

基于集对理论的矢量地图水印算法评价模型

杨 猛

(江苏自动化研究所 连云港 222061)

摘要 随着地理信息系统的发展,矢量地图在生活的各个领域已经得到广泛的应用。作为一种数字化数据,矢量地图具有易被非法复制、篡改、传播等特点,从而导致数字矢量地图的版权保护问题日益严峻。目前学者们提出了众多针对矢量地图版权保护的数字水印算法,但由于鲁棒性标准各不相同,因此很难得到一个比较公正的评价。通过分析现有的矢量地图水印算法,总结矢量地图水印攻击类型,结合集对理论,提出一种比较通用的适用于矢量地图水印算法的评价模型。

关键词 矢量地图,数字水印,鲁棒性,攻击分类,集对理论

中图法分类号 TP309 **文献标识码** A

Evaluation Model of Vector Map Watermarking Algorithms Based on Set Pair Analysis

YANG Meng

(Jiangsu Automation Research Institute, Lianyungang 222061, China)

Abstract With the development of geographic information system(GIS), vector map has been widely used in all areas of life. As a kind of digital data, vector map risks of illegal copying, tampering, communication, etc. It makes the copyright protection of vector maps an increasingly serious problem. Currently, a number of scholars have proposed kinds of watermarking algorithms. However, without uniform standard, it is very difficult to give a more fair assessment of the watermarking algorithms. This paper proposed a more general evaluation of vector map watermarking model by analyzing the existing vector map watermarking algorithm, summarizing vector map watermarking attack classification and combining theory of the set pair analysis.

Keywords Vector map, Digital watermarking, Robustness, Classification of attack, Theory of set pair analysis

随着地理信息系统的发展,矢量地图已被广泛应用于生活的各个领域,如城市地图、汽车导航及军事指挥等。由于数字产品便于编辑修改,数据很容易被非法篡改、复制、传播,因此数据的版权问题也越来越突出。为有效保护矢量地图数据的版权,许多关于矢量地图的水印算法被提出,但这些水印算法都只是对某些攻击具有较好的鲁棒性,不具有很好的实用价值。

矢量地图数字水印攻击技术是指针对矢量地图数字水印系统,按照矢量地图特点进行各种攻击,以检验其鲁棒性。矢量地图水印攻击分类提供了一系列的攻击方法,试图让水印提取者无法提取出水印,以达到减弱或消除水印的目的。集对理论^[1-2]作为一种新型的处理模糊和不确定性知识的数学工具,能够有效分析和处理不精确、不一致、不完整等各种不确定信息。本文将集对理论和矢量地图水印攻击分类相结合,提出一种合理的水印算法评价模型。

1 矢量地图数据的特点

矢量地图数据由几何数据和属性数据组成,空间数据包括几何数据和关系数据。几何数据是描述空间实体的空间特征数据,如点的坐标。关系数据是描述空间实体之间关系的数据,如空间实体的相邻关系等,主要是指拓扑关系。属性数

据是描述空间实体的属性特征数据,例如类型、名称、时间等,其中也包括描述空间特征的数据。

根据被描述实体的形态特征,矢量地图数据表示为点、线、多边形等图形。点图由在某种特定坐标系下的坐标点组成,线、多边形则是由一串有序的坐标点组成。所有地图对象都是由许多有某种特定联系的点构成,矢量地图中的空间定位数据实际上是由某种坐标系下的点所构成的。由于属性数据是用来描述地图对象的属性,不能轻易修改,因此水印一般都被冗余地嵌入到几何数据中。

2 矢量地图水印嵌入要求及算法分析

2.1 矢量地图数字水印嵌入要求

数字水印的目的是保证嵌入载体的数据是在任何操作下都能够被正确地提取出来。依照这一准则,水印应该主要满足以下要求。

(1)不可感知性:水印嵌入前、后,地图数据的变化不能引起观测者的觉察或测量误差,即水印的嵌入不影响该矢量地图的视觉效果或者不影响矢量地图中的元素。

(2)鲁棒性:在矢量地图的使用过程中经常遇到图元增减、旋转或非法使用者的破坏等操作,这就要求嵌入的水印具有较强的抵抗这些操作的能力。

(3)容量:嵌入的水印信息要能够唯一确定矢量地图的版权信息,这就要求水印算法能够嵌入足够的水印信息。

2.2 矢量地图水印嵌入算法分析

2.2.1 基于空域的水印算法

基于空域的水印嵌入算法在地图数据精度允许的范围内直接对地图数据的坐标值进行修改,它具有能够嵌入的水印信息量大、算法简单、不可感知性好等特点。文献[3]通过调整多边形线段的角度参数实现水印信息的嵌入,不能抵抗数据的更新操作。文献[4]提出基于统计量检测的水印算法,该算法对于地图简化及多数几何变换具有较好的鲁棒性,但在网络的交界处有可能产生不自然的形状扰动,而且水印的嵌入效率不高。文献[5-7]将地图进行分块嵌入,文中算法可以抵抗数据简化和小幅度的随机噪声攻击。文献[8]先将数据按照空间关系及位置进行排序,将水印信息嵌入在预处理后的数据中,文中算法对乱序攻击具有较强的鲁棒性,但对数据更新的鲁棒性不强。文献[9]利用道格拉斯-普克算法提取矢量地图曲线的特征点,在容错值范围内嵌入水印信息。

2.2.1 基于频域的水印算法

基于频域的水印算法首先对地图数据进行数学变换,在变换域上嵌入水印信息,然后经反变换输出。它具有鲁棒性强等特点,但算法相对复杂,并且对原始数据的修改较大。基于离散傅里叶变换(Discrete Fourier Transform, DFT)的矢量地图数字水印算法^[10-13]首先提取出图形的定位点(或特征点)用于构成一个实数序列,对实数序列进行 DFT 变换,利用傅里叶描述子的几何变换的不变性把水印信息嵌入在傅里叶系数的幅值或相位上。基于傅里叶变换的数字水印可以有效地抵抗图形的平移、缩放、旋转等几何变换。傅里叶变换是一种全局变换,局部很小的改动就会引起几乎全部傅里叶系数的变化,它对局部修改不具有鲁棒性。基于 DCT 的数字水印方案^[14]对地图的影响很大。基于小波变换^[15-17]的矢量地图数字水印算法可以很好地反映局部性质,但是小波变换不能很好地抵抗平移操作。

3 矢量地图水印攻击的分类

为了给矢量地图水印算法提供公平的评价平台,矢量地图数字水印基准采用各种攻击来进行测试,用全面综合的攻击分类以及科学的评估标准来促使人们设计出更具鲁棒性的水印算法。数字水印攻击分析^[18]通过采用各种攻击来检验数字水印系统的鲁棒性,根据攻击测试结果分析其易受攻击的原因并改进水印算法,攻击分析的目的是使水印检测者无法检测到水印的存在。

矢量地图数字水印攻击的基本前提是在地图坐标所能允许的精度范围之内进行攻击,否则攻击也就失去了其应有的价值。矢量地图水印算法的优劣主要通过水印鲁棒性来评价,故水印攻击应从水印的鲁棒性角度出发。从鲁棒性角度考虑,攻击主要包括有意攻击和无意攻击^[19-20]。有意攻击者的目标是去除、减弱甚至检测提取到数字水印;无意攻击的目标不是移除数字水印,而是指带有水印的地图在通用的操作中不可避免地会遇到如图元增加、删除、修改等操作,致使数字水印无法被检测。分类的整体情况如图 1 所示。

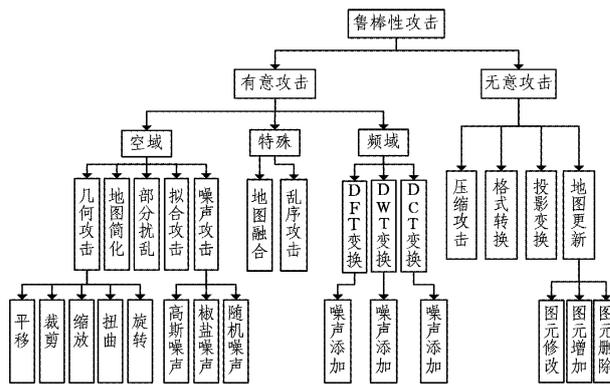


图 1 矢量地图水印攻击的分类

3.1 有意攻击

矢量地图的有意攻击可分为基于空域的有意攻击、基于频域的有意攻击以及基于矢量地图特点的特殊攻击。基于空域的有意攻击包括几何攻击、地图部分扰乱、地图简化及基于空域的各种噪声添加等;基于频域的攻击主要包括对地图进行某种数学变换之后在变换域中添加各种噪声的攻击;基于矢量地图特点的特殊攻击包括地图融合、数据乱序等。

几何攻击包括裁剪、缩放、平移、旋转等基本的几何操作;地图部分扰乱主要是对空域中的局部坐标进行修改;地图简化是对地图中一些特征不明显的图元进行简化;空域噪声添加主要是指直接在地图空域上添加高斯噪声、随机噪声、椒盐噪声等。

基于频域的攻击主要是对矢量地图数据进行一些特殊变换,如 DCT(离散余弦变换)、DWT(小波变换)、DFT(傅里叶变换)等,并在变换后的数据域中添加各种噪声(如高斯噪声等),以达到去除频域中添加的水印的目的。

地图融合是将多幅矢量地图融合到一幅矢量地图中的操作,由于各层的图元混合到了一层,因此水印提取者很难找到水印同步嵌入位置;矢量地图数据的显示主要是依靠数据的拓扑关系而非数据的存储位置,但是许多水印的添加是基于存储位置的,乱序操作主要是对数据的存储位置进行置乱操作,使得水印提取者很难找到正确的起始嵌入位置。

3.2 无意攻击

结合矢量地图的特点,无意攻击包括地图更新和压缩、格式转换及投影变换等操作。由于地图实体经常会发生变化(如地图中某个地方新添加了一个建筑物,某个地方新建了一条路等),因此不可避免地要对地图进行更新操作,如图元的添加、修改、删除等。在数据的传输过程中,为了能够以更小的空间传输更多的数据,需要对地图按照一定的方式进行有损压缩操作,这就在一定程度上对原数据进行了简化。格式转换是将地图按照一定的格式转换之后再转换为原有地图的格式,但在格式转换过程中可能会丢失一些数据信息,进而对原地图数据造成影响。投影变换是将地图按照所需要的坐标标准进行投影以达到适应不同应用效果的目的,但造成了数据的损失。

4 集对分析

4.1 集对分析理论

集对分析(Set Pair Analysis, SPA)理论是一种新型的处

理模糊和不确定知识的数学工具。集对分析把具有一定联系的两个集合视为一个集对,如某一产品质量与该产品的质量标即可被视为一个集对。集对分析是指在特定背景下对两个集合特性做相同、相反和相异分析并加以定量描述,获得两个集合的联系度表达式,在联系度的基础上再进行深入研究^[21]。其中,联系度表达式为:

$$\mu(\omega) = \frac{S}{N} + \frac{F}{N}i + \frac{P}{N}j \quad (1)$$

其中, $N=S+F+P$ 为所讨论的集对具有的特性总数, S 为集对中两个集合共同具有的特性总数, P 为集对中两个集合相互对立的特性总数, F 是集对中两个集合既不相同又不互相对立的特性总数, $\frac{S}{N}$, $\frac{F}{N}$ 和 $\frac{P}{N}$ 分别称为在某种背景下集对的一度、差异度、对立(相反)度; j 为对立度的系数,规定取值为 -1 ; i 为差异度的系数,规定在 $[-1, 1]$ 区间视不同情况取值。

令 $a = \frac{S}{N}$, $b = \frac{F}{N}$, $c = \frac{P}{N}$, 则式(1)可以简写为:

$$\mu = a + bi + cj \quad (2)$$

因 $N=F+S+P$, 故:

$$a + b + c = 1 \quad (3)$$

4.2 多属性决策集对分析模型

设所需评价的单级多属性决策问题方案集为 S_i ($i=1, 2, \dots, m$), 评价指标集为 E_t ($t=1, 2, \dots, n$), W_t 为 E_t 的权重系数,且满足 $W_t > 0$, $\sum_{t=1}^n W_t = 1$ 。令 D_{it} 为方案 S_i 关于指标 E_t 的属性值,由 D_{it} 构成决策矩阵:

$$D = \begin{bmatrix} D_{11} & \dots & D_{1n} \\ \dots & \dots & \dots \\ \dots & \dots & \dots \\ D_{m1} & \dots & D_{mn} \end{bmatrix} = (D_{it})_{m \times n} \quad (4)$$

其中, d_{it} 是每个 E_t 的各个指标值。为了规范化处理,按照式(5)对每个指标值进行归一化:

$$d_{it} = \frac{D_{it} - \min(D_{it})}{\max(D_{it}) - \min(D_{it})}, i=1, \dots, m; t=1, \dots, n \quad (5)$$

设在同一属性中,归一化后的指标值的范围是 $[x_i, y_i]$, x_i, y_i 分别是某属性下的上、下极限值, $\frac{d_{it}}{x_i}, \frac{y_i}{d_{it}}$ 分别表示 d_{it} 与 x_i 和 y_i 的接近程度。比较空间 $[x_i, y_i]$ 中的同一度 a_i 、对立度 c_i 分别为:

$$a_i = \frac{d_{it}/x_i}{1 + y_i/x_i} \quad (6)$$

$$c_i = \frac{y_i/d_{it}}{1 + y_i/x_i} \quad (7)$$

结合权重进行扩展,可得:

$$a_i = \sum_{t=1}^n W_t a_{it}, c_i = \sum_{t=1}^n W_t c_{it} \quad (8)$$

贴近度 r_i 为:

$$r_i = \frac{a_i}{a_i + c_i} \quad (9)$$

最终根据各方案的贴近度的大小,得出评价方案的优劣顺序。

5 基于集对理论的矢量地图水印评价模型实例分析

本节结合第 3 节和第 4 节的分析,提出一种基于集对理

论的矢量地图水印评价模型,其主要流程如图 2 所示。

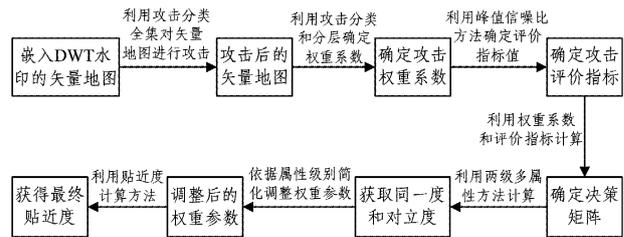


图 2 模型的主要流程图

5.1 基于集对理论的水印算法评价模型的主要流程

主要步骤如下:

步骤 1 依据第 3 节总结的攻击分类对嵌入水印的矢量地图进行全集攻击;

步骤 2 依据攻击分类使用频率确定各攻击权重系数;

步骤 3 使用峰值信噪比方法构建百分制评价矩阵,并以此作为依据,利用式(4)、式(5)建立决策矩阵;

步骤 4 采取低级多属性向高级多属性递归的方法计算多级联系度相关系数;

步骤 5 依据属性级别简化调整权重参数;

步骤 6 利用式(8)、式(9)计算贴近度,进而获得最终贴近度值。

最终获取的贴近度即可作为某矢量地图数字水印算法的评价依据,数值范围为 $[0, 1]$,数值越大则说明该数字水印算法抗各种攻击的综合能力越强,算法鲁棒性越强。

在后续内容中,将结合最流行的现有矢量地图水印算法(DWT 水印算法)来给出该模型评价的主要过程,并对主要过程进行详细描述。

5.2 确定权重系数

权重系数表征的是某种属性在整个属性集中所处的地位,即该属性在属性集中的重要程度。本评价模型根据某种攻击对一般矢量地图水印嵌入算法的攻击效果或该种攻击在矢量地图一般操作中的使用频率或地图本身的特性来确定权重系数。本文给出了一种权重系数确定方案,权重系数确定者可以根据使用目的进行权重系数的调整。

一级评价指标体系:

$$\omega = (\omega_1, \omega_2) = (0.3, 0.7)$$

二级评价指标体系:

$$E_t (t=1, 2), W_{E1} = (0.6, 0.4), W_{E2} = (0.25, 0.25, 0.25, 0.25)$$

三级评价指标体系:

$$E_{k} (k=1, 2, 4), W_{E11} = (0.1, 0.1, 0.1, 0.4, 0.3) \\ W_{E12} = (0.3, 0.4, 0.3), W_{E24} = (0.3, 0.4, 0.3)$$

四级评价指标体系:

$$E_{l} (l=1, \dots, 5), W_{E114} = (0.2, 0.2, 0.2, 0.2, 0.2) \\ W_{E115} = (0.3, 0.4, 0.3), W_{E121} = (0.3, 0.4, 0.3) \\ W_{E122} = (0.3, 0.4, 0.3), W_{E123} = (0.3, 0.4, 0.3)$$

5.3 构建决策矩阵

由 4.2 节可知,首先获取各属性的评价指标值构成决策矩阵,指标值采用百分制作为打分规则。本文通过计算提取出的水印的峰值信噪比 PSNR 参数并将其换算成百分制来

作为指标取值。各属性的评价指标如表1—表4所列。

表1 特殊攻击评价价值

评价	特殊攻击	
	地图融合	扰乱攻击
	评价1	43
评价2	40	23
评价3	37	25

表2 空域评价价值

评价	属性 简化	地图 拟合	拟 部分 扰乱	噪声攻击			几何攻击				
				高斯	椒盐	随机	平移	剪切	缩放	旋转	扭曲
评价1	62	98	86	45	42	40	49	20	96	99	48
评价2	60	97	85	48	40	38	48	21	97	97	46
评价3	58	99	83	43	43	45	40	25	96	95	43

表3 频域评价价值

评价	属性	DCT+噪声			DWT+噪声			DFT+噪声		
		高斯	椒盐	随机	高斯	椒盐	随机	高斯	椒盐	随机
评价1		85	86	86	93	95	94	81	84	86
评价2		83	83	84	96	97	93	83	82	81
评价3		80	81	85	91	93	91	85	80	82

表4 无意攻击评价价值

评价	属性	压缩 攻击	格式 转换	投影 变换	地图更新		
					图元修改	图元增加	图元删除
评价1		71	95	97	85	84	79
评价2		72	96	96	83	86	76
评价3		70	98	93	81	81	73

以上根据式(4)、式(5)确定了规范化后的矩阵。为了方便,首先对较低级别的多属性进行处理,将有意攻击的子集进行合并。有意攻击决策矩阵为:

$$D_1 = \begin{bmatrix} 0.53 & 0.99 & 0.84 & 0.32 & 0.28 & 0.25 & 0.37 & 0 & 0.96 & 1 & 0.35 & 0.82 & 0.84 & 0.84 & 0.92 & 0.95 & 0.94 & 0.77 & 0.81 & 0.84 & 0.29 & 0 \\ 0.51 & 0.97 & 0.82 & 0.35 & 0.25 & 0.26 & 0.32 & 0.01 & 0.97 & 0.97 & 0.33 & 0.79 & 0.79 & 0.81 & 0.96 & 0.97 & 0.92 & 0.79 & 0.78 & 0.77 & 0.25 & 0.04 \\ 0.48 & 1 & 0.79 & 0.29 & 0.29 & 0.32 & 0.25 & 0.06 & 0.96 & 0.95 & 0.29 & 0.76 & 0.77 & 0.85 & 0.89 & 0.92 & 0.89 & 0.82 & 0.76 & 0.78 & 0.22 & 0.06 \end{bmatrix}$$

无意攻击决策矩阵为:

$$D_2 = \begin{bmatrix} 0.04 & 0.89 & 0.96 & 0.54 & 0.5 & 0.32 \\ 0.07 & 0.93 & 0.93 & 0.46 & 0.57 & 0.21 \\ 0 & 1 & 0.82 & 0.39 & 0.39 & 0.11 \end{bmatrix}$$

5.4 计算多级联系度的相关系数

前面介绍了单级多属性的相关系数多属性决策方案,多级多属性的决策方案是单级多属性的一种延伸扩展,由最低级逐级按照单级多属性方法向上递归,最终即可得到决策方案。在5.3节中已经进行了简化,现在仅需求取两级多属性的相关系数,求取方法如式(6)和式(7)所示,同一度为:

$$a_{111} = 0.53 / (1+0) = 0.53 \tag{10}$$

$$a_{121} = 0.51 / (0.97+0.01) \approx 0.51 \tag{11}$$

同理,可以求取其他同一度:

$$c_{111} = (1 \times 0) / ((1+0) \times 0.53) = 0 \tag{12}$$

$$c_{121} = (0.97 \times 0.01) / ((0.97+0.01) \times 0.51) \approx 0.025 \tag{13}$$

同理,可以求取其他所有对立度,此处从略。

5.5 获得贴近期

根据5.3节中的属性级别的简化,相应的各属性权值也应该做相应的调整,如噪声中的高斯噪声新权值为:

$$W_{高斯} = W_{高斯权值} \times W_{噪声权值} \times W_{空域权值} = 0.4 \times 0.3 \times 0.5 = 0.06 \tag{14}$$

按照新权值的计算方法,根据式(8)和式(9)可得该算法的贴近期 $r = a / (a+c) = 0.724$ 。用同种方法对基于网格的水印算法进行了评价,所得贴近期 $r' = 0.52, r > r'$,故通过结合集对理论与矢量水印攻击分类可以较客观地得出: r 所采用的水印算法整体上优于 r' 所采用的水印算法。

结束语 本文通过分析矢量地图数字水印算法和矢量地图自身的特点,总结了针对矢量地图水印较为全面的攻击分类,并在此基础上结合集对理论,提出一种矢量地图数字水印算法的评价模型。通过实例分析结果表明,该模型能够从综合客观的角度来分析评价矢量地图水印算法,为矢量地图水印系统提供了一种比较通用的评价模型。

参考文献

[1] 赵克勤. 集对分析对不确定性的描述和处理[J]. 信息与控制,

1995,24(3):162-166.
 [2] JIANG Y L, XU C F, LIU Y. A New Approach for Representing and Processing Uncertainty Knowledge[C]// The 2003 IEEE Intl. Conf on Information Reuse and Integration. 2003:466-470.
 [3] XU Z, YU R, PAN X Z. Watermark Embedded in Polygonal Line for Copyright Protection of Contour Map [J]. International Journal of Computer Science and Network Security, 2006, 6(7B):202-205.
 [4] SHAO C Y, WANG X T, AN J L. A Robust Algorithm for Watermarking 2D Vector Maps with Low Shape-Distortions[J]. Journal of China Ordnance, 2006, 2(3):233-236.
 [5] GERRIT S, VOIGT M. A High Capacity Watermarking System for Digital Maps[C]//Proceedings of the 2004 Multimedia and Security Workshop on Multimedia and Security. Germany, Magdeburg, 2004:180-186.
 [6] 王伟. 一种鲁棒性的2D矢量图形水印算法[J]. 中国图像图形学报, 2007, 12(2):200-205.
 [7] 李媛媛, 许录平. 用于矢量地图版权保护的数字水印[J]. 西安电子科技大学学报(自然科学版), 2004, 31(5):719-723.
 [8] 李安波. GIS矢量数据产品版权保护的关键技术研究[D]. 南京: 南京师范大学, 2007.
 [9] 陈晓光, 李岩. 针对二维矢量图形数据的盲水印算法[J]. 计算机应用, 2011, 31(8):2174-2177.
 [10] SOLACHIDIS V, PITAS I. Watermarking polygonal lines using fourier descriptors[C]//Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing. Turkey, Istanbul, 2000:1955-1958.
 [11] KITAMURA I, KANA S, KISHINAMI T. Copyright Protection of Vector Map using Digital Watermarking method based on Discrete Fourier Transform[C]//Proceedings of IEEE 2001 International Geoscience and Remote Sensing. 2001:1191-1193.
 [12] GLANNOULA A, NIKOLAIDID N, PITAS I. Watermarking of sets of polygonal lines using fusion techniques[C]//Proceedings of the 2002 IEEE International Conference on Multimedia and Expo. Switzerland, Lausanne, 2002:26-29.
 [13] 王奇胜. 基于DFT的矢量地理空间数据数字水印技术研究[D]. 郑州: 信息工程大学测绘学报, 2008.
 [14] VOIGT M, YANG B, BUSCH C. Reversible watermarking of 2D-vector data [C]// Proceeding of the 2004 Multimedia and

Security, Magdeburg, Germany, 2004: 160-165.

- [15] KITAMURA I, KANAN S, KISHINAMI T. Digital watermarking method for vector map based on wavelet transform [C]// Proceedings of the Geographic Information Systems Association, Tokyo, Japan, 2000: 417-421.
- [16] LI Y, XUL A. A blind watermarking of vector graphics images [C]// Proceedings of the Fifth International Conference on Computational Intelligence and Multimedia Applications, 2003: 27-30.
- [17] 李媛媛, 许录平. 矢量图形中基于小波变换的盲水印算法[J]. 光子学报, 2004, 33(1): 97-100.

- [18] 易开祥, 石教英, 孙鑫. 数字水印技术研究进展[J]. 中国图像图形学报, 2001, 6(2): 111-117.
- [19] 孙建国. 矢量地图数字水印技术研究[M]. 北京: 人民邮电出版社, 2012: 20-80.
- [20] VASSAUX B, NGUYEN P, BAUDRY S, et al. A survey on attacks in image and video watermarking [C]// Proceedings of the SPIE-The International Society for Optical Engineering, 2002: 169-179.
- [21] 成科扬. 论集对分析在软件质量评价体系建立中的运用[J]. 计算机应用与软件, 2004, 21(3): 19-21.

(上接第 341 页)

表 3 几种特征码提取方法的分析比较/%

方法	成功率(SR)	误报率(FPR)
3-Gram	90.76	2.51
4-Gram	91.82	2.73
CVSAE	94.23	2.84
CSR	81.23	3.21
改进的变长 N-Gram	93.21	1.65

由表 3 可知, 与其他几种特征码提取方法相比, 改进的 N-Gram 特征提取算法可以成功地提取出病毒特征码, 并且误报率较低。

结束语 本文采用变长 N-Gram 作为提取特征, 同时考虑特征频率的作用, 利用特征浓度有向选择生成数据字典, 选择有效特征, 并将其加入到病毒 N-Gram 特征库中, 与待测病毒样本进行匹配; 借助 N-Gram 语言模型提取病毒特征码, 实现病毒特征码的自动提取工作。从实验结果可以看出, 本文提出的基于 N-Gram 改进的病毒特征码提取算法能够有效地提取出计算机病毒特征码, 且其性能优于传统的 N-Gram 算法。下一步将研究更为有效的特征选择方法, 通过组合使用多个特征码来表示一类病毒, 实现病毒的分类检测。

参 考 文 献

- [1] YEGNESWARAN V, GIFFIN J T, BARFOD P, et al. An architecture for generating semantics-aware signatures [C]// Conference on Usenix Security Symposium, USENIX Association, 2004: 7-7.
- [2] LEE H, KIM W, HONG M. Biologically Inspired Computer Virus Detection System [J]. Lecture Notes in Computer Science, 2004, 3141: 153-165.
- [3] KIJEWski P. Automated Extraction of Threat Signatures from Network Flows [OL]. <http://www.first.org/conference/2006/papers/kijewski-piotr-paper.pdf>.
- [4] KREIBICH C, ROWCROFT J. Honeycomb: creating intrusion detection signatures using honeypots [J]. Acm Sigcomm Computer Communication Review, 2015, 34(1): 51-56.

- [5] 张小康, 帅建梅, 史林. 基于加权信息增益的恶意代码检测方法 [J]. 计算机工程, 2010, 36(6): 149-151.
- [6] KEPHART J O, ARNOLD W C. Automatic extraction of computer virus signatures [C]// 4th Virus Bulletin International Conference, 1994.
- [7] 张福勇. 基于 n-gram 词频的恶意代码特征提取方法 [J]. 网络安全技术与应用, 2015(11): 88-89.
- [8] 白金荣, 王俊峰, 赵宗渠. 基于 PE 静态结构特征的恶意软件检测方法 [J]. 计算机科学, 2013, 40(1): 122-126.
- [9] RAFF E, ZAK R, COX R, et al. An investigation of byte n-gram features for malware classification [J]. Journal of Computer Virology & Hacking Techniques, 2016: 1-20.
- [10] 曾键, 赵辉. 一种基于 N-Gram 的计算机病毒特征码自动提取方法 [J]. 计算机安全, 2013(10): 2-5.
- [11] 李沁蕾, 王蕊, 贾晓启. OSN 中基于分类器和改进 n-gram 模型的跨站脚本检测方法 [J]. 计算机应用, 2014, 34(6): 1661-1665.
- [12] DHAYA R, POONGODI M. Detecting software vulnerabilities in android using static analysis [C]// International Conference on Advanced Communication, Control and Computing Technologies, 2014.
- [13] O'KANE P, SEZER S, MCLAUGHLIN K. N-gram density based malware detection [C]// Computer Applications & Research, IEEE, 2014: 1-6.
- [14] SHABTAI A, MOSKOVITCH R, FEHER C, et al. Detecting unknown malicious code by applying classification techniques on OpCode patterns [J]. Security Informatics, 2012, 1(1): 1-22.
- [15] SANTOS I, BREZO F, UGARTE-PEDRERO X, et al. Opcode sequences as representation of executables for data-mining-based unknown malware detection [J]. Information Sciences, 2013, 231(9): 64-82.
- [16] 吴军. 数学之美 [M]. 北京: 人民邮电出版社, 2012.
- [17] 恶意代码网站 [OL]. <http://vxheaven.org>.
- [18] 金雄斌. 计算机病毒特征码自动提取技术的研究 [D]. 武汉: 华中科技大学, 2011.
- [19] TANG Y, XIAO B, LU X. Using a bioinformatics approach to generate accurate exploit-based signatures for polymorphic worms [J]. Computers & Security, 2009, 28(8): 827-842.