

Detection and Defense of Network Virus using Data Mining Technology

Zhijun Li (✉ ju42781@163.com)

Hebei Normal University for Nationalities

Xuedong Jiang

Hebei Normal University for Nationalities

Research

Keywords: Data mining, network virus, application programming interface call, feature selection, virus detection

Posted Date: May 18th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-508107/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Detection and defense of network virus using data mining technology

2 Zhijun Li^{1*}, Xuedong Jiang²

3 ¹Integrated Business Department of Assets and Laboratory Management Center, Hebei Normal
4 University for Nationalities, Chengde, Hebei 067000, China

5 ²Information Center, Hebei Normal University for Nationalities, Chengde, Hebei 067000, China

6 Corresponding address: Assets and Laboratory Management Center, Hebei Normal University for
7 Nationalities, Chengde, Hebei 067000, China

8 Email: ju42781@163.com

9 **Abstract:**

10 The spread of network viruses has posed a serious threat to the security of the network; therefore, it
11 is necessary to detect and defend them effectively. This paper used Debug application programming
12 interface (API) technology to obtain the features of API calls as viruses, filtered API calls according
13 to information entropy, and finally used the support vector machine (SVM) model for virus
14 detection. The experimental results showed that when the number of API was 1200, the algorithm
15 had the best virus detection performance, with an average true positive rate (TPR) of 95.2%, a false
16 positive rate (FPR) of 3.31%, and an overall accuracy of 95.42%; compared with the K-means
17 algorithm and Naive Bayes algorithm, the SVM algorithm had the best performance. The results
18 show that the proposed method is effective in virus detection and defense and can be further
19 promoted and applied in practice.

20 Keywords: Data mining, network virus, application programming interface call, feature selection,
21 virus detection

22

23

24 **1. Introduction**

25 The rapid development of the network has brought great conveniences to people's study, life, work,
26 and entertainment (Gao et al., 2021), but at the same time, the security problem of the network has
27 become increasingly prominent (Ma, 2017). Network virus refers to a group of codes that can
28 destroy the function or data of computers and also has the ability of self-replication, i.e., it can
29 spread quickly among the networks to affect the normal use of computers. Once the virus code is
30 executed in computers, it will be inserted into various programs to multiply itself. If it is not handled
31 in time, these viruses will spread to other uninfected computers through all platforms on the network,
32 making more computers lose the ability to work normally and even causing the paralysis of the
33 network system. In addition, the viruses have strong concealment, sometimes can not be detected
34 by anti-virus software. Before triggering, the viruses will lurk in the program; once triggered, they
35 will wantonly destroy the network. The increasingly rampant network virus severely challenges
36 computer security (Zuo, 2018). The high-speed spread of information and the rampant virus bring
37 a huge threat to people's information security. Therefore, how to realize the detection and defense
38 of network viruses and reduce the damage of viruses to computers to ensure the normal work of the
39 network has been widely concerned by researchers. With the emergence and development of data
40 mining technology, it plays an increasingly important role in network virus defense. Binh et al.
41 (2016) divided the control flow graph (CFG) extracted from a binary executable into specific
42 subgraphs to extract malicious behaviors. The method could reduce the complexity of verification
43 by incremental verification strategy. The experiment showed that the method could significantly
44 improve the performance. Nguyen et al. (2015) proposed a maximum entropy-based method. During
45 training, the features in the virus database were extracted to establish the upper model. During
46 detection, the upper model was used to identify viruses based on the corresponding features of a

47 checked file. Saudi et al. (2015) designed an efficient Trojan horse detection model based on the
48 sequential minimum optimization (SMO) algorithm and found through experiments that the method
49 had a 98.2% correct rate and a 1.7% error rate. Jerbi et al. (2020) evolved a group of application
50 programming interface (API) calls combined with the artificial malware patterns and based on the
51 genetic algorithm and found through the test on different data sets that the method could
52 significantly improve the detection rate. In this paper, based on the support vector machine (SVM),
53 the detection and defense methods of network viruses were studied, and its performance was tested
54 on the data set to verify the possibility of the method in practical application. This work makes some
55 contributions to the further development of network security.

56

57 **2. The SVM-based detection and defense method**

58 **2.1 Network virus feature screening**

59 In the Windows system, PE file is the standard format of executable file (Belaoued and
60 Mazouzi, 2018). Therefore, the virus is written according to the format of the PE file in order to
61 spread better. In the operation process, Windows calls various service functions to realize some
62 function, and the functions are called API functions. To read a file on disk, first of all, the Readfile
63 function in Kernel32.dll is called; if the permission meets the requirement, the NtReadFile function
64 in Ntdll.dll is called; then, CPU is switched to the operating system kernel mode through the
65 KiSystemService function; finally, the NtReadFile kernel function in ntoskrnl.exe is called. After
66 file reading, the NtReadFile function in ntoskrnl.exe will return in the same way. This is an API call.

67 All programs on the Windows platform need API calls to implement functions, including the
68 degree of viruses. Therefore, network viruses can be detected by capturing API calls, such as
69 Windows debug API (Sigalov et al., 2017) and APIHOOK technology (Case et al., 2019). However,
70 APIHOOK technology can only be implemented when every API prototype is known, and not every
71 function can be effectively captured. Therefore, this study used Windows Debug API technology to
72 obtain API calls.

73 Debug API technology can load a program or bind itself to the program to facilitate debugging.
74 If there are debugging-related events, the debugger will be triggered, i.e., the monitoring process
75 will be activated. The steps of debug API technology can be described as follows: (1) the monitor
76 program is started after inputting samples; (1) the export table of the library file called by the
77 debugging program is analyzed; (3) the breakpoint is set at the entry address of the function; (4) if
78 an interrupt occurs at the breakpoint, the debug process information will be obtained; (5) the
79 debugging process is cycled until the end, and a behavior report is obtained.

80 In the collected API calls, a large part of the content does not play a great role in virus detection
81 but will cause information redundancy and increase the calculation of subsequent virus detection.
82 Therefore, it is necessary to screen the collected features and select the features with a high
83 differentiation degree and a high important degree. This paper used the method of information
84 entropy. Regarding whether a program calls related API as information and whether a program is a

85 virus or not as an event, the contribution of information to an event was judged through the
 86 calculation of information entropy, so as to determine the importance of the feature.

87 The entropy of a random variable X is set as: $H(X) = -\sum_{i=1}^k p_i(X = v_i) \log_2 p_i(X = v_i)$, where X
 88 has k values and p_i represents the probability that the value of X is v_i . It is assumed that under
 89 the condition that a random variable Y has been known, the conditional entropy of X is:
 90 $H(X|Y) = -\sum_{j=1}^m p_j(Y = v_j) \times H(X|Y = v_j)$, where $H(X|Y)$ refers to the uncertainty of X . Then,
 91 under the condition that Y has been known, the information added value of X , i.e., the information
 92 grain, can be written as: $IG(X|Y) = H(X|Y) - H(X)$; the greater the value is, the larger the function
 93 of the API call is. Features with large information gains are reserved to establish a virus feature set.
 94 Every feature vector is represented by a boolean vector. If the feature appears in the program, it is
 95 assigned as 1; otherwise, it is assigned as 0. A program is set as B_i , and a feature attribute is set as
 96 A_i . The feature set can be represented by matrix $M_{A \times B}$:

$$97 \quad M_{A \times B} = \begin{matrix} & A_1 & A_2 & \Lambda & A_n \\ \begin{matrix} B_1 \\ B_2 \\ \vdots \\ B_m \end{matrix} & \begin{pmatrix} b_{11} & b_{12} & \Lambda & b_{1n} \\ b_{21} & b_{22} & \Lambda & b_{2n} \\ \text{M} & \text{M} & \Lambda & \text{M} \\ b_{m1} & b_{m2} & \Lambda & b_{mn} \end{pmatrix} \end{matrix}, (1)$$

98 where $b_{ij} = \begin{cases} 1, \text{Program } B_i \text{ has feature } A_j \\ 0, \text{Program } B_i \text{ does not have feature } A_j \end{cases}$.

99

100 2.2 The SVM-based detection mode

101 SVM is a method based on structural risk minimization (Shi et al., 2015). It uses quadratic
 102 programming to find the optimal solution, avoiding the local minimization problem of neural
 103 networks, and solves the dimension problem through the kernel function. Therefore, it has been
 104 widely used in many fields. Virus detection can be regarded as a classification problem. Therefore,
 105 this paper uses SVM to detect viruses.

106 It is assumed that a sample set is (x_i, y_i) , where $i = 1, 2, \Lambda, n$, $y \in \{+1, -1\}$, x_i represents the
 107 feature of the program, and y_i represents the class that the program belongs to. The corresponding
 108 classification surface can be written as $f(x) = wx + b$, i.e., satisfying $y_i(wx_i + b) \geq 0$, where w

109 stands for a weight and b stands for a bias. The class interval is $\frac{2}{\|w\|}$, the optimal classification
 110 surface can be written as: $\min \varphi(w) = \frac{\|w\|^2}{2}$, and the constraint condition is: $f(x) \geq 1, i = 1, 2, \dots, n$. Then,
 111 the decision function is obtained:

$$112 \quad g(x) = \text{sgn}(wx_i + b). \quad (2)$$

113 In order to solve equation (2), Lagrange multiplier a is introduced to transform it into a dual
 114 problem:

$$115 \quad Q(a) = \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j=1}^n a_i a_j y_i y_j (x_i \cdot x_j). \quad (3)$$

116 The constraint condition is:

$$117 \quad \begin{cases} \sum_{i=1}^n y_i a_i = 0 \\ a_i \geq 0 \end{cases}. \quad (4)$$

118 Finally, the SVM model for virus detection can be written as:

$$119 \quad f(x) = \text{sgn} \left(\sum_{i=1}^n a_i^* y_i k(x_i \cdot x) + b^* \right)$$

120 where k is a kernel function. The common kernel functions are:

$$121 \quad (1) \text{ linear kernel function: } K(x_i, x_j) = x_i^T x_j;$$

$$122 \quad (2) \text{ polynomial kernel function: } K(x_i, x_j) = (\gamma x_i^T x_j + r)^d, \gamma > 0;$$

$$123 \quad (3) \text{ S-shaped kernel function: } K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r);$$

$$124 \quad (4) \text{ radial basis function (RBF) kernel function: } K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \gamma > 0.$$

125 γ , r , and d are all nuclear parameters.

126 In selecting the kernel function, because of the large number of samples and features required
 127 in virus detection, the RBF kernel function can avoid the dimension disaster better. Therefore, this
 128 paper chose the RBF kernel function to build the SVM model.

129

130

131 3. Experimental analysis

132 Virus samples were collected from some well-known forums and laboratories, and normal
133 samples were also selected from the Windows system. Finally, 1256 samples were obtained,
134 including 692 virus samples and 564 normal samples. Then, features were extracted by debug API
135 technology, and the obtained features were stored in the MySQL database.

136 The purpose of the experiment was to analyze the virus detection and defense ability of the
137 model designed in this paper. The judgment result was “yes” or “no”, as shown in Table 1.

138 Table 1 Classification results

139

Results	Normal	Virus
Normal file	TP	FN
Virus files	FP	TN

140

141 The evaluation indexes of the model included:

142 true-positive rate: $TPR = \frac{TP}{TP + FN}$;

143 false-positive rate: $FPR = \frac{FP}{FP + TN}$;

144 total accuracy: $Total_Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$.

145 Firstly, the influence of the number of features on the performance of virus detection was
146 analyzed. After extracting API calls by using debug API technology, a total of 2167 API calls were
147 obtained. Through feature screening, different numbers of APIs were selected for virus detection.
148 The overall accuracy of the model is shown in Table 2.

149 Table 2 Changes of overall accuracy

150

API quantity/n	Overall accuracy/%
100	80.79
300	86.78
600	92.34
900	93.62

1200	95.27
1500	91.88
1800	91.36
2100	90.27

151

152 It was seen from Table 2 that with the increase of the number of APIs, the overall accuracy of
 153 the model rose rapidly, from about 80% to about 90%. When the number of APIs reached 1200, the
 154 overall accuracy was the highest, 95.27%. Then, with the continuous growth of the number of APIs,
 155 the accuracy of the algorithm began to decline. When the number of APIs reached 2100, the overall
 156 accuracy of the algorithm was 90.27%, which decreased by 5%. Therefore, the number of selected
 157 APIs was set as 1200.

158 The method of ten-fold cross-validation was adopted to determine the performance of the SVM
 159 model, as shown in Table 3.

160 Table 3 Virus detection results of the SVM model

161

Verification times	TPR/%	FPR/%	Overall accuracy/%
1	95.64	3.27	96.12
2	94.87	4.12	95.65
3	96.33	2.36	96.11
4	94.39	3.25	94.36
5	95.85	2.12	95.18
6	96.21	3.51	94.82
7	92.58	3.64	95.88
8	93.69	3.98	96.09
9	97.08	4.26	94.39
10	95.39	2.61	95.64
Average	95.20	3.31	95.42

162

163 It was seen from Table 3 that the average TPR and FPR of the SVM model were 95.2% and
 164 3.31% respectively, and the overall accuracy rate reached 95.42%, which indicated that the feature
 165 screening and detection model designed in this paper had a high accuracy rate in the detection and
 166 defense of network viruses. In order to further prove the effectiveness of the model, it was compared
 167 with the K-means algorithm and the Naive Bayes algorithm. The experimental results are shown in
 168 Figure 1.

169 It was seen from Figure 1 that the TPR values of the three algorithms were 93.22%, 94.26%, and
170 95.2%, respectively, i.e., the TPR of the SVM model was 1.98% and 0.94% higher than that of the
171 former two algorithms; the FPR values of the three algorithms was 6.32%, 4.67%, and 95.2%,
172 respectively, i.e., the FPR of the SVM model was 3.01% and 1.36% lower than that of the former
173 two algorithms. The overall accuracy of the three algorithms were 88.67%, 92.16%, and 95.42%,
174 respectively. The SVM model had the highest accuracy, which was 6.75% and 3.26% higher than
175 the former two algorithms, which verified the reliability of the SVM model.

176

177 **4. Discussion**

178 The emergence of more and more new and unknown viruses not only will have a huge impact
179 on society but also may lead to irreparable economic losses (Wang, 2017); therefore, virus detection
180 and defense has become a key and difficult problem in the field of network security (Tang et al.,
181 2020), and more new and effective methods are urgently needed. At present, the commonly used
182 methods include the behavior detection method (Choi et al., 2019), the heuristic scanning method
183 (Zakeri et al., 2015), the intelligent detection method, etc. The intelligent detection method refers to
184 applying intelligent algorithms such as data mining and machine learning to realize the detection
185 and defense of viruses. Many data mining algorithms have been successfully applied in virus
186 detection, such as decision trees and neural networks. It can be found that data mining has great
187 potential and prospects in virus detection and defense.

188 In this study, API calls were extracted by debug API technology as virus features and screened
189 based on information entropy, and the SVM model was used as a detection model to study the
190 network virus detection. First of all, the number of extracted features had an impact on the virus
191 detection performance. When the number of features was too large, the redundant information
192 contained will not be conducive to virus detection. According to Table 2, when the number of APIs
193 used was 1200, the algorithm had the best performance. In the ten-fold cross test, the SVM model
194 showed 95.2% TPR, 3.31% FPR, and 95.42% overall accuracy, indicating that the SVM model had
195 an excellent performance in virus detection and could effectively distinguish normal and virus
196 programs. Finally, compared with the other algorithms, the SVM model had good reliability in virus
197 detection.

198 In this paper, although some results have been obtained from the research of network virus
199 detection and defense methods, there are still some shortcomings. In future research, the author will:

200 (1) analyze more data mining methods;

201 (2) further improve the performance of virus detection;

202 (3) further subdivide the network viruses and compare Trojan horse viruses with worm viruses.

203

204 **5. Conclusion**

205 Based on the SVM model, this paper analyzed the method of network virus detection and defense
206 and collected virus programs and normal programs as data sets for experiments. The results showed
207 that the performance was the best when 1200 API calls were selected for virus detection. The
208 average TPR of the SVM model was 95.2%, the average FPR was 3.31%, and the overall accuracy
209 was 95.42%. Compared with the K-means algorithm and Bayesian algorithm, the SVM model had
210 a better performance; therefore, the SVM model can be further promoted and applied in practice.

211

212

213 List of abbreviations

214 API: application programming interface

215 SVM: support vector machine (SVM)

216 FPR: false positive rate

217 TPR: true positive rate

218 CFG: control flow graph

219 SMO: sequential minimum optimization

220

221 Declaration Statements

222 Availability of data and material

223 The datasets used and/or analysed during the current study are available from the corresponding
224 author on reasonable request.

225

226 Competing interests

227 The authors declare that they have no competing interests.

228

229 Funding

230 Not applicable.

231

232 Authors' contributions

233 ZJL conceived the idea. ZJL and XDJ designed the experiment and carried out the experiment. XDJ
234 analyzed the data. XDJ wrote the first draft. ZJL revised the manuscript.

235

236

237 Acknowledgements

238 Not applicable.

239

240

241 **References**

242 A Case, MM Jalalzai, M Firoz-Ul-Amin, RD Maggio, A Ali-Gombe, MX Sun, GG Richard,
243 HookTracer: A System for Automated and Accessible API Hooks Analysis. Digit. Invest. 29, S104-
244 S112 (2019). doi:10.1016/j.diin.2019.04.011.

245 C Choi, C Esposito, M Lee, J Choi, Metamorphic malicious code behavior detection using
246 probabilistic inference methods. Cogn. Syst. Res. 56(AUG.):142-150 (2019).
247 doi:10.1016/j.cogsys.2019.03.007.

248 C Zuo, Defense of Computer Network Viruses Based on Data Mining Technology. Int. J. Netw.
249 Secur. 20(4):805-810 (2018). doi:18163548-201807-201807180001-201807180001-805-810.

250 DA Sigalov, AV Razdobarov, AA Petukhov, Detecting DOM-based XSS vulnerabilities using
251 debug API of the modern web-browser. Prikl. Diskr. Mat. 35, 63 - 75 (2017).
252 doi:10.17223/20710410/35/6. doi:10.17223/20710410/35/6.

253 J Shi, WJ Lee, Y Liu, Y Yang, P Wang, Forecasting Power Output of Photovoltaic Systems Based
254 on Weather Classification and Support Vector Machines. IEEE T. Ind. Appl. 48(3):1064-1069
255 (2015). doi:10.1109/TIA.2012.2190816.

256 M Belaoued, S Mazouzi, A Real-Time PE-Malware Detection System Based on CHI-Square Test
257 and PE-File Features. IFIP Adv. Inform. Commun. Tech., 456, 416-425 (2018). doi:10.1007/978-
258 3-319-19578-0_34.

259 M Jerbi, ZC Dagdia, S Bechikh, M Makhlof, LB Said, On the use of artificial malicious patterns
260 for android malware detection. Comput. Secur. 92(May):101743.1-101743.22 (2020).
261 doi:10.1016/j.cose.2020.101743.

262 M Zakeri, FF Daneshgar, M Abbaspour, A static heuristic approach to detecting malware targets.
263 Secur. Commun. Netw.8(17):3015-3027 (2015). doi:10.1002/sec.1228.

264 MM Saudi, AM Abuzaid, BM Taib, ZH Abdullah, Designing a new model for Trojan horse
265 detection using sequential minimal optimization. Lect. Notes Electr. Eng. 315, 739-746 (2015).
266 doi:10.1007/978-3-319-07674-4_69.

267 NT Binh, QT Tho, HM Ngoc, NM Hai, Incremental verification of ω -regions on binary control
268 flow graph for computer virus detection. 2016 3rd National Foundation for Science and Technology
269 Development Conf. on Information and Computer Science (NICS), 68-73 (2016).
270 doi:10.1109/NICS.2016.7725670.

271 NT Nguyen, PV Huong, BC Le, DT Le, TH Van Le, A New Method of Virus Detection Based on
272 Maximum Entropy Model. Adv. Intell. Syst. Comput. 358, 151-161 (2015). doi:10.1007/978-3-319-
273 17996-4_14.

274 W Tang, YJ Liu, YL Chen, YX Yang, XX Niu, SLBRS: Network Virus Propagation Model based
275 on Safety Entropy. Appl. Soft Comput. 97, 106784 (2020). doi:10.1016/j.asoc.2020.106784.

276 X Gao, L Liu, X Zhu, Research on the Main Threat and Prevention Technology of Computer
277 Network Security. IOP Conf. Ser. Earth Environ. Sci. 632(5):052065 (7pp) (2021).
278 doi:10.1088/1755-1315/632/5/052065.

279 XJ Ma, Research and Implementation of Computer Data Security Management System. Proc. Eng.
280 174(Complete):1371-1379 (2017). doi:10.1016/j.proeng.2017.01.290.

281 Y Wang, Food Information Management and Security Strategy of Computer Network. Adv. J. Food
282 Sci. Tech.11(12):792-794 (2017). doi:10.19026/ajfst.11.2793.

283

284

285

286 Figure 1 Performance comparison of different algorithms

287

Figures

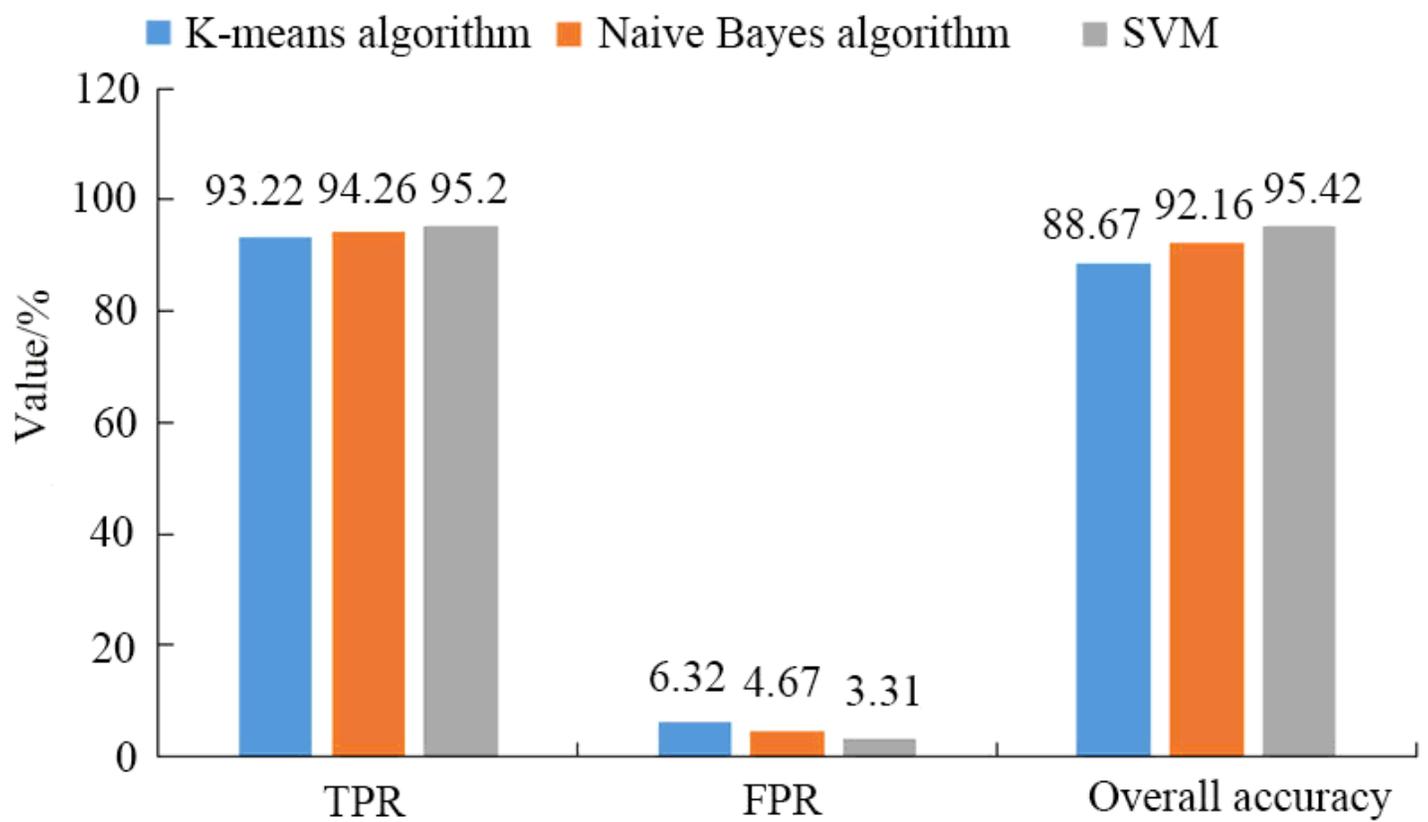


Figure 1

Performance comparison of different algorithms