

# 基于图象内容的数字水印模型

刘瑞祯 王蕴红 谭铁牛

(中国科学院自动化研究所模式识别国家重点实验室, 北京 100080)

**摘要** 为了提高水印算法的鲁棒性、安全性,提出了一种基于图象内容的水印模型(CBWM: Content-Based Watermarking Model),其基本原理是水印向量被嵌入到图象的特征向量中,且两者相互正交.在此模型中,静止图象的水印检测被归结于 Neymann-Pearson 准则下的统计模型.将该模型应用于通用的 DCT 水印算法,从实验结果看,算法性能有很大地提高.

**关键词** 数字水印 正交 基于内容 Neymann-Pearson 准则 统计模型

中图法分类号: TN919 文献标识码: A 文章编号: 1006-8961(2001)06-0558-05

## Image Content-Based Watermarking Model

LIU Rui-zhen, WANG Yun-hong, TAN Tie-niu

(National Laboratory of Pattern Recognition Institute of Automation, Chinese Academy of Sciences, Beijing 100080)

**Abstract** In this paper, we proposed a new general additive watermarking model based on the content of digital images, called as CBWM (Content-Based Watermarking Model). It provides a common basis to study many existing watermarking algorithms that are based on unitary transforms and to evaluate their performance and characteristics. CBWM is designed to address two important issues of digital watermarking. One is the requirement of robustness. In order to improve robustness and security, the embedded watermark is designed to be orthogonal to the feature vector of the original image, which means that watermarking casting is image content dependent. The second issue is watermark detection. CBWM presents a statistical approach to watermark detection based on the Neymann-Pearson criterion and describes a method for computing the probability of false alarm and missing rejection during watermark detection. In experimental tests, CBWM is applied to the popular discrete cosine transform and promising results are obtained.

**Keywords** Digital watermarking, Orthogonal, Content-based, Neymann-Pearson criterion, Statistical model

## 0 概 况

当今信息媒体的数字化为信息的存取提供了极大的便利,同时也显著地提高了信息表达的效率 and 准确度,特别是随着计算机网络通讯技术的发展,使数据的交换和传输变成了一个相对简单的过程,如今人们借助于计算机、数字扫描仪、打印机等电子设备已可方便、迅捷地将数字信息传达到世界各地,但随之而来的副作用是通过网络传输的数据文件或作品,使有恶意的个人或团体有可能在没有得到作品

所有者的许可下,进行拷贝和传播有版权的内容,因此如何在网络环境中实施有效的版权保护(Copyright Protection)和信息安全(Information Security)手段已成为一个迫在眉睫的现实问题,而数字水印(Digital watermarking)技术则为上述问题提供了一个潜在的解决方案.

所谓数字水印就是向欲被保护的多媒体数据嵌入某种信息(即水印),以满足保护所有者的权益等要求,如今图象的水印技术根据水印嵌入的方式可以大致分为:空间域技术(水印被直接嵌入在图象的亮度值上)<sup>[1,2]</sup>和变换域技术(即首先将图象做某种数学变

换,然后将水印嵌入到变换系数中)<sup>[3]</sup>两类.从目前的情况看,变换域方法的应用正变得日益普遍,因为变换域方法通常都具有很好的鲁棒性,其对图象压缩、常用的图象滤波以及噪声均有一定的抵抗力,而且其中绝大多数均使用了酉变换,如离散余弦变换<sup>[3,4]</sup>、离散傅立叶变换<sup>[5]</sup>和离散小波变换<sup>[6,7]</sup>等等.

大家知道,一个有效的水印算法至少应满足下面两个特性:①不可觉察性(Imperceptibility),即添加水印后的数字文档与原始文档对人的感觉器官应是一样的;②鲁棒性(Robustness),即给定一个水印文档,非授权的个人或团体在使文档可用的情况下无法剔除水印.由此可见,这两个特性是互相冲突的,所以很多水印算法均是首先从原始图象的变换域系数中提取一个子集,然后将水印嵌入其中.为了保证水印算法的鲁棒性,水印应该嵌入到图象的最显著分量中<sup>[3]</sup>,或者说,水印应该和图象的特征集捆绑在一起,即选择被嵌入的系数集合通常也可以被看成为图象的特征向量集.目前的大多数水印算法中,水印的选择一般均与图象内容无关,这实际上对算法的安全性和鲁棒性是不利的,因为当带水印的图象遭受到有意的(如恶意的破坏或删除水印)或无意的攻击行为(如图象压缩、滤波、扫描与复印、噪声污染、尺寸变化等等)时,如果嵌入的水印与图象无关,攻击者往往可以很容易地在不破坏图象基本质量的情况下而去掉水印;而当水印与图象特征捆绑在一起后,则在水印被破坏的同时,图象的特征也被破坏,从而再没有应用价值.

## 1 CBWM 的基本原理

CBWM 是基于变换域的通用数字水印模型,从线性代数的角度讲,一幅观察到的灰度图象可看成是由非负元素构成的矩阵.一般地,将大小为  $M \times N$  的图象定义为  $A = \{a_{ij}\} \in F^{M \times N}$  ( $i = 1, 2, \dots, M; j = 1, 2, \dots, N$ ), 其中,  $F$  为实数域或复数域.  $\tilde{A} = \{\tilde{a}_{ij}\} \in F^{M \times N}$  是  $A = \{a_{ij}\}$  的变换域矩阵,即对于映射算子  $T$ ,有  $T: A \rightarrow \tilde{A}$ . 向量  $C = \{c_k\} \in F^K$  为变换域矩阵  $\tilde{A} = \{\tilde{a}_{ij}\} \in F^{M \times N}$  的一个子集,即  $\{c_k\} \in \{\tilde{a}_{ij}\}$ , 其中,  $k = 1, 2, \dots, K$ . 为了保证水印算法的鲁棒性,一般应要求水印嵌入到图象最显著的分量上,这样可把向量  $\{c_k\}$  看成是图象的特征向量,因此水印  $W = \{w_k\} \in F^K$  需被嵌入到图象的特征向量  $\{c_k\}$  中.

### 1.1 水印嵌入算法

当前大多数的水印算法都采用了扩展谱通信的方法<sup>[3,6]</sup>,其中,原始数据的模型均假设为非周期、零均值、宽平稳的随机过程,但事实上这种假设不总是正确,且常会导致一些错误结果.本文描述的 CBWM 对原始图象数据不做任何假设,它是一个加性水印嵌入模型,且选择的水印与原始图象的特征向量正交,因此是一个基于图象内容的水印模型.

若将水印  $\{w_k\}$  嵌入到特征向量  $\{c_k\}$  中,则新的特征向量  $\tilde{C} = \{\tilde{c}_k\} \in F^K$  由下式得到

$$\tilde{c}_k = c_k + aw_k, \quad k = 1, 2, \dots, K \quad (1)$$

其中,尺度因子  $a$  是控制水印嵌入强度或能量的常数.这里将式(1)定义的水印嵌入模型称为可加性模型,其优点是在给定某种图象失真评价准则(如 PSNR: 峰值信噪比)后,可以精确计算水印的嵌入能量和嵌入容量(关于这一点将另文阐述).

用新得到的特征向量  $\{\tilde{c}_k\}$  替换原来的特征向量  $\{c_k\}$ , 并从  $\tilde{A}$  得到新的变换域矩阵  $\tilde{B}$ , 如果  $T$  是双射的(表明它有逆算子  $T^{-1}$ , 那么即可用  $\tilde{B}$  来重构水印图象  $B$ , 即  $T^{-1}: \tilde{B} \rightarrow B$ ).

设  $A = \{a_{ij}\} \in F^{M \times N}$  为原始图象,  $B = \{b_{ij}\} \in F^{M \times N}$  为水印图象,那么, CBWM 所描述的通用水印算法步骤如下:

(1)对原始图象做二维前向离散变换

$$T: A \rightarrow \tilde{A} \quad (2)$$

(2)从图象  $A$  的变换域中选择一个特征向量  $C = \{c_k\} \in \tilde{A}$ ;

(3)将水印  $W = \{w_k\}$  按式(1)嵌入到  $\{c_k\}$  中,以得到新的特征向量  $\tilde{C} = \{\tilde{c}_k\} \in \tilde{B}$ ;

(4)用新的变换域矩阵  $\tilde{B}$  重构水印图象  $B$

$$T^{-1}: \tilde{B} \rightarrow B \quad (3)$$

以上算法中,水印的选择与确定是关键所在.

### 1.2 水印检测框架

根据在水印检测过程中是否使用原始图象,水印算法大致可分为不需要使用原始图象检测水印<sup>[8~10]</sup>和需要使用原始图象检测水印<sup>[11,12]</sup>两类. CBWM 检测算法是否需要使用原始图象,取决于具体采用的图象变换方法以及特征向量的选取.当给定原始图象  $A = \{a_{ij}\} \in F^{M \times N}$ 、带水印的图象  $B = \{b_{ij}\} \in F^{M \times N}$  和水印  $W = \{w_k\}$  后,即得到特征向量  $\tilde{C} = \{\tilde{c}_k\} \in \tilde{B}$ (注意到对某些变换域和特征向量的提取方法,不需要使用原始图象就能得到  $\tilde{C}$ . 举一个简单的例子就是当对图象做 DCT 的块分解时,若

用 zig-zag 扫描来确定被叠加水印的系数,即特征向量  $C$ ,这样,在水印检测时,根据已知特征向量在块变换域中的位置,即可直接提取  $\tilde{C} = \{\tilde{c}_k\} \in \tilde{B}$ ,而不需要使用原始图象),然后按下式计算标量  $Z$  的值

$$\begin{aligned} Z &= \sum_k \tilde{c}_k \tau w_k = \sum_k (c_k + a \tau w_k) \tau w_k \\ &= \sum_k c_k \tau w_k + \sum_k a \tau w_k^2 \end{aligned} \quad (4)$$

给定阈值  $Z_T$ ,则可得到  $0 \sim 1$  决策,以确定水印的存在与否.假设检验为

$$H_1: Z > Z_T \quad H_0: Z < Z_T \quad (5)$$

其中,  $H_1$  表示水印存在,  $H_0$  表示水印不存在.

### 1.3 水印的选择

为了保证水印算法的鲁棒性和安全性,一个自然的选择是将水印与图象特征结合起来,即算法应该是基于图象内容的.对一幅图象的特征向量  $C = \{c_k\}$ ,一般均可找到相应的与  $C$  正交的水印空间  $F_W \subset F^K$

$$F_W = \{W: \sum_k c_k \tau w_k = 0, k = 1, 2, \dots, K\} \quad (6)$$

即对任何属于  $F_W$  的向量  $W = \{\tau w_k\} \in F_W$ ,有

$$C^H W = 0 \quad (7)$$

那么,式(4)可变为

$$Z = \sum_k \tilde{c}_k \tau w_k = a \sum_k \tau w_k^2 = aE \quad (8)$$

其中,变量  $E$  可看成是嵌入的水印能量.

由此可见,满足  $W \in F_W$  的水印有无穷多个,而水印空间  $F_W$  的确定,实际上是一个受约束的优化过程,即

$$W = \arg \min_{W \in F^K} |C^H W| \quad (9)$$

其中,  $C$  为原始图象的特征向量.

其优化的中止条件为

$$|C^H W| < \delta \quad (10)$$

式中,  $\delta$  ( $\delta > 0$ ) 是优化过程中预先定义的阈值.

若考虑使用有意义的水印(如文本、图象、公司标志或用户自定义的数字序列),则从安全和鲁棒性考虑,需要先将有意义的水印进行随机化或加密处理,以得到  $\tilde{W}$ ,这里,密码可以是随机种子或加密密钥等.那么优化过程就可表示为

$$D_{key} = \arg \min_{D_{key} \in R} |C^H \tilde{W}| \quad (11)$$

其中,密码  $D_{key}$  是欲搜索的优化变量,  $|C^H \tilde{W}|$  为目标函数,  $\tilde{W}$  定义为

$$\tilde{W} = p(W, D_{key}) \quad (12)$$

这里,加密函数  $p(\cdot)$  实际上是将有意义的水印空

间映射到与特征向量正交的向量空间中,即

$$F_{\tilde{W}} = \{\tilde{W} \in F^K: C^H \tilde{W} < \delta\} \quad (13)$$

从提高水印算法的安全性出发,加密函数  $p(\cdot)$  一般是单向和非对称的<sup>[13,14]</sup>.

## 2 统计分析

如果应用 CBWM 给出的水印算法,那么就可作为水印检测给出一个很好的统计模型.与很多水印方法不同的是,CBWM 不需要对原始文档和水印的分布进行某种统计假设.

### 2.1 统计模型

上面所描述的水印算法应称为水印检测,而不是水印解码,因为水印检测只给出一个水印存在与否的二值决策,即在给定一个预定义的正数  $\delta$  作为阈值的情况下,若图象的特征向量  $C$  和给出的水印  $W$  满足式(10),那么,水印检测用式(4),并得到

$$Z = \sum_k \tilde{c}_k \tau w_k \quad (14)$$

事实上,如果图象包含水印,那么则有

$$\begin{aligned} Z &= \sum_k \tilde{c}_k \tau w_k = \sum_k c_k \tau w_k + \\ &\quad \sum_k a \tau w_k^2 \leq \delta + m \cong m \end{aligned} \quad (15)$$

若图象不包含水印,则式(14)变为

$$Z = \sum_k \tilde{c}_k \tau w_k = \sum_k c_k \tau w_k \leq \delta \quad (16)$$

因为式(15)和(16)中,  $\delta$  非常小,可以忽略不计,所以水印检测可以化为如下的假设检验

$$H_1: Z = m + e(t), H_0: Z = e(t) \quad (17)$$

其中,  $m$  是常数,  $e(t)$  表示由图象的失真而产生的误差.这里,把由图象失真(如图象滤波、添加噪声、几何变换等等)产生的误差看成是随机噪声,并认为服从高斯分布  $e(t) \sim N(0, \sigma^2)$ .

根据 Neymann-Pearson 准则来计算决策阈值  $Z_T$ ,若  $Z > Z_T$ ,则判定水印存在,若  $Z < Z_T$ ,则判定水印不存在.这里只讨论适合于静止图象水印技术的 Neymann-Pearson 准则,因为在此条件下,先验概率未知;而对于视频或音频流的数字水印技术,当已知先验概率的情况下,就可应用其他的决策规则,如 Bayesian 规则、最小错误率规则等.

### 2.2 Neymann-Pearson 决策

对静止图象的水印检测,若事先未知水印的存在与否,这就意味着先验概率未知,所以人们希望在

给定虚警概率(probability of false alarm)  $P_F$  的情况下(即虽检测到水印,但水印实际上不存在),其漏报概率(probability of missing detection)  $P_M$ (即水印存在,但未被检测到)最小。

根据 Neymann-Pearson 准则,可以计算在固定  $P_F$  时,  $P_M$  为最小。根据式(17)描述的检测模型,概率密度似然函数为

$$f(Z|H_1) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{(Z-m)^2}{2\sigma^2}\right]$$

$$f(Z|H_0) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{Z^2}{2\sigma^2}\right] \quad (18)$$

此时,决策准则为

若  $Z \geq Z_T$ , 则  $H_1$  为真,表示水印存在;  
若  $Z < Z_T$ , 则  $H_0$  为真,表示水印不存在;  
阈值  $Z_T$  的值可由虚警概率  $P_F$  计算

$$P_F = \int_{Z_T}^{\infty} f(Z|H_0) dZ$$

$$= \int_{Z_T}^{\infty} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{Z^2}{2\sigma^2}\right) dZ \quad (19)$$

而漏报概率  $P_M$  可用下式计算

$$P_M = \int_{-\infty}^{Z_T} f(Z|H_1) dZ$$

$$= \int_{-\infty}^{Z_T} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(Z-m)^2}{2\sigma^2}\right) dZ$$

$$= 1 - f_{\text{err}}\left(\frac{Z_T - m}{\sigma}\right) \quad (20)$$

其中,误差函数  $f_{\text{err}}(\cdot)$  定义为

$$f_{\text{err}}(x) = \frac{1}{\sqrt{2\pi}} \int_x^{\infty} e^{-\frac{t^2}{2}} dt \quad (21)$$

### 3 实验结果

CBWM 是一个通用的水印算法,它可和其他任何基于变换域的水印方法相结合。该算法有两个关键之处:其一是如何选择原始图象的特征向量,因为特征向量的选取决定了算法性能的好坏,且在某种程度上决定了水印算法是否需要使用原始图象来检测水印;二是如何选择合适的水印,由于水印选择是一个受约束的优化过程,且目标函数是非光滑连续函数,因此应该仔细选择优化算法。

本文将 CBWM 应用于基于 DCT(离散余弦变换)变换的水印算法,其类似于 Cox 方法<sup>[3]</sup>。为了评价 CBWM 的鲁棒性能,试验时还应用了由 Petitcolas 等人提出并设计的水印攻击测试工具 Stir-

Mark<sup>[15,16]</sup>。

本文对各种不同的图象进行了测试,这里给出了用大小为  $256 \times 256$  的灰度图象 Lena 做实验所得到的结果,其测试方法与 Cox 方法类似,即先计算原始图象的 DCT 变换,其所获得的特征向量  $C = \{c_i\}$  即为 DCT 变换域上幅值最大的前 256 个系数(不包含直流分量),然后根据式(10)的要求计算出水印,该水印为与  $C = \{c_i\}$  正交的一个  $256 \times 1$  的二值向量,即  $W = \{w_i\} \in \{1, -1\}, i=1, 2, \dots, 256$ 。从实验结果可见,由于水印空间  $F_W = \{W \in F^K; C^H W < \delta\}$  是一个线性子空间,因而满足条件的水印很容易搜索到。这里仅采用简单的随机搜索方法,若要嵌入有意义的水印,那么可考虑采用稍微复杂的优化算法(如遗传算法等),但需要注意的是,目标函数是非连续和非光滑的,因此不能使用传统的梯度优化算法。实验时,尺度因子设为  $a=0.2$ ,然后叠加水印到向量  $C$  中,得到新的特征集  $\tilde{C} = \{\tilde{c}_k\}$ ,再做反变换即得带水印的 Lena 图象。原始图象和水印图象之间的峰值信噪比 PSNR 为 38.06dB。

然后再用 StirMark 来模拟 88 种图象的失真情形,以完成对算法的鲁棒性测试,并用 Neymann-Pearson 准则进行水印测试。由于 DCT 方法在选取特征向量的时候,实际上已隐含了排序过程,因此在检测过程中需要使用原始图象来提取特征向量  $\tilde{C} = \{\tilde{c}_k\}$ 。

为了对比本文方法的有效性,用 Stirmark 对 Cox 方法进行了鲁棒性测试,表 1 列出了实验结果。

表 1 Cox 方法与 CBWM 方法的实验结果对比

测试类型	测试数目	Cox 方法 正确检测数目	CBWM 方法 正确检测数目
对称及非对称地移去图象的行和列	5	5	5
滤波(中值、高斯、FMLR、锐化)	6	5	6
JPEG 压缩	12	12	12
中心裁剪	9	2	3
通用线性几何变换	3	2	3
改变 $x$ - $y$ 轴的显示比例	8	8	8
带裁剪、不带尺度变换的旋转	16	8	12
带裁剪和尺度变换的旋转	16	8	12
尺度变换	6	6	6
$x$ - $y$ 方向的修剪	6	3	6
Stirmark 随机弯曲	1	1	1
合计	88	60	74

从表 1 可见,基于 DCT 的 CBWM 不能正确检测中心裁剪和图象旋转两类图象失真,即如果水印图象的裁剪超过 10%或旋转超过 15°后,那么水印即不能被正确检测到,而对其他情形,嵌入的水印均可被正确检测;而 Cox 方法则对多种情形均不能检测到水印,除了裁剪和旋转外,4×4 滤波器的中值滤波、几何变换以及  $x$ - $y$  方向的修剪均可导致 Cox 方法失效.从总的正确检测率来看,Cox 方法为  $60/88=68.18\%$ ,小于 CBWM 的  $74/88=84.09\%$ .

## 4 结 论

本文提出了一个基于图象内容的通用数字水印模型 CBWM,在该模型中,水印被嵌入到图象的特征向量中,且与特征向量正交,而且该模型可以和常用的基于变换域的水印算法相结合.另外,对静止图象,水印的检测实际上可归纳为信号的假设检验,并可用 Neymann-Pearson 准则来检测水印的存在与否.本文还用 Stirmark 作为测试工具来检验 CBWM 的鲁棒性.从实验结果看,将 CBWM 与 DCT 变换结合起来,其所得到的结果优于 Cox 方法.

## 参 考 文 献

- 1 Nikolaidis N, Pitas I. Robust image watermarking in the spatial domain. *Signal Processing*, 1998,66(3):385~403.
- 2 Darmstadter V, Delaigle J F, Quisquater J J *et al.* Low cost spatial watermarking. *Computers & Graphics*, 1998, 22(4): 417~424.
- 3 Cox I J, Kilian J, Leighton F T *et al.* T. Secure spread spectrum watermarking for multimedia. *IEEE Trans. on Image Processing*, 1997,6(12):1673~1687.
- 4 Barni M, Bartolini F, Cappellini V *et al.* A DCT-domain system for robust image watermarking. *Signal Processing*, 1998,66(3): 357~372.
- 5 O'Ruanaidh J J K, Csurka G. A bayesian approach to spread spectrum watermark detection and secure copyright protection for digital image libraries. In:Proc. of CVPR'99, Fort Collins, Colorado, 1999,1:207~212.
- 6 Hsu C T, Wu J L. Multiresolution watermarking for digital images. *IEEE Trans. on Circuits and Systems II-Analog and Digital Signal Processing*, 1998,45(8):1097~1101.
- 7 Swanson M D, Zhu B, Tewfik A H. Multiresolution scene-based video watermarking using perceptual models. *IEEE Journal on Selected Areas in Communications*, 1998, 16(4): 540~550.
- 8 Chug T Y, Hong M S, Oh Y N *et al.* Digital watermarking for copyright protection of MPEG2 compressed video. *IEEE Trans.*

- on *Consumer Electronics*, 1998,44(3):895~901.
- 9 Voyatzis G, Pitas I. Digital image watermarking using mixing systems. *Computers & Graphics*, 1998, 22(4):405~416.
- 10 Barni M, Bartolini F, Cappellini V *et al.* Copyright protection of digital images by embedded unperceivable marks. *Image and Vision Computing*, 1998,16(12-13):897~906.
- 11 Kundur D, Hatzinakos D. A robust digital image watermarking method using wavelet-based fusion. In:Proc. of ICIP'97, Washington DC,1997,1:544~547.
- 12 Hsu C T, Wu J L. Hidden digital watermarks in images. *IEEE Trans. on Image Processing*, 1999,8(1):58~68.
- 13 Craver S, Memon N, Yeo B L *et al.* Resolving rightful ownerships with invisible watermarking techniques: Limitations, attacks, and implications. *IEEE Journal on Selected Areas in Communications*, 1998,16(4):573~586.
- 14 Qiao L T, Nahrstedt K. Watermarking schemes and protocols for protecting rightful ownership and customer's rights. *Journal of Visual Communication and Image Representation*, 1998,9(3): 194~210.
- 15 Kutter M, Petitcolas F. A fair benchmark for image watermarking systems. In:Electronic Imaging'99, Security and Watermarking of Multimedia Contents, San Jose, CA. 1999,3657:226~239.
- 16 Petitcolas F, Anderson R J, Kuhn M G. Attacks on copyright marking systems. In:Proceeding of Conference on Information Hiding'98,Portland Oregon, U. S. A,1998:218~238.

刘瑞祯 1969 年生,1990 年获北京航空航天大学学士学位,1998 年入中国科学院自动化所模式识别国家重点实验室攻读博士学位.主要从事数字水印、图形、图象及视频处理、计算机视觉等领域的研究.

王蕴红 1968 年生,1989 年获西北工业大学学士学位,1998 年获南京理工大学博士学位.现为中国科学院自动化所模式识别国家重点实验室副研究员,主要从事图象处理、模式识别、神经网络及信号处理等领域的研究.

谭铁牛 1964 年生,1989 年获英国帝国理工学院博士学位.现为中国科学院自动化研究所所长,模式识别国家重点实验室主任、研究员、博士生导师,中国“五四”青年奖章获得者.主要从事图象处理、计算机视觉和模式识别等相关领域的研究工作.