information security personnel cannot identify any associations from
52 large amounts of disorganized alert information. Nevertheless, an IDS that incorporates data mining
53 techniques facilitate identifying accurate and appropriate data characteristics from large quantities
54 of alert information by examining the association between multiple characteristics. The adequacy of
55 the data mining technique applied determines the usefulness of the data classification result, which
56 then affects the accuracy and effectiveness of intrusion detection. Thus, special attention must be
57 paid to the applicability of data mining algorithms under different scenarios; this explains why
58 further analysis should be performed on existing data mining algorithms.

59 [1-4] utilized an algorithm that augments existing scenarios with alerts generated by
60 heterogeneous IDSs; this algorithm considers each new alert by estimating the probability that this
61 new alert belongs to a given scenario and ascribes each alert to the most likely candidate scenario.
62 This involves real-time data mining on the basis of ideal scenarios. In addition, [1, 4] evaluated the
63 proposed algorithm by comparing it with two other probability estimation techniques; because this
64 method only requires simple data entry, it can serve as a more appropriate probability estimation
65 technique and is thereby more likely to obtain desired results. [5-9] employed hidden Markov
66 models (HMMs) to detect complex network attacks. Although a typical HMM is more effective in
67 detecting such complex intrusion events than are decision trees and neural networks, it has certain
68 drawbacks. Collaborative network attacks are often conducted over a long period (e.g., a few weeks
69 or months). Therefore, developing methods for detecting such slow-paced attacks in advance is one
70 of the primary research topics in this field. These methods must also be capable of detecting attacks
71 that involve highly complex methodologies, and must feature a set of procedures to detect
72 slow-paced attacks, and must be able to counteract attacks that have been detected. Sensor failures
73 or interference in transmission channels (such as open wireless networks), can affect the detection of
74 sequential attacks and cause errors. [10-12] proposed a sequential pattern mining technique that
75 incorporated alert information associated with distributed denial-of-service (DDoS) attacks and
76 reduce the repetition and complexity of abundant alert information received by the IDS, thereby
77 expediting the identification of the DDoS attack types. To identify the alert sequence of DDoS
78 attacks, [10] employed a particular HMM to compare algorithms. The comparison results indicated
79 that sequential pattern mining algorithms are more solid than other algorithms are because they do
80 not result in mutual exclusivity of correlation data during a DDoS attack and can increase the
81 accuracy of multi-step intrusion detection. The application of that HMM in experimental
82 comparison of algorithms verified whether the proposed algorithm exhibited favorable effects for
83 identifying the sequences of DDoS attacks and ensured effective reduction in alert complexity.

84 This paper is organized as follows. Section 2 elucidates the employed data mining algorithms.
85 Section 3 presents the data mining algorithms deemed as adequate for application to network
86 security. Section 4 presents the conclusions of this paper.

87 2. Data Mining Algorithms

88 In general, data mining involves converting collected data into preferred formats for analysis
89 and employing algorithms to derive useful intrusion detection models for identifying known or
90 unknown intrusion types. Data mining studies do not have a long history; they emerged mainly
91 from a study by Agrawal et al. (1993) [13], which stimulated widespread academic enthusiasm in
92 this topic.

93 Research manuscripts reporting large datasets that are deposited in a publicly available
94 database should specify where the data have been deposited and provide the relevant accession
95 numbers. If the accession numbers have not yet been obtained at the time of submission, please state
96 that they will be provided during review. They must be provided prior to publication.

97 Association rule mining is a procedure used to discover interesting associations or relevant
 98 connections between itemsets in abundant data. As increasing amounts of data were collected and
 99 stored, database users began to demonstrate interest in discovering association rules in their
 100 databases. Discovery of interesting association rules is conducive to formulating commercial
 101 policies, such as design categorization and cross selling. The purpose of such discovery (or mining)
 102 is to identify any interesting relationships between items within a given data set. Support and
 103 confidence are interestingness measures used for evaluating association rules; they reflect the
 104 usefulness and certainty of association rules, respectively. Interestingness measures can also be used
 105 to assess a model's simplicity, certainty, utility, and novelty. Association rule mining is typically
 106 performed with two steps:

- 107 (1) Find all frequent itemsets: Frequent itemsets are defined as itemsets with a high probability of
 108 appearance in a database. If the obtained itemsets meet this definition, they are considered as
 109 frequent itemsets.
- 110 (2) Generate strong association rules from qualified frequent itemsets: If the obtained itemsets meet
 111 the criteria of minimum support and minimum confidence, they can be defined as strong
 112 association rules.

113 The first step is clarified as follows. The Apriori algorithm is one of the most influential
 114 algorithms used for association rule mining; it serves as the basis for comparing various algorithms
 115 and is a basic technique for deriving different algorithms. This study designated the Apriori
 116 algorithm as the first algorithm to be thoroughly analyzed because it shares several key terms
 117 with the other algorithms discussed herein. The Apriori algorithm requires verification on the basis
 118 of the characteristics of frequent itemsets; to meet the criterion of this algorithm, all nonempty
 119 subsets in a frequent itemset must be frequent. The Apriori algorithm involves layer-by-layer
 120 iteration and mines $(k+1)$ -itemsets by using k -itemsets. This mining process can be illustrated as
 121 follows: First, find all frequent 1-itemsets and label them as L_1 . Subsequently, L_1 is used to find
 122 frequent 2-itemsets (L_2), which are then used to find L_3 . This process is repeated until no more
 123 frequent k -itemsets can be found. This algorithm scans the database when each search of L_k is
 124 completed and generates candidate itemsets from the itemset L obtained in the previous search. The
 125 candidate itemset meeting the Apriori criterion is then used as the itemset L for the following search.

126 To increase the efficiency of layer-by-layer iteration used for mining frequent itemsets, reducing
 127 research space is a crucial goal of the Apriori algorithm. Figure 1 presents the pseudocode used in
 128 the Apriori algorithm, where *Input* represents the transaction item of D , *min_sup* denotes the
 129 minimum support count, and *Output* is the frequent itemset L in the database.

130 Once the first step has been completed by finding all frequent itemsets from the transaction
 131 data in database D , the second step (generate strong association rules) can be performed
 132 immediately. The identification of strong association rules requires meeting the criteria of minimum
 133 support and minimum confidence. The following equation can be used to calculate confidence:

$$134 \text{confidence}(A \Rightarrow B) = P(A|B) = \frac{\text{support_count}(A \cup B)}{\text{support_count}(A)} \quad (1)$$

135 where $\text{support_count}(A \cup B)$ denotes the number of transactions containing itemset $A \cup B$ and
 136 $\text{support_count}(A)$ denotes the number of transactions containing itemset A .

```

Method:
1:  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;
2: for ( $k=2$ ;  $L_{k-1} \neq \emptyset$ ;  $k++$ ){
3:    $C_k = \text{apriori\_gen}(L_{k-1}, \text{min\_sup})$ ;
4:   for each transaction  $t \in D$  { //scan D for counts
5:      $C_t = \text{subset}(C_k, t)$ ; //get the subset of t that are candidates
6:     for each candidate  $c \in C_t$ 
7:        $c.\text{count}++$ ;
8:   }
9:    $L_k = \{c \in C_t \mid c.\text{count} = \text{min\_sup}\}$ 
10: }
11: return  $L = \cup_k L_k$ ;

Procedure apriori_gen( $L_{k-1}$ :frequent( $k-1$ )-itemsets;  $\text{min\_sup}$ :minimum support threshold)
1: for each itemset  $l_1 \in L_{k-1}$ 
2:   for each itemset  $l_2 \in L_{k-1}$ 
3:     if ( $l_1[1]=l_2[1]$ ) ? ( $l_1[2]=l_2[2]$ ) ? ... ?
       ( $l_1[k-2]=l_2[k-2]$ ) ? ( $l_1[k-1]=l_2[k-1]$ ) then {
4:        $C = l_1 \cup l_2$ ; //join step: generate candidates
5:       if has_infrequent_subset( $c, L_{k-1}$ ) then
6:         delete  $c$ ; //prune step: remove unfruitful candidate
7:       else add  $c$  to  $C_k$ ;
8:     }
9:   return  $C_k$ ;

Procedure has_infrequent_subset( $c$ : candidate  $k$ -itemsets;  $L_{k-1}$ : frequent( $k-1$ )-itemsets);
1: for each ( $k-1$ )-subset  $s$  of  $c$ 
2:   if  $s \in L_{k-1}$  then
3:     return TRUE;
4: return FALSE;

```

137

138

139

Figure 1. Apriori algorithm employed to discover frequent itemsets during Boolean association rule mining

140

141

142

143

Because all rules are generated from frequent itemsets, each of these rules automatically satisfy the minimum support count. Frequent itemsets and their support levels are included in the list of items to ensure they can be immediately mined. An example of applying the Apriori algorithm to identify frequent itemsets in a database is presented as follows:

TID	List of item_IDs
T 100	11,12,15
T 200	12,14
T 300	12,13
T 400	11,12,14
T 500	11,13
T 600	12,13
T 700	11,13
T 800	11,12,13,15
T 900	11,12,13

144

145

Figure 2. Data from the transaction database of Company Branch X

146 In this example, item IDs 11 and 12 are assumed to represent attacks in the field of network
 147 security (they could just as easily be applied to represent other types of data, e.g., diapers and beer in
 148 a supermarket). Frequent itemsets were generated from the following steps:

149 Step 1: In iteration 1, consider every item as a member of the candidate 1-itemset, C_1 . Scan all
 150 transactions roughly to count all items that appear (Figure 3).

 C_1

Item set	Sup count
{I1}	6
{I2}	7
{I3}	6
{I4}	2
{I5}	2

151

152

Figure 3. Members of C_1

153 Step 2: Assume the minimum support count to be 2 (i.e., $min_sup=2/9=22\%$) so that the frequent
 154 1-itemset, L_1 , consists of candidate 1-itemsets meeting the threshold of the minimum support
 155 count.

 L_1

Item set	Sup count
{I1}	6
{I2}	7
{I3}	6
{I4}	2
{I5}	2

156

157

Figure 4. Members of L_1

158 Step 3: Search for the frequent 2-itemset L_2 . Use the algorithm to generate the set of candidate
 159 2-itemset, C_2^6 , where C_2 consists of $C_{L_1}^2$ 2-itemsets (Figure 5).

 C_2

Itemset
{I1, I2}
{I1, I3}
{I1, I4}
{I1, I4}
{I2, I3}
{I2, I4}
{I2, I5}
{I3, I4}
{I3, I5}
{I4, I5}

160

161

Figure 5. Members of C_2

162

Step 4: Scan the transactions in database D . Calculate the level of support of every candidate itemset

163

in C_2 (Figure 6).

C_2

Itemset	Sup count
{I1, I2}	4
{I1, I3}	4
{I1, I4}	1
{I1, I4}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2
{I3, I4}	0
{I3, I5}	1
{I4, I5}	0

164

165

Figure 6. Support count of C_2

166

Step 5: Confirm the frequent 2-itemset, L_2 , which consists of candidate 2-itemsets in C_2 that satisfy

167

the criterion of greater than or equal to minimum support count.

L_2

Itemset	Sup count
{I1, I2}	4
{I1, I3}	4
{I1, I5}	2
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2

168

169

Figure 7. Set L_2

170

Step 6: First, let $C_3 = \{\{I1, I2, I3\}, \{I1, I2, I5\}, \{I1, I3, I5\}, \{I2, I3, I4\}, \{I2, I3, I5\}, \{I2, I4, I5\}\}$. According to the

171

characteristics of the Apriori algorithm, all subsets in frequent itemsets should appear

172

frequently. Because any two items in the last four candidate itemsets do not appear in L_2 , these

173

are not frequent items, and thus these four candidate itemsets must be deleted from C_3 .

174

Therefore, their count must no longer be calculated when scanning D and confirming L_3 .

175

Notably, because the Apriori algorithm involves layer-by-layer iteration, once k itemsets are

176

given, the mining process starts by only checking whether their $(k-1)$ subsets are frequent.

177

Step 7: Scan the transaction in database D to confirm whether L_3 consists of candidate 3-itemsets in

178

 C_3 that satisfy the minimum support count (Figure 8).

179

180

C_3		C_3	
Itemset		Itemset	Sup count
{I1, I2, I3}	→	{I1, I2, I3}	2
{I1, I2, I5}		{I1, I2, I5}	2

Figure 8. Set C_3

181

182

183

184

Step 8: Use the algorithm to generate the candidate 4-itemset, C_4 . Because subset $\{I2, I3, I5\}$ did not appear frequently in the $\{\{I1, I2, I3, I5\}\}$ itemset, a result of C_4 , the subset was removed. Accordingly, C_4 was identified to be equal to an empty set, thereby terminating the algorithm and identified all frequent itemsets in L_3 (Figure 9).

 L_3

Itemset	Sup count
{I1, I2, I3}	2
{I1, I2, I5}	2

185

186

Figure 9. Set L_3

187

188

189

190

191

Now the example database is used to calculate which of the itemsets contain strong association rules. The third data set derived from the example data, including frequent itemset $l=\{I1, I2, I5\}$, indicates that the nonempty subsets of l are $\{I1, I2\}$, $\{I1, I5\}$, $\{I2, I5\}$, $\{I1\}$, $\{I2\}$, and $\{I5\}$. The resulting association rules are obtained from the aforementioned calculation method, each listed with its confidence:

192

$I1 \wedge I2 \Rightarrow I5$, confidence = $2/4 = 50\%$

193

$I1 \wedge I5 \Rightarrow I2$, confidence = $2/2 = 100\%$

194

$I2 \wedge I5 \Rightarrow I1$, confidence = $2/2 = 100\%$

195

$I1 \Rightarrow I2 \wedge I5$, confidence = $2/6 = 33\%$

196

$I2 \Rightarrow I1 \wedge I5$, confidence = $2/7 = 29\%$

197

$I5 \Rightarrow I1 \wedge I2$, confidence = $2/2 = 100\%$

198

If the minimum confidence is 70%, then only the second, third, and sixth rules satisfy the threshold of $\geq 70\%$. Thus, only these three items are strong association rules.

200

201

202

203

204

205

206

207

208

209

210

211

In [14, 17, 19], the main drawback of Apriori algorithm is that it must scan the database each time when generating an association file. If the database contains abundant data, this method easily increases the system load. To increase the efficiency of the Apriori algorithm, scholars have modified this algorithm using specific methods. For instance, the hash-based technique can be used to calculate the count of hash itemsets and reduce the size of the candidate k -itemset, C_k ($k > 1$). After each transaction in the database is scanned, the frequent 1-itemset L_1 is generated from the candidate 1-itemset in C_1 , and subsequently 2-itemsets can be generated from each transaction and then hashed (i.e., mapped) into the buckets in the hash table, after which the corresponding bucket counts can be increased (Figure 10). 2-itemsets with corresponding bucket counts in the hash table that are lower than the support threshold are deemed as nonfrequent and are excluded from the candidate set. The hash-based technique can largely reduce the number of k -itemsets examined (particularly when $k = 2$) [19].

bucket address	0	1	2	3	4	5
bucket count	2	2	4	2	2	4
bucket contents	{I1, I4} {I3, I5}	{I1, I5} {I1, I5}	{I2, I3} {I2, I3} {I2, I3} {I2, I3}	{I2, I4} {I2, I4}	{I2, I5} {I2, I5}	{I1, I2} {I1, I2} {I1, I2} {I1, I2}

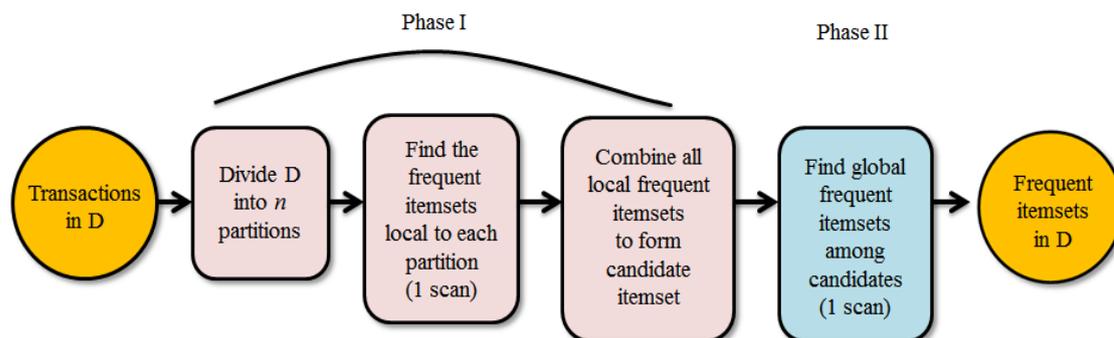
212

213

Figure 10. Hash table, H_2 , for candidate 2-itemsets

214 Note: This table was generated by scanning the transactional database when determining L_1 by C_1 . (If
215 the minimum support count is 3, then the itemsets in buckets 0, 1, 3, and 4 are nonfrequent and thus
216 should be removed from C_2).

217 Transaction reduction can also be used to improve the Apriori algorithm. This method can
218 reduce the number of transactions scanned in future iterations. A transaction without any frequent
219 k -itemsets cannot contain any frequent $(k+1)$ -itemsets. Thus, such transactions can be labeled or
220 removed; hence, subsequent database scans for j -itemsets ($j > k$) will no longer consider such
221 transactions. One variation of the Apriori algorithm involves partitioning the data to find candidate
222 itemsets. The partitioning technique requires only two database scans to mine frequent itemsets
223 (Figure 11). It consists of two phases. In phase I, the algorithm partitions the transactional data of D
224 into n nonoverlapping partitions. If the minimum support threshold in D is set as min_sup , then the
225 minimum support count for a partition should be $min_sup \times \text{the number of transactions in the partition}$.
226 For each partition, all local frequent itemsets are found [17, 19].



227

228

Figure 11. Mining by the partitioning technique

229 Obtaining all local frequent k -itemsets requires only one database scan. All local frequent
230 itemsets can be treated as candidate itemsets for D . Therefore, collecting frequent itemsets from all
231 partitions enables generating the global candidate itemsets for D . In phase II, a second scan of D
232 is performed to evaluate the actual support count of each candidate to identify the global frequent
233 itemsets. This phase determines the partition size and the number of partitions so that each partition
234 can fit into main memory and thus be read only once in each phase [19].

235 Sampling can also be utilized to mine a subset of the given data. This technique involves
236 selecting a random sample S of the given database D and searching for frequent itemsets in S rather
237 than D . This method trades off certain degree of accuracy but increases efficiency. Because the
238 frequent itemsets are searched in S instead of D , some of the global frequent itemsets might be lost.
239 To reduce that possibility, a support threshold lower than the minimum support is adopted to
240 identify the frequent itemsets local to S (denoted L_S). The rest of the database is then used to
241 calculate the actual frequencies of each itemset in L_S to determine whether all the global frequent
242 itemsets are included in L_S . If L_S actually contains all the frequent itemsets in D , then only one scan of
243 D required. Otherwise, a second scan can be conducted to find any frequent itemsets that might have
244 been missed in the first scan. This sampling technique is particularly adequate when efficiency is
245 prioritized. Scholars have also proposed dynamic itemset counting as a method for enhancing the

246 Apriori algorithm. This approach adds candidate itemsets at various points during a scan and
247 partitions the database into blocks marked by start points. In contrast to standard Apriori, which
248 confirms new candidate itemsets only before each complete database scan, this variation can add
249 new candidate itemsets at any start point. This technique dynamically evaluates the support count of
250 all itemsets examined. If all subsets in an itemset are confirmed as frequent, then this itemset will be
251 considered as a new candidate itemset. Accordingly, this technique requires fewer scans than
252 standard Apriori does[14, 19].

253 Under various circumstances, the candidate itemsets generated by Apriori require an
254 examination method that will substantially reduce the size of the candidate itemsets, thereby
255 achieving favorable effectiveness of the Apriori algorithm. However, this algorithm must consider
256 the following two conditions:

257 (1) It may be required to produce a large number of candidate itemsets. For instance, if there are 10^4
258 frequent 1-itemsets, the Apriori algorithm will be required to generate more than 10^7 candidate
259 2-itemsets, during which the frequencies of all the candidate itemsets must be counted and
260 examined. In addition, to find frequent itemsets of length 100, such as $\{a_1, \dots, a_{100}\}$, a total of 2^{100}
261 $\approx 10^{30}$ candidate itemsets must be generated.

262 (2) It may be required to repeatedly scan the database and examine a large set of candidate itemsets
263 by pattern checking. Thus, a large amount of time is required to scan a database containing
264 abundant data.

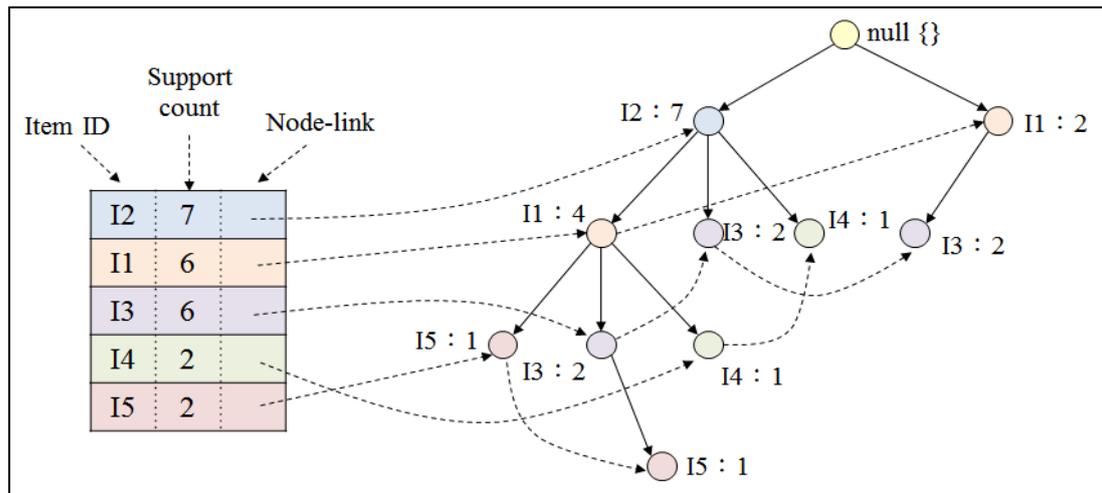
265 [14-17] have improved previous algorithms by proposing an algorithm that mines the complete
266 set of frequent itemsets without requiring generation of candidate itemsets; the improved algorithm
267 involves using frequent pattern growth (FPG) and frequent pattern tree (FP-tree) methods. Because
268 this algorithm (hereinafter referred to as the FPG algorithm) possesses a higher efficiency than
269 conventional algorithms possess, the present study examined it further.

270 The FPG algorithm employs a divide-and-conquer strategy, illustrated as follows: First, it
271 compresses the database containing frequent itemsets into a FP-tree, which retains the itemset
272 association information. Subsequently, it partitions the compressed database into a set of conditional
273 databases, each of which is a special projected database that is associated with one frequent item.
274 Then it mines each database separately. The same example data were used to illustrate the algorithm
275 proposed by [14-17].

276 The first database scan of this algorithm is identical to that of the Apriori algorithm, which
277 derives the set of frequent items (1-itemsets) and their support (regarding count and frequency). Let
278 the minimum support count be 2. The set of frequent items is ranked in the order of descending
279 support count and the resulting set or list is denoted by L , thereby obtaining $L=[I_2:7, I_1:6, I_3:6, I_4:2,$
280 $I_5:2]$. An FP-tree is subsequently constructed. The root of the tree is created and labeled with "null."
281 Perform a second scan on database D . The items in each transaction are then arranged in L order
282 (i.e., ranked in the order of descending support count), with a branch created for each individual
283 transaction.

284 The scan of the first transaction, " $T_{100}:I_1, I_2, I_5$," which contains three items (I_2, I_1 , and I_5 ,
285 sorted in L order), leads to the construction of the first branch of the tree with three nodes, $(I_2:1)$,
286 $(I_1:1)$, and $(I_5:1)$, where I_2 is linked as a child to the root, I_1 is linked to I_2 , and I_5 is linked to I_1 . The
287 second transaction, " $T_{200}:I_2, I_4$," which contains the items I_2 and I_4 in L order, leads to a branch
288 where I_2 is linked to the root and I_4 is linked to I_2 . However, this branch shares a common prefix, I_2 ,
289 with the existing path for T_{100} . Therefore, the count of node I_2 is increased by 1, and a new node,
290 $(I_4:1)$, is created; it is linked as a child to $(I_2:2)$. When adding a branch for a transaction, the count of
291 each node along a common prefix is generally increased by 1, and nodes for the items following the
292 prefix are created and linked accordingly.

293 To facilitate tree traversal, an item header table is constructed so that each item points to nodes
294 with the same ID in the tree through a chain of node-links. The tree obtained after scanning all the
295 transactions is shown in Figure 12, where all the associated node-links are illustrated. Therefore, the
296 problem of mining frequent patterns in databases is transformed into that of FP-tree mining.



297
298

Figure 12. FP-tree that contains compressed and frequent pattern information

Method:

1. The FP-tree is constructed in the following steps.

(a) Scan the transaction database D once. Collect the set of frequent items F and their supports. Sort F in support descending order as L , the list of frequent items.

(b) Create the root of an FP-tree, and label it as "null". For each transaction $Trans$ in D do the following.

Select and sort the frequent items in $Trans$ according to the order of L . Let the sorted frequent item list in $Trans$ be $[p|P]$, where p is the first element and P is the remaining list. Call $insert_tree([p|P], T)$, which is performed as follows. If T has a child N such that $N.item-name = p.item-name$, then increment N 's count by 1; else create a new node N , and let its count be 1, its parent link be linked to T , and its node-link to the nodes with the same item-name via the node-link structure. If P is nonempty, call $insert_tree(P, N)$ recursively.

2. Mining of an FP-tree is performed by called $FP_growth(FP_tree, null)$, which is implemented as follows.

procedure $FP_growth(Tree, \alpha)$

1: if $Tree$ contains a single path P then

2: for each combination (denoted as β) of the nodes in the path P

3: generate pattern $\beta \cup \alpha$ with support = minimum support of nodes in ;

4: else for each a_i in the header of $Tree$ {

5: generate pattern $\beta = a_i \cup \alpha$ with support = a_i . support;

6: construct β , s conditional pattern base and then β 's conditional

$FP_tree Tree_{\beta}$;

7: if $Tree_{\beta} \neq \emptyset$ then

8: call $FP_growth(Tree_{\beta}, \beta)$;

299
300
301

Figure 13. Summary of the FPG algorithm, which mines frequent itemsets without generating candidate itemsets

302 FP-tree mining is performed starting from each frequent length-1 pattern. Construct a
303 conditional pattern base (a sub-database that comprises the set of prefix paths in the FP-tree
304 co-occurring with the suffix pattern) for each pattern. Then, construct its conditional FP-tree, and
305 perform mining recursively on the tree. The pattern growth is attained by connecting the suffix
306 pattern with the frequent patterns generated from a conditional FP-tree.

307 Take I5 as an example. I5 occurs in two branches of the FP-tree displayed in Fig. 12. The paths
308 formed by these branches are $(I2\ I1\ I5:1)$ and $(I2\ I1\ I3\ I5:1)$. Therefore, the two corresponding prefix
309 paths for I5 are $(I2\ I1:1)$ and $(I2\ I1\ I3:1)$, which form its conditional pattern base. The I5-conditional
310 FP-tree contains only a single path, $(I2:2\ I1:2)$; I3 is included because its support count of 1 is less than

311 the minimum support count of 2. This single path generates all the combinations of frequent
 312 patterns, which are (I2, I5: 2), (I1, I5: 2), and (I2, I1, I5: 2). The same process repeats when using other
 313 items. The FPG algorithm is summarized in Figure 13.

314 In terms of mining long and short frequent patterns, the FPG algorithm exhibits higher
 315 effectiveness, flexibility, and efficiency than the Apriori algorithm does.

316 3. Data Mining Algorithms Applicable in Network Security

317 Data related to network security require special attention in certain respects. For instance, they
 318 may contain time packets so that the employed data mining algorithm may require adjustment
 319 accordingly. Therefore, data mining algorithms applicable to network security were analyzed in the
 320 present study to facilitate their subsequent application. Dain and Cunningham [1] proposed an
 321 algorithm that organizes alerts generated by heterogeneous IDSs into scenarios; this algorithm
 322 computes a new alert by estimating the probability that this new alert belongs to a given scenario
 323 and adds alerts to the most likely candidate scenario. In a publication of the IEEE, Ourston et al. [5]
 324 proposed a method that utilized HMMs to detect complex network attacks. The HMM is more
 325 effective in detecting such complex intrusion events than are decision trees and neural networks.

326 This section uses [10, 18] to illustrate the application of data mining in network security. These
 327 two studies proposed a sequential pattern mining technique that incorporated alert information
 328 associated with DDoS attacks and reduced the repetition and complexity of abundant alert
 329 information received by the IDS, thereby identifying the type of the DDoS attack. In addition, [10]
 330 employed a particular HMM to compare algorithms. This method comprises four algorithms.
 331 Algorithm 1 is the main mining loop and the other three algorithms are various extensions of
 332 Algorithm 1.

333 In Algorithm 1, four variables are first inputted. The first variable is set as alertSet, which
 334 denotes an alert set consisting of m raw alerts. The second to fourth variables are address prefix
 335 length n , time constraint t , and support threshold s , respectively. The output variables are alerts that
 336 have undergone correlation analysis and DDoS scenarios. Therefore, alertSet, address prefix length
 337 n , and time constraint t are introduced as parameters into genSequence, a function that generates
 338 sequential patterns. The output values produced from the calculation using genSequence are then
 339 used as the parameters of genSequence and introduced into Algorithm 3, where sequenceSet and the
 340 minimum support threshold s are inputted into function minePattern to identify the alert sample set
 341 with the longest sequential pattern, labelled lspSet. Subsequently, lspSet and n are introduced into
 342 the corrSPattern function to generate the desired DDoS scenarios. This entire process is completed
 343 using input parameters and three functions (i.e., the following algorithms to be discussed).

Algorithm 1 corrAlerts
Input: alertSet containing m raw alerts, address prefix
Length n , time constraint t , support threshold s
Output: correlated alerts, DDoS scenarios
 1: $sequenceSet \leftarrow genSequence(alertSet, t, n)$
 2: $lspSet \leftarrow minePatten(sequenceSet, s)$
 3: $scenarios \leftarrow corrSPattern(lspSet, n)$
 4: *output DDoS scenarios and the alerts after correlation*

344
 345

Figure 14. Algorithm 1

346 The match (IP1, IP2, 1) employed by the following algorithms is used as an example for
 347 clarification. Let the addresses of IP1 and IP2 be 135.13.216.191 and 135.13.16.189, respectively. A
 348 comparison of these addresses reveals that the third items of the two IP addresses are different.
 349 Therefore, this indicates that n is the length of the first two identical IP addresses, namely $n = 16$.

```

Algorithm 2 genSequence
Input: alertSet containing  $m$  alerts,  $t$ ,  $n$ 
Output: sequenceSet
1: sequence, seqArray  $\leftarrow \varnothing$ 
2:  $k = 0$ 
3: for  $all = alert_1$  to  $alert_m$  do
4:   if ( $all$  is not marked)
5:     for  $al_2 = all$  to  $alert_m$  do
6:       if (( $al_2$  is not marked)  $\wedge$ 
            $match(all.srcIP, al_2.srcIP, n) \wedge$ 
            $match(all.dstIP, al_2.dstIP, n)$ )
7:          $curAlert = al_2$ 
8:          $sequence \leftarrow sequence \cup al_2$ 
9:         mark  $al_2$ 
10:        for  $al_3 = al_2 + 1$  to  $alert_m$  do
11:          if (( $al_3$  is not marked)  $\wedge$ 
              ( $al_3.time - curAlert.time = t$ ))
12:             $flag = false$ 
13:            if ( $match(al_2.srcIP, al_3.srcIP, n) \wedge$ 
                 $match(al_2.dstIP, al_3.dstIP, n) \wedge$ 
                 $al_3.stage == curAlert.stage + 1$ )
14:               $flag = true$ 
15:            else if ( $curAlert.stage == 2 \wedge$ 
                     $al_3.stage == 3$ )
16:               $flag = true$ 
17:            end if
18:            if ( $flag$ )
19:               $curAlert = al_3$ 
20:               $sequence \leftarrow sequence \cup al_3$ 
21:              mark  $al_3$ 
22:            end if
23:          end for
24:           $seqArray \leftarrow seqArray \cup sequence$ 
25:           $sequence \leftarrow \varnothing$ 
26:        end if
27:         $sequenceSet[k++] \leftarrow seqArray$ 
28:         $seqArray \leftarrow \varnothing$ 
29:      end for
30:    output sequenceSet

```

Figure 15. Algorithm 2

350
351

352 Because real DDoS attacks at the final stage often contain packets generated from address
353 errors, adding alerts to a sequence facilitating counteracting this problem. The objective of
354 Algorithm 2 is to establish alert sequences by inputting variables m , t , and n to generate the desired
355 sequence. This algorithm consists of three loops. The first two loops define an alert group by
356 comparing IP addresses, namely comparing match (IP1, IP2, 1) with the host address. Then, the third
357 loop extracts alerts during various attack stages from the previously obtained alert group to
358 construct alert sequences.

359 According to lines 6 and 7 of the following algorithm, all alerts that have been sorted are
360 marked, starting from alert 2. If an alert is not marked, its srcIP and dstIP will be compared with all
361 other srcIP and dstIP values. After the comparison, the original sequence and the alert will be

362 redefined into a new sequence, and this alert will be marked. For alerts that have already been
 363 marked, alert time is used to determine whether they will be marked as false alerts, as suggested by
 364 lines 11 and 12. If the result of subtracting an alert time from the current alert time is \leq the inputted
 365 variable t , then this alert is marked as a false alert. If the alert time meets the aforesaid criterion, then
 366 whether the stage of this alert is higher than the current alert stage by one level will be checked, as
 367 indicated by lines 13 and 14. If the alert stage is lower than the current stage, the alert is marked as a
 368 false alert. The same rule applies to the rest of the alerts.

369 Algorithm 3 searches the sequential alert patterns from the alert groups created in algorithm 2.
 370 Alerts found in each pattern are regarded as correlated. The longest sequential pattern is reserved in
 371 the group because normal activities of a legitimate user may sometimes exhibit less frequent
 372 patterns that are incorrectly captured into alert sequences. Such alerts are removed from further
 373 analysis by this algorithm.

374 Algorithm 3 is as follows: Create a $lspSet$ variable and set it as an empty set. Create a variable sp
 375 and extract a portion of sequence to sp , as indicated in line 3. The sequence of sp is first used to find
 376 long alert sequence (labelled lsp). Alert sequences meeting this lsp criterion are defined as $srcIP$ and
 377 $dstIP$. Finally, the identified lsp variable is added into the alert group, which is redefined as a new
 378 group. The same process repeats until all $lspSet$ sets found by the algorithm have been outputted.

```

Algorithm 3 minePattern
Input: sequenceSet ( $p$  alert groups), support  $s$ 
Output: longest sequential alert pattern set  $lspSet$ 
1:  $lspSet \leftarrow \varnothing$ 
2: for  $i = 0$  to  $p$  do
3:    $sp \leftarrow SPAM(sequenceSet[i], s)$ 
4:   if ( $sp = \varnothing$ )
5:     correlate all alerts covered by  $sp$ 
6:      $lsp =$  the longest sequential alert pattern in  $sp$ 
7:      $lsp.srcIP \leftarrow$  the  $srcIP$  of the first alert in  $lsp$ 
8:      $lsp.dstIP \leftarrow$  the  $dstIP$  of the last alert in  $lsp$ 
9:      $lspSet \leftarrow lspSet \cup lsp$ 
10:  end if
11: end for
12: output  $lspSet$ 
  
```

379
 380

Figure 16. Algorithm 3

381 Algorithm 4 generates DDoS scenarios mainly using the obtained sequential alert pattern
 382 sample. The alert sample set with the longest sequential pattern is inputted into the algorithm and
 383 loops of sequential patterns are used to generate a scenario catalog. In Algorithm 4, the first
 384 $lspatternn$ is required to end the alert at stage2 or stage3 particularly when producing correlation
 385 data.

```

Algorithm 4 corrSPattern
Input: longest sequential pattern set lspSet (containing n
lspatterns), address prefix length l
Output: DDoS scenarios
1: corrlist, scenariolist ← ∅
2: for lsp1 = lspattern1 to lspatternn do
3:   if (lsp1 is not marked )
4:     currentlsp ← lsp1
5:     corrlist.addToTail(lsp1)
6:     marklsp1
7:     for lsp2 = lsp1 to lspatternn do
8:       if (( lsp2 is not maked) ^
           match(currentlsp.dstIP, lsp2.srcIP, n) ^
           currentlsp.lastAlert.stage = 2)
9:         currentlsp ← lsp2
10:        corrlist.addToTail(lsp2)
11:        mark lsp2
12:       end if
13:     end for
14:     scenario ← fuse the lspatterns in corrlist
15:     scenariolist = scenariolist ∪ scenario
16:     corrlist ← ∅
17:   end if
18: end for
19: output DDoS scenarios in scenariolist

```

386
387

Figure 17. Algorithm 4

388 The test results confirmed that sequential pattern mining algorithms are more reliable than
389 other types of algorithm because they do not result in mutual exclusivity of correlation data during a
390 DDoS attack and can increase the accuracy of multi-step intrusion detection. Applying HMMs to
391 compare algorithms through comprehensive and experimental methods verifies whether the
392 proposed algorithm exhibits favorable effects for identifying the sequence of DDoS attacks and
393 ensures effective reduction in alert complexity.

394 We launched attacks with three tools, Nmap, Pizzaviat DDoS, Nikto, and compared them with
395 used algorithms and unused algorithms. Table 1 shows the number of times the alarm was detected.
396 The method of unused algorithms will repeatedly detect alerts, resulting in a much larger number
397 of times. The method of used algorithms will greatly reduce the number of alarms for the same
398 event.

399

Table 1. The Number of alerts

	Tool	Used Algorithms	Unused Algorithms
1	Nmap	2	6
2	Pizzaviat DDoS	3	15
3	Nikto	5	45

400

401 4. Conclusions

402 The importance of data mining in network security lies in its ability to automatically extract
403 hidden but crucial information from the alert and log data of large systems, thereby providing a
404 helpful tool for information security personnel to make decisions. Therefore, data mining has been
405 gradually introduced to network security studies that aim to investigate methods for detecting

406 unknown attacks before hazards occur and to minimize potential hazards caused by such attacks.
407 Section 3 indicates that sequential patterns can be used to generate a desired DDoS scenario. If the
408 effectiveness of the algorithm can be enhanced, a future version of the algorithm may be able to
409 construct multi-step attack scenarios; such an improvement will be beneficial to future practice.
410 Accordingly, a number of key algorithms related to this topic were analyzed in this study to
411 facilitate analyzing massive data in the future.
412

413 **Author Contributions:** Kai Chain conceived and designed the experiments; performed the experiments and
414 analyzed the data and wrote the paper.

415

416 **Acknowledgments:** This research was supported by MOST 106-2221-E-145-002.

417 **Conflicts of Interest:** The author declares no conflict of interest.

418 References

- 419 1. Dain OM and Cunningham RK. Fusing a heterogeneous alert stream into scenarios. In: Proceedings of the
420 2001 ACM Workshop on Data Mining for Security Applications, 14th September 2001; pp. 1-13.
- 421 2. Eckmann ST, Vigna G and Kemmerer RA. Statl: An attack language for statebased intrusion detection. *J*
422 *Computer Security* **2002**; 10: 71-103.
- 423 3. Totel E, Vivinis B and Me L. A language driven IDS for event and alert correlation. In: SEC 2004; pp.
424 209-224.
- 425 4. Wang L, Ghorbani A and Li Y. Automatic multistep attack pattern discovering. *Int J Network Security*. **2010**;
426 10: 142-152.
- 427 5. Ourston D, Matzner S, Stump W and Hopkins B. Applications of Hidden Markov Models to detecting
428 multi-stage network attacks. In: Proceedings of the 36th Ann. Hawaii International Conference on System
429 Science 2003; pp. 334.
- 430 6. Somayaji A and Forrest S. Automated response using system-call delays. In: Proceedings of the 9th
431 USENIX Security Symp, Denver, Colorado, 14 -17 August 2000; pp. 185-198.
- 432 7. Lee W, Stolfo SJ and Chan PK. Learning patterns from unix process execution traces for intrusion
433 detection. In: Proceedings of the AAAI97 workshop on AI Approaches to Fraud Detection and Risk
434 Management, AAAI Press, 1997; pp. 50-56.
- 435 8. Warrender C, Forrest S and Pearlmutter BA. Detecting intrusions using system calls: Alternative data
436 models. In: Proceedings of the IEEE Symp. on Security and Privacy, 1999; pp. 133-145.
- 437 9. Forrest S, Perelson AS, Allen L and Cherukuri R. Self-nonsel self discrimination in a computer. In: SP'94: Proc.
438 of the 1994 IEEE Computer Society, 1994; pp. 202.
- 439 10. Xiang G, Dong X and Yu G. Correlating Alerts with a Data Mining Based Approach. In: Proceedings of
440 the IEEE International Conference on e-Technology, e-Commerce and Service, 29 March-1 April 2005.
- 441 11. AL-mamory SO and Zhang HL. A survey on IDS Alerts processing techniques. In: Proceedings of the 6th
442 WSEAS International Conference on Information Security and Privacy, Tenerife, Spain, 14-16 December
443 2007; pp. 69-78.
- 444 12. Cheng B, Liao G and Huang C. A novel probabilistic matching algorithm for multi stage attack forecasts.
445 *IEEE J on Selected Areas in Comm.* **2011**; 29: 1438-1448.
- 446 13. Agrawal R, Imielinski T and Swami AN. Mining association rules between sets of items in large databases.
447 In: Proceedings of the 1993 ACM SIGMOD Intl. Conference on Management of Data, 1993; pp. 207-216.
- 448 14. Han J, Pei J and Yin Y. Mining frequent patterns without candidate generation. In: Proceedings of the
449 2000 ACM SIGMOD international conference on Management of data, 2000; pp. 1-12.
- 450 15. Leung CKS, Mateo MAF and Brajczuk DA. A tree-based approach for frequent pattern mining from
451 uncertain data. In: Proceedings of the Pacific-Asia Conference on Knowledge Discovery and Data Mining,
452 2008; pp. 653-661.
- 453 16. Zaki MJ, Li F and Yi K. Finding frequent items in probabilistic data. In: Proceedings of the 2008 ACM
454 SIGMOD international conference on Management of data, 2008; pp. 819-832.
- 455 17. Han J and Chang KCC. Data mining for web intelligence. *Computer*. **2002**; 35: 64-70.

- 456 18. Stroeh K, Madeira ERM and Goldenstein SK. An approach to the correlation of security events based on
457 machine learning techniques. *J Int Services and Applications*. **2013**; 4: 1-16.
458 19. Han J, Pei J and Kamber M. *Data Mining Concepts and Techniques*. **2011**.