

版权保护中的组合安全码及相关问题

献给朱烈教授 80 华诞

范金萍¹, 顾玉杰^{2*}, 缪莹³

1. 上海师范大学数理学院, 上海 200234;

2. Faculty of Information Science and Electrical Engineering, Kyushu University, Fukuoka 819-0395, Japan;

3. Faculty of Engineering, Information and Systems, University of Tsukuba, Ibaraki 305-8573, Japan

E-mail: jpfan@shnu.edu.cn, gu@inf.kyushu-u.ac.jp, miao@sk.tsukuba.ac.jp

收稿日期: 2022-04-30; 接受日期: 2022-07-22; 网络出版日期: 2022-09-28; * 通信作者

国家自然科学基金(批准号: 12101414)和日本学术振兴会科研项目(批准号: 21K13830 和 18H01133)资助项目

摘要 现代科技的快速发展给数据传播和流通提供了便捷,同时也对数据内容的版权保护构成了巨大的威胁. 本文聚焦版权保护中的对抗合谋攻击的数学理论及其最新研究进展. 针对广播加密和多媒体指纹识别等不同场景的应用, 本文提出可追踪分配方案和防诬陷分配方案的统一的数学模型. 在此基础上, 本文介绍具体的可追踪组合安全码和防诬陷组合安全码, 以及关于它们的最大码字个数的上下界和具体构造等研究的组合方法、最新结果和公开问题. 此外, 本文也将介绍版权保护与群试理论和多用户通信领域相关组合问题的联系.

关键词 版权保护 组合安全码 可追踪分配方案 防诬陷分配方案 集合系 群试理论 多用户通信

MSC (2020) 主题分类 05B20, 05D05, 05D40, 68P30, 94A60, 94B25

1 版权保护中组合安全码的研究历史和发展

1.1 研究背景简介

随着科学技术的发展, 多媒体技术日益成熟, 常见的视频和图片等多媒体内容在我们的日常生活中得到大量的流通和使用. 然而, 由于多媒体内容复制便利和传播手段简单等特点, 授权用户 (authorized user) 在获取多媒体内容之后可能对内容进行复制并传播给非授权用户而从中获取利益. 这种非法复制和传播的行为对多媒体内容的版权保护而言是极大的威胁和挑战.

1994年, Chor等^[18]首次在广播加密 (broadcast encryption) 的应用场景下提出了对抗合谋攻击 (collusion attack) 的可追踪分配方案的数学方法. 拥有数据版权的发行商对数据进行加密并对密文进行广播, 任何用户都可以收到广播的密文. 一个授权用户在购买版权后会被分配到一个唯一的个人解

英文引用格式: Fan J P, Gu Y J, Miao Y. Combinatorial secure codes for copyright protection and related problems (in Chinese). Sci Sin Math, 2023, 53: 123–150, doi: 10.1360/SSM-2022-0079

密密钥, 可以用来解密密文并得到原始数据. 由于每个授权用户被分配到的个人解密密钥都具有唯一性和身份认证的性能, 因而单个授权用户不会将个人解密密钥直接进行非法复制利用, 但是多个授权用户会利用他们所拥有的不同的个人解密密钥进行合谋攻击, 制造出一个海盗版 (pirated copy) 解密密钥, 并进行非法使用或交易. Chor 等^[18]可追踪分配方案主要考虑如何合理地设计授权用户的个人解密密钥分配问题, 以及在抓获海盗版解密密钥后如何快速地追踪到至少一名合谋攻击者的问题. 1998 年, Boneh 和 Shaw^[11]提出了在广播加密应用场景下如何保护无辜用户 (innocent user) 不会被多个用户的合谋攻击者集团诬陷的问题. 此外, 近些年在多媒体传播等领域, 基于数字水印 (digital watermarking) 技术, 具有对抗合谋攻击能力的多媒体指纹识别码 (multimedia fingerprinting codes) 也相继被提出 (参见文献 [65, 88]). 值得注意的是, 在广播加密和数字指纹识别 (digital fingerprinting) 等不同的应用场景下, 授权用户的个人解密密钥或者数字指纹的分配方式不同, 合谋者所采取的攻击方式也可能不同, 因此已经存在的组合安全码 (combinatorial secure codes) 和集合系 (set systems) 等也不尽相同. 在接下来的第 1.2 小节中, 将针对这些不同应用场景中的对抗合谋攻击问题提出统一的数学刻画模型.

1.2 研究问题

由于密钥/指纹分配方案需要满足的安全性不同, 我们以如下的可追踪分配方案和防诬陷分配方案两大类来进行阐述.

定义 1.1 (可追踪分配方案) 一个拥有 M 个用户 $[M] := \{1, 2, \dots, M\}$ 且能够防止至多 t ($\leq M$) 个用户合谋攻击的可追踪分配方案是一个四元组 $(\text{Enc}, \text{Tra}; \mathcal{A}, t)$, 其中,

- (1) $\text{Enc} : [M] \rightarrow \Gamma^n$ 是一个密钥/指纹分配方案, 其中 Γ 为一个有限集合;
- (2) $\mathcal{A} : 2^{\Gamma^n} \rightarrow \Upsilon^n$ 是一个攻击策略, 即合谋团体 $S \subseteq [M]$ 基于 \mathcal{A} 可以制造出一个海盗版 $\mathcal{A}(\text{Enc}(S)) \in \Upsilon^n$, 其中 Υ 为一个集合, 通常包含 Γ , 且 $\mathcal{A}(\text{Enc}(S))$ 可以不唯一;
- (3) $\text{Tra} : \Upsilon^n \rightarrow 2^{[M]}$ 是一个追踪算法, 即对于任何的至多 t 个合谋者 $S \subseteq [M]$ 以及海盗版 $\mathcal{A}(\text{Enc}(S)) \in \Upsilon^n$, 满足 $\text{Tra}(\mathcal{A}(\text{Enc}(S)))$ 不为空集, 且 $\text{Tra}(\mathcal{A}(\text{Enc}(S))) \subseteq S$.

上述可追踪分配方案可以保证至少追踪到部分或者全部的合谋攻击者, 而接下来介绍的防诬陷分配方案将考虑一种相对较弱的安全性质, 即它一般不具有可追踪能力, 但是可以保证无辜的用户或者团体不被合谋攻击集团所构陷.

定义 1.2 (防诬陷分配方案) 一个拥有 M 个用户 $[M]$ 且具有防诬陷性质的分配方案是一个四元组 $(\text{Enc}; \mathcal{A}, t, s)$, 其中,

- (1) $\text{Enc} : [M] \rightarrow \Gamma^n$ 是一个密钥/指纹分配方案, 其中 Γ 为一个有限集合;
- (2) $\mathcal{A} : 2^{\Gamma^n} \rightarrow \Upsilon^n$ 是一个攻击策略, 即合谋团体 $S \subseteq [M]$ 基于 \mathcal{A} 可以制造出一个海盗版 $\mathcal{A}(\text{Enc}(S)) \in \Upsilon^n$, 其中 Υ 为一个集合, 通常包含 Γ , 且 $\mathcal{A}(\text{Enc}(S))$ 可以不唯一;
- (3) Enc 具有防诬陷性质, 即对于任何不相交用户集合 $S, S' \subseteq [M]$ 且 $|S| \leq t, |S'| \leq s$, 以及他们能够制造的任何海盗版 $\mathcal{A}(\text{Enc}(S)), \mathcal{A}(\text{Enc}(S'))$, 都有 $\mathcal{A}(\text{Enc}(S)) \neq \mathcal{A}(\text{Enc}(S'))$.

从上述两个定义的描述中可以看出, 一个可追踪分配方案一定具有防诬陷的性质, 然而反之未必成立. 组合安全码研究的核心问题如下:

问题 1.1 设计可以容纳尽可能多的用户的可追踪分配方案和防诬陷分配方案, 以及相应的高效快速的追踪/译码算法.

在下面的第 1.4 和 1.5 小节中, 将参照定义 1.1 和 1.2 中的统一模型, 分别对广播加密和多媒体指纹识别的应用场景进行介绍. 为了叙述方便, 首先列举本文使用的一些记号.

1.3 一些记号

本文采用以下数学符号和定义:

- $[n] := \{1, 2, \dots, n\}$, $2^{[n]} := \{X : X \subseteq [n]\}$ 表示 $[n]$ 的幂集, $\binom{[n]}{w} := \{B \subseteq [n] : |B| = w\}$.
- \mathbb{R}^n 表示 n 维 Euclid 空间, \mathbb{R}^n 中任意两个点 \mathbf{z} 与 \mathbf{z}' 之间的 Euclid 距离为 $\|\mathbf{z} - \mathbf{z}'\|$.
- 令 n 和 q 为正整数, $Q = \{0, 1, \dots, q-1\}$, $Q^n = \{\mathbf{x} = (x(1), x(2), \dots, x(n)) : x(i) \in Q, 1 \leq i \leq n\}$ 为所有基于 Q 且长度为 n 的向量所构成的集合.
 - 对于 $\mathbf{x} = (x(1), x(2), \dots, x(n)) \in Q^n$ 和 $\mathbf{y} = (y(1), y(2), \dots, y(n)) \in Q^n$, \mathbf{x} 与 \mathbf{y} 之间的 Hamming 距离定义为 $d(\mathbf{x}, \mathbf{y}) := |\{i \in [n] : x(i) \neq y(i)\}|$, \mathbf{x} 与 \mathbf{y} 之间的内积为 $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x(i)y(i)$.
 - 用 $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\} \subseteq Q^n$ 表示一个基于 Q 的码长为 n 、大小为 M 的 q 元码, 记为 (n, M, q) 码, \mathcal{C} 中的每一个 $\mathbf{c}_j = (c_j(1), c_j(2), \dots, c_j(n))$ 称为码字.
 - 用 $\mathcal{B} = \{B_1, B_2, \dots, B_M\} \subseteq \binom{[n]}{w}$ 表示一个 (w, M, n) 集合系, 其中每一个 $B_j \subseteq [n]$ 被称为区组 (block), 且区组大小为 $|B_j| = w$.
 - 对于任意的码 $\mathcal{C} \subseteq Q^n$, 记 $\mathcal{C}(i) := \{c(i) \in Q : \mathbf{c} = (c(1), c(2), \dots, c(n)) \in \mathcal{C}\}$ ($1 \leq i \leq n$) 为码 \mathcal{C} 的第 i 个分量的元素所构成的集合.
 - 对于任意的码 $\mathcal{C} \subseteq Q^n$, 定义 \mathcal{C} 的后代码 (descendant code) 为

$$\text{desc}(\mathcal{C}) := \mathcal{C}(1) \times \mathcal{C}(2) \times \dots \times \mathcal{C}(n) = \{(x(1), x(2), \dots, x(n)) \in Q^n : x(i) \in \mathcal{C}(i), 1 \leq i \leq n\}.$$
 - 令 $t \geq 2$ 为正整数. 对于任意 $\mathbf{w} \in Q^n$, 定义 \mathbf{w} 在 \mathcal{C} 中的限定父代码集 (set of bounded parent codes) 为

$$\mathcal{P}_t(\mathbf{w}) := \{\mathcal{C}' \subseteq \mathcal{C} : \mathbf{w} \in \text{desc}(\mathcal{C}'), |\mathcal{C}'| \leq t\},$$

其中, $\mathcal{P}_t(\mathbf{w})$ 中的每一个 \mathcal{C}' 都是 \mathbf{w} 在 \mathcal{C} 中的一个限定父代码.

- 令 $t \geq 2$ 为正整数. 对于任意 $S \subseteq Q^n$, 定义 S 在 \mathcal{C} 中的限定父代码集和父代码集 (set of parent codes) 分别为

$$\mathcal{P}_t(S) := \{\mathcal{C}' \subseteq \mathcal{C} : \text{desc}(\mathcal{C}') = S, |\mathcal{C}'| \leq t\}, \quad \mathcal{P}(S) := \{\mathcal{C}' \subseteq \mathcal{C} : \text{desc}(\mathcal{C}') = S\}, \quad (1.1)$$

其中, $\mathcal{P}_t(S)$ 中的每一个 \mathcal{C}' 是 \mathcal{C} 中能产生 S 的一个限定父代码, $\mathcal{P}(S)$ 中的每一个 \mathcal{C}' 是 \mathcal{C} 中能产生 S 的一个父代码.

- 对于一个 (n, M, q) 码 \mathcal{C} , 定义 \mathcal{C} 的码率 (code rate) 为 $(\log_q M)/n$. 记 $M(n, q)$ 为 n 长的 q 元码的最大码字个数, 定义其最大渐近码率 (largest asymptotic code rate) 为

$$R(q) := \limsup_{n \rightarrow \infty} \frac{\log_q M(n, q)}{n}.$$

- 对于两个关于正整数 n 取正值的函数 f 和 g , 如果 $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$, 则记 $f = o(g)$; 如果存在常数 c 和 N , 使得对于任意 $n \geq N$ 都有 $f(n) \leq cg(n)$, 则记 $f = O(g)$; 如果存在常数 c 和 N , 使得对于任意 $n \geq N$ 都有 $f(n) \geq cg(n)$, 则记 $f = \Omega(g)$; 如果 $f = O(g)$ 与 $f = \Omega(g)$ 同时成立, 则记 $f = \Theta(g)$.

- $\text{Poly}(n)$ 表示关于变量 n 取正值的多项式函数.

为避免记号混淆, 需要指出 $\mathcal{P}_t(\mathbf{w})$ 与 $\mathcal{P}_t(\{\mathbf{w}\})$ 之间的区别: $\mathcal{P}_t(\mathbf{w})$ 是后代码包含 \mathbf{w} 的 \mathcal{C} 的大小不超过 t 的子码的集合, 而 $\mathcal{P}_t(\{\mathbf{w}\})$ 是后代码恰好是 $\{\mathbf{w}\}$ 的 \mathcal{C} 的大小不超过 t 的子码的集合, 它们之间的关系为 $\mathcal{P}_t(\{\mathbf{w}\}) \subseteq \mathcal{P}_t(\mathbf{w})$.

1.4 广播加密中的版权保护

在广播加密的应用场景中, 可追踪分配方案可以分为针对内容和针对密钥的加密^[83]. 下面对这两种情形进行具体说明.

1.4.1 针对内容加密的模型

Chor 等^[18] 可追踪分配方案考虑了对数据内容的加密. 如图 1 所示, 首先, 原始数据内容被划分成 n 块, 即 $\text{data}_1, \text{data}_2, \dots, \text{data}_n$; 然后, 数据发布者利用均匀且独立地随机 (uniformly and independently at random) 选取的 qn 个密钥 $\{K_{i,j} : i \in Q, j \in [n]\}$, 对这 n 块数据进行加密并广播, 其中每块数据 data_j 被加密成 q 个备份: $\{E_{K_{i,j}}(\text{data}_j) : i \in Q\}$. 如果用户想要恢复原始数据内容, 需要解密得到所有的 n 块数据. 每个合法用户 $l \in [M]$ 在付费后将分配到一个码字 $\mathbf{c}_l = (c_l(1), c_l(2), \dots, c_l(n)) \in Q^n$, 对应到一个人解密密钥 $(K_{c_l(1),1}, K_{c_l(2),2}, \dots, K_{c_l(n),n})$, 其中 $K_{c_l(j),j}$ 可以用来解密 $E_{K_{c_l(j),j}}(\text{data}_j)$ 得到第 j 块数据 data_j (为了叙述方便, 这里假设使用的是对称密钥加密, 即加解密的密钥是相同的).

此时, 在定义 1.1 的模型下, 密钥分配方案可以表述为 $\text{Enc} : [M] \rightarrow Q^n$. 若不超过 t 个用户 $S \subseteq [M]$ 进行合谋攻击, 则根据嵌入假设 (marking assumption)^[11], 他们可以制造出一个海盗版解密密钥 $(K_{x(1),1}, K_{x(2),2}, \dots, K_{x(n),n})$, 对应于

$$\mathcal{A}(\text{Enc}(S)) \in \{(x(1), x(2), \dots, x(n)) \in Q^n : x(j) = c_s(j), s \in S\}.$$

换句话说, 海盗版解密密钥中用于对每块数据 data_j 解密的密钥必须由合谋团体 S 中的某一个合谋者来提供. 在此攻击模型下, 有如下两种经典的追踪算法:

(T1) Chor 等^[18, 19] 提出了比较海盗版解密密钥 $\mathcal{A}(\text{Enc}(S))$ 和所有合法用户解密密钥的 Hamming 距离的方法来追踪合谋者, 即

$$\text{Tra}(\mathcal{A}(\text{Enc}(S))) = \left\{ l \in [M] : d(\mathbf{c}_l, \mathcal{A}(\text{Enc}(S))) = \min_{k \in [M]} d(\mathbf{c}_k, \mathcal{A}(\text{Enc}(S))) \right\}.$$

显然 $\text{Tra}(\mathcal{A}(\text{Enc}(S)))$ 不是空集. 若要进一步保证此追踪算法的输出一定是合谋者, 则设计的密钥分配方案需满足 Chor-Fiat-Naor 可追踪性质, 即构造可追踪码 (traceability code, TAC, 参见定义 2.1).

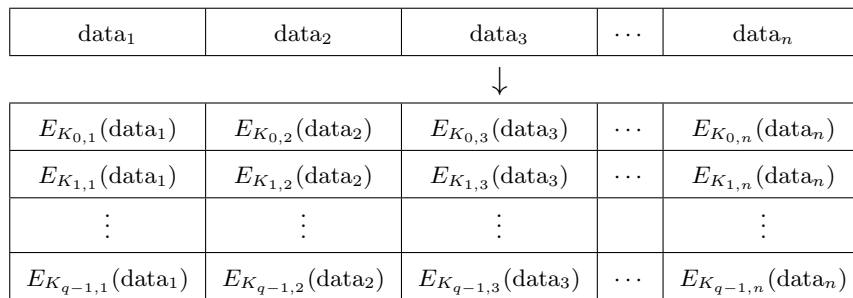


图 1 针对内容加密的模型

(T2) 上述 Chor-Fiat-Naor 可追踪码, 当合谋攻击者人数不超过 t 时, 可以快速识别出至少一个合谋者, 但是可追踪码的性质较强, 在码长和基础集的大小确定的情形下, 很难构造出具有较大码字个数的可追踪码, 因此很难应用于实际授权用户人数较多的情形. 为了解决这个难题, Hollmann 等^[54] 提出了父代可识别码 (identifiable-parent-property code, IPPC), 而后这个概念被 Staddon 等^[82] 推广至一般情形 (参见定义 2.1). 基于父代可识别码的追踪算法的思想是, 首先, 根据海盗版 $\mathcal{A}(\text{Enc}(S))$ 查找其所有的限定父代码; 然后, 对识别出的所有的限定父代码求交集, 作为追踪算法的输出. 基于父代可识别码的密钥分配方案可以保证该追踪算法的输出不为空集, 且所有输出一定是制造出海盗版 $\mathcal{A}(\text{Enc}(S))$ 的真实的合谋攻击者. 相较于可追踪码, 一个父代可识别码可应用于具有更多授权用户的系统, 而且在合谋攻击者人数不超过 t 时同样能够识别出至少一个合谋攻击者, 但是在一般情形下其追踪速率不如可追踪码的快.

除此以外, 对于防诬陷分配方案也有两类较为经典的组合安全码被提出. 1998 年, Boneh 和 Shaw^[11] 提出了防诬陷码 (frameproof code, FPC) 的概念 (参见定义 2.1), 证明了 t -防诬陷码在数字指纹识别中虽不具有追踪任何一个合谋攻击者的能力, 但是能够防止一个没有参与合谋攻击的授权用户被其他人不超过 t 的合谋攻击用户所构陷, 因此在数字指纹识别中 t -防诬陷码能够提供一定的安全性. 在防诬陷码的研究基础上, Stinson 等^[84] 提出了 t -安全防诬陷码 (secure frameproof code, SFPC) 的概念 (参见定义 2.1), 并证明了 t -安全防诬陷码和 t -防诬陷码一样不能追踪到任何一个合谋攻击者, 但是能够提供比 t -防诬陷码更强的安全性. 最初, 防诬陷码和安全防诬陷码的概念是以二元码的形式被提出和研究的, 后来 Staddon 等^[82] 将其一般化并开启了对于 q 元防诬陷码和安全防诬陷码的研究.

1.4.2 针对密钥加密的模型

1998 年, Stinson 和 Wei^[85] 对 Chor-Fiat-Naor 模型进行了扩展. 如图 2 所示, 在广播加密系统中, Stinson-Wei 模型考虑利用会话密钥 (session key) K 对数据 Data 的整体进行加密, 得到密文 $E_K(\text{Data})$; 并利用 (n, w) 阈值秘密共享方案 (threshold secret sharing scheme)^[75] 对密钥 K 进行加密形成一个短的能动区块 (enabling block), 该区块包含由 n 个基本密钥 (base key) 对 K 共享加密生成的 n 个密文. 数据发行商将密文 $E_K(\text{Data})$ 和能动区块一起广播. 一个合法的用户会被分配到 w 个基本密钥, 用来解密能动区块中的 w 个密文, 根据 (n, w) 阈值秘密共享方案的原理, 可进一步恢复会话密钥 K , 从而能够对密文 $E_K(\text{Data})$ 进行解密得到原始数据 Data. (这里同样假设使用对称密钥对原始数据 Data 进行加密, 其加解密的密钥是相同的.)

在此模型下, 密钥分配方案可以表述为 $\text{Enc} : [M] \rightarrow \{0, 1\}^n$, 其中用户 l 对应的 $\text{Enc}(l)$ 是一个 Hamming 重量为 w 的 0-1 向量, 记为 $\mathbf{c}_l \in \{0, 1\}^n$. 值得一提的是, \mathbf{c}_l 可以等价表示为一个大小为 w 的区组, 其中区组的元素是该向量的非零位置的指标. 此时, 若有不超过 t 个用户的团体 $S \subseteq [M]$ 合谋攻击, 则他们可以制造出一个海盗版解密密钥

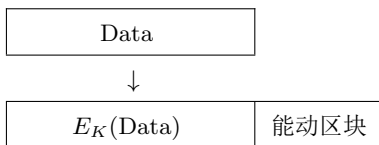


图 2 针对会话密钥加密的模型

$$\mathcal{A}(\text{Enc}(S)) \in \left\{ (x(1), (x(2), \dots, x(n)) \in \{0, 1\}^n : \sum_{j=1}^n x(j) = w \text{ 且 } x(j) \leq \max_{s \in S} c_s(j) \right\},$$

并利用 (n, w) 阈值秘密共享方案来进行解密得到会话密钥 K . 这里值得注意的是, 用户在使用其被分配到的基本密钥来解密 K 并进而解密原始数据的过程是在一个黑盒下完成的, 即用户不能看到中间的会话密钥 K , 因而合谋者不能直接得到 K , 其合谋攻击只能是利用其拥有的基本密钥来实现. 类似于第 1.4.1 小节中的追踪算法的思路, 也有如下两类经典的组合结构被提出:

(1) Stinson 和 Wei^[85] 提出了利用基于区组的可追踪分配方案 (traceability scheme, TS, 参见定义 2.2) 来进行密钥分配方案的设计, 其对应的追踪算法是通过找到与海盗版 $\mathcal{A}(\text{Enc}(S))$ 的 Hamming 距离最近的合法解密密钥及其对应的用户, 且能够保证这些被追踪到的用户一定是真正的合谋攻击者.

(2) Collins^[22] 提出了父代可识别集合系 (identifiable-parent-property set system, IPPS, 参见定义 2.2) 的概念, 其可以看作是第 1.4.1 小节中的父代可识别码的扩展形式. 基于父代可识别集合系来设计的密钥分配方案也可以保证能追踪到至少一个真实的合谋攻击者, 其中追踪算法的思想与第 1.4.1 小节中的 (T2) 类似, 即, 首先根据海盗版 $\mathcal{A}(\text{Enc}(S))$ 查找其所有的大小不超过 t 的可能的父代集合, 然后对识别出的这些父代集合求交集, 作为追踪算法的输出. 父代可识别集合系可以保证该追踪算法的输出不为空集且输出一定是制造海盗版 $\mathcal{A}(\text{Enc}(S))$ 的真实的合谋攻击者. 相比于 Stinson 和 Wei^[85] 的可追踪分配方案, 父代可识别集合系可应用于具有更多授权用户的系统, 但其追踪算法的速率一般来说更慢一些.

在针对密钥加密的模型中, 除了上述两种具有追踪性质的组合结构以外, 也有对于防诬陷分配方案的组合结构的相关研究. Kautz 和 Singleton^[60] 在 1964 年提出的叠加码 (superimposed code) 可以用来设计密钥分配方案, 且能够保证任何用户不会被其他不超过 t 个合谋用户所构陷. 到目前为止, 叠加码在通信、信息安全、组合学等不同的研究领域均有丰富的研究和应用. 在组合学中, 叠加码又常被以无覆盖族 (cover-free family, CFF, 参见定义 2.2) 的名称来研究. 在此基础上, Gu 和 Miao^[49] 提出了并交限定族 (union-intersection-bounded family, UIBF, 参见定义 2.2) 的概念, 其可以看作是无覆盖族的一般推广. 当 d 不超过某个阈值时, 基于 $(s, t; d)$ 并交限定族来设计得到的密钥分配方案, 可以保证任何 s 个用户都不会被其他 t 个合谋攻击者所诬陷. 这两类具有防诬陷性质的组合集合族都不具有可追踪的性质.

1.5 多媒体内容的版权保护

多媒体指纹识别是为了保护多媒体内容版权而发展起来的一种技术 (参见文献 [65, 88]). 在向 M 个授权用户发送多媒体内容 $\mathbf{x} \in \mathbb{R}^m$ 前, 版权所有者首先为 M 个用户设计一个指纹分配方案 $\text{Enc} : [M] \rightarrow \Gamma^n$, 其中 $\text{Enc}(l)$ 是针对用户 l 的指纹分配方案, 记为 $\mathbf{c}_l = (c_l(1), c_l(2), \dots, c_l(n)) \in \Gamma^n$, 然后使用基于 n ($\leq m$) 个类似噪声的标准正交信号 (noise-like orthonormal signals) $\mathcal{F} = \{\mathbf{f}_i \in \mathbb{R}^m : 1 \leq i \leq n\}$ 的线性调制方案构造一组水印, 也就是指纹 (fingerprints), 并将其嵌入内容中, 再发布给授权用户. 用户 l 实际接收到的多媒体内容为 $\mathbf{y}_l = \mathbf{x} + \mathbf{w}_l$, 其中 $\mathbf{w}_l = \sum_{i=1}^n c_l(i) \mathbf{f}_i$ 是分配给用户 l 的指纹, 用户自己并不知晓. 当 $\Gamma = \{0, 1\}$ 时, 指纹分配方案的设计问题与组合学有着紧密联系, 且得到了广泛的研究. 根据多媒体嵌入假设 (multimedia marking assumption)^[11, 34] 用户无法操控和修改 \mathcal{F} , 且合谋攻击仅限于线性攻击 (linear attack). 假设不超过 t 个用户 $S \subseteq [M]$ 进行合谋, 在线性攻击模型 \mathcal{A} 下伪造出来一个海盗版

$$\mathcal{A}(\text{Enc}(S)) \in \left\{ \sum_{s \in S} \lambda_s \mathbf{y}_s \in \mathbb{R}^m : \lambda_s > 0 \text{ 且 } \sum_{s \in S} \lambda_s = 1 \right\},$$

其中 λ_s 代表合谋用户 $s \in S$ 在线性攻击中的“贡献”, 由合谋用户来决定, 每个合谋用户的贡献要大于 0, 并且为了海盗版能够正常使用, 这里要求 $\sum_{s \in S} \lambda_s = 1$. 当 λ_s 都等于 $1/|S|$ 时, \mathcal{A} 称为平均攻击 (averaging attack) [88]. 在线性攻击模型下, 海盗版内容 $\mathcal{A}(\text{Enc}(S))$ 中也会包含一个指纹, 也就是 $\mathcal{A}(\text{Enc}(S)) - \mathbf{x}$. 版权所有者在获取海盗版之后可以从中提取出指纹, 并通过与 \mathcal{F} 中的正交信号作内积计算得到 $\mathbf{r} = (r(1), r(2), \dots, r(n)) \in \mathbb{R}^n$, 其中 $r(i) := \langle \mathcal{A}(\text{Enc}(S)) - \mathbf{x}, \mathbf{f}_i \rangle$. 根据线性攻击的特点, 可以计算得到 $\mathbf{r} = \sum_{s \in S} \lambda_s \mathbf{c}_s$.

最开始, 人们在设计追踪算法时都是考虑利用 \mathbf{r} 离散化后的结果 $\mathcal{D}_{\mathbf{r}} = \mathcal{D}(1) \times \mathcal{D}(2) \times \dots \times \mathcal{D}(n)$ 来追踪, 其中, $\mathcal{D}(i) = \{0, 1\}$ ($0 < r(i) < 1$), 否则 $\mathcal{D}(i) = \{r(i)\}$. 不难发现, $\mathcal{D}_{\mathbf{r}}$ 即为集合 $\{\mathbf{c}_s : s \in S\}$ 的后代码. 通常来看, 具有可追踪性质的指纹分配方案主要有以下两大类:

(1) 对于完全可追踪的指纹分配方案, Cheng 和 Miao [17] 在 2011 年首先提出了可分离码 (separable code, SC) 的概念 (参见定义 2.3), 其追踪算法的思想是根据 $\mathcal{D}_{\mathbf{r}}$ 找出其限定父代码, 可分离码的性质能够保证算法的输出一定是唯一的, 且该输出对应的是合谋用户集 S . 值得一提的是, 根据第 1.4.1 小节, 防诬陷码在针对内容加密的模型中并不具备可追踪能力. 然而, Cheng 和 Miao [17] 证明了防诬陷码在多媒体指纹识别中具有完全可追踪的能力, 其追踪算法的过程是通过对比每个授权用户的指纹分配方案和 $\mathcal{D}_{\mathbf{r}}$, 找出所有满足 $\mathbf{c}_s \in \mathcal{D}_{\mathbf{r}}$ 的授权用户 s , 防诬陷码的性质可以保证该算法输出一定是所有参与合谋攻击的用户. 一般而言, 防诬陷码的追踪算法的速率远高于可分离码, 但是可分离码的性质比防诬陷码的更弱, 因此, 可分离码更能适用于具有较多授权用户的系统. 为了融合可分离码和防诬陷码这两类指纹分配方案的优点, Jiang 等 [56] 提出了强可分离码 (strongly separable code, SSC) 的概念 (参见定义 2.3). 其追踪算法的过程分为两步, 第一步, 与防诬陷码一样, 找出所有满足 $\mathbf{c}_l \in \mathcal{D}_{\mathbf{r}}$ 的授权用户 l , 不同于防诬陷码的是, 这里找出的是一个包含合谋攻击用户集 S 的更大的集合; 第二步, 再从这个集合中找出参与合谋攻击的用户集 S , 详细算法过程参见第 2.2 小节. 最近, 基于两步追踪的想法, Gu 等 [51] 提出了带有序列译码特性的安全码 (secure codes with list decoding) 的概念, 并证明了它可以同时具有较快速的追踪/译码算法和较大的码率.

(2) 对于不完全可追踪的指纹分配方案, Cheng 等 [15] 提出了多媒体父代可识别码 (multimedia identifiable-parent-property code, MIPPC, 参见定义 2.3) 的概念, 其追踪算法的思想与第 1.4 小节中的两类具有父代可识别性质的组合安全码相似, 因此其追踪算法的速率也不高. 为了提高追踪算法的速率, Jiang 等 [57] 基于强可分离码的追踪算法思想提出了一种具有更强性质的多媒体父代可识别码 (strongly multimedia identifiable-parent-property code, SMIPPC, 参见定义 2.3), 其追踪算法的过程和强可分离码的相似, 不同的是, 这里第一步找出的是包含至少一个合谋攻击用户 (不一定包含所有合谋攻击用户) 的授权用户集合, 第二步再从该集合中找出合谋攻击用户 (参见 2.2 小节).

2019 年, Egorova 等 [34] 提出了签名码 (signature code) 的概念 (参见定义 2.3). 其追踪算法的思想与可分离码的相似, 因而追踪算法的速率相同, 不同的是, 利用签名码追踪时所依据的是 \mathbf{r} 而不是离散化的结果 $\mathcal{D}_{\mathbf{r}}$ (参见 2.2 小节). 相较于可分离码, 签名码的性质更弱, 能够适用于具有更多授权用户的系统.

1.6 相关领域问题

版权保护中的组合安全码一直是组合学与信息科学这两个领域的交叉研究热点, 多类组合安全码

不仅与集合论、组合设计理论、图论中的某些结构等具有高度相似性, 而且概率方法、代数方法、极值组合方法和组合设计方法等也常被用来研究组合安全码的最大码率和构造问题. 不仅如此, 组合安全码与编码理论、信息安全和理论计算机科学等领域的诸多问题也有着紧密的联系. 第 5 节将就群试理论 (group testing) 和多用户通信系统 (multi-user communication system) 应用中的相关组合学问题以及它们与版权保护中的组合安全码的研究之间的联系进行阐述说明.

1.7 论文结构

本文余下内容的结构框架如下: 第 2 节介绍组合安全码的数学定义及其译码算法与码之间的关系等. 第 3 和 4 节分别介绍关于组合安全码最大码率上界和下界的研究方法. 第 5 节介绍组合安全码与群试理论和多用户通信等领域的联系. 第 6 节总结全文.

2 组合安全码及其译码算法

本节主要介绍广播加密和多媒体指纹识别中对抗合谋攻击的组合安全码及其追踪/译码算法的思想. 其中, 第 2.1 小节将分别介绍广播加密中的组合安全码和集合系, 以及它们在对抗合谋攻击时所发挥的作用; 第 2.2 小节将介绍多媒体指纹识别中的几类多媒体指纹码及其追踪/译码算法思想; 第 2.3 小节提出与组合安全码的设计相关的公开问题.

2.1 广播加密中的组合安全码

本小节主要介绍可追踪码、父代可识别码、防诬陷码和安全防诬陷码这几类抗合谋攻击的组合安全码, 并讨论它们之间的关系以及它们在实际对抗合谋攻击时各自所发挥的作用.

2.1.1 针对内容加密的组合安全码

本小节给出几类组合安全码的定义, 如可追踪码^[9, 18, 89]、父代可识别码^[54, 82]、防诬陷码^[6, 11]和安全防诬陷码^[84].

定义 2.1 (广播加密中的组合安全码) 设 \mathcal{C} 是一个 (n, M, q) 码, 对于任意 $\mathcal{C}_0 \subseteq \mathcal{C}$, $|\mathcal{C}_0| \leq t$ 和任意 $\mathbf{w} \in \text{desc}(\mathcal{C}_0)$,

(i) 称 \mathcal{C} 为 t -可追踪码 (traceability code, TAC), 记作 $t\text{-TAC}(n, M, q)$, 如果存在 $\mathbf{y} \in \mathcal{C}_0$ 使得对于任意 $\mathbf{z} \in \mathcal{C} \setminus \mathcal{C}_0$ 都有

$$d(\mathbf{w}, \mathbf{y}) < d(\mathbf{w}, \mathbf{z}); \quad (2.1)$$

(ii) 称 \mathcal{C} 为 t -父代可识别码 (identifiable-parent-property code, IPPC), 记作 $t\text{-IPPC}(n, M, q)$, 若

$$\bigcap_{\mathcal{C}' \in \mathcal{P}_t(\mathbf{w})} \mathcal{C}' \neq \emptyset; \quad (2.2)$$

(iii) 称 \mathcal{C} 为 t -安全防诬陷码 (secure frameproof code, SFPC), 记作 $t\text{-SFPC}(n, M, q)$, 若对于任意 $\mathcal{C}_1 \subseteq \mathcal{C}$, $|\mathcal{C}_1| \leq t$ 且 $\mathcal{C}_0 \cap \mathcal{C}_1 = \emptyset$ 都有

$$\text{desc}(\mathcal{C}_0) \cap \text{desc}(\mathcal{C}_1) = \emptyset; \quad (2.3)$$

(iv) 称 \mathcal{C} 为 t -防诬陷码 (frameproof code, FPC), 记作 $t\text{-FPC}(n, M, q)$, 若

$$\text{desc}(\mathcal{C}_0) \cap \mathcal{C} = \mathcal{C}_0. \tag{2.4}$$

例 2.1 令

$$\mathcal{C} = \begin{matrix} & \mathbf{c}_1 & \mathbf{c}_2 & \mathbf{c}_3 & \mathbf{c}_4 \\ \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \end{matrix}$$

为一个 $(3, 4, 2)$ 码. 根据定义 2.1, 可以证明:

(1) \mathcal{C} 是一个 2-SFPC. 例如, 对于 \mathcal{C} 的两个不交的子集 $\mathcal{C}_0 = \{\mathbf{c}_1, \mathbf{c}_2\}$ 和 $\mathcal{C}_1 = \{\mathbf{c}_3, \mathbf{c}_4\}$, 不难验证 $\text{desc}(\mathcal{C}_0) = \{0, 1\} \times \{0, 1\} \times \{0\} = \{\mathbf{c}_1, \mathbf{c}_2, (0, 0, 0), (1, 1, 0)\}$ 和 $\text{desc}(\mathcal{C}_1) = \{0, 1\} \times \{0, 1\} \times \{1\} = \{\mathbf{c}_3, \mathbf{c}_4, (0, 1, 1), (1, 0, 1)\}$, 因此 $\text{desc}(\mathcal{C}_0) \cap \text{desc}(\mathcal{C}_1) = \emptyset$.

(2) \mathcal{C} 也是一个 2-FPC. 例如, 对于 (1) 中的 \mathcal{C}_0 和 \mathcal{C}_1 , 有 $\text{desc}(\mathcal{C}_0) \cap \mathcal{C} = \mathcal{C}_0$ 且 $\text{desc}(\mathcal{C}_1) \cap \mathcal{C} = \mathcal{C}_1$.

(3) \mathcal{C} 不是 2-IPPC. 这是因为对于 (1) 中的 \mathcal{C}_0 和 $\mathbf{w} = (1, 1, 0) \in \text{desc}(\mathcal{C}_0)$, 可以得到 \mathbf{w} 在 \mathcal{C} 中的限定父代码集 $\mathcal{P}_2(\mathbf{w}) = \{\mathcal{C}_0, \{\mathbf{c}_1, \mathbf{c}_4\}, \{\mathbf{c}_2, \mathbf{c}_4\}\}$, 但是 $\mathcal{C}_0 \cap \{\mathbf{c}_1, \mathbf{c}_4\} \cap \{\mathbf{c}_2, \mathbf{c}_4\} = \emptyset$.

(4) \mathcal{C} 也不是 2-TAC. 这是因为对于 (3) 中的 \mathcal{C}_0 以及 \mathbf{w} , 存在 $\mathbf{c}_4 \in \mathcal{C} \setminus \mathcal{C}_0$, 使得 $d(\mathbf{w}, \mathbf{c}_4) = 1 = d(\mathbf{w}, \mathbf{c}_1) = d(\mathbf{w}, \mathbf{c}_2)$, 不满足 2-TAC 的条件.

由定义 2.1 不难证明 TAC、IPPC、SFPC 和 FPC 之间的关系, 感兴趣的读者可参见文献 [82, 引理 1.3, 例 1.4] 和 [84, 定理 2.2].

定理 2.1 设 \mathcal{C} 是一个 (n, M, q) 码, 则有如下结论:

- (1) (参见文献 [82, 引理 1.3 和例 1.4]) 如果 \mathcal{C} 是 t -TAC, 则 \mathcal{C} 也是 t -IPPC; 反之, 则不一定成立.
- (2) 如果 \mathcal{C} 是 t -IPPC, 则 \mathcal{C} 也是 t -SFPC; 反之, 则不一定成立.
- (3) (参见文献 [84, 定理 2.2]) 如果 \mathcal{C} 是 t -SFPC, 则 \mathcal{C} 也是 t -FPC; 反之, 则不一定成立.

根据定理 2.1, 利用图 3 简要刻画 TAC、IPPC、SFPC 和 FPC 之间的关系, 其中 $A \Rightarrow B$ 表示若一个码 \mathcal{C} 具有 A 的性质, 则 \mathcal{C} 也具有 B 的性质. 因此, 容易推出以下结论, 其中 SSC、SC、SMIPPC、MIPPC 和签名码的概念以及它们之间的关系将在第 2.2 小节中介绍.

推论 2.1 设 $R_{\text{TAC}}(q, t)$ 、 $R_{\text{IPPC}}(q, t)$ 、 $R_{\text{SFPC}}(q, t)$ 和 $R_{\text{FPC}}(q, t)$ 分别为 q 元 t -TAC、 t -IPPC、 t -SFPC 和 t -FPC 的最大渐近码率, 则

$$R_{\text{TAC}}(q, t) \leq R_{\text{IPPC}}(q, t) \leq R_{\text{SFPC}}(q, t) \leq R_{\text{FPC}}(q, t).$$

当将定义 2.1 中的这几类组合安全码实际应用于数字指纹识别时, 一个 (n, M, q) 码 \mathcal{C} 中的 M 个码字代表被分配给 M 个授权用户的数字媒体中的指纹. 假设有不超过 t 个授权用户 (记为集合 \mathcal{C}_0) 非

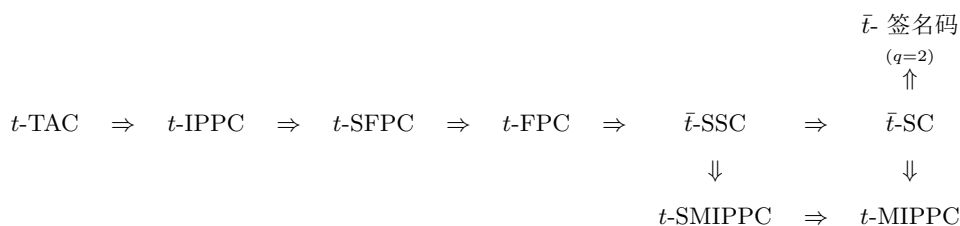


图 3 组合安全码之间的关系

合法复制数字媒体, 则 $\text{desc}(C_0)$ 代表 C_0 中的合谋攻击者在违法复制品中可能植入的所有数字指纹的集合 (即所有可能的盗版数字指纹的集合). 若合谋攻击者在违法复制品中植入的指纹是 $\mathbf{w} \in Q^n$, 则 $\mathcal{P}_t(\mathbf{w})$ 中每一个 C' 是 \mathbf{w} 的一个限定父代码, 代表了 C 中能产生 \mathbf{w} 且人数不超过 t 的授权用户的集合. 显然, $C_0 \in \mathcal{P}_t(\mathbf{w})$. 我们的目的是利用码 C 和 \mathbf{w} 来追踪 C_0 中的成员. 根据定义 2.1, 有如下结论:

- 如果 C 是 t -TAC, 当发行商从盗版产品中提取出数字指纹 \mathbf{w} 时, 可以对比 \mathbf{w} 和每一个授权用户的数字指纹之间的 Hamming 距离, 由 (2.1) 可以确定 C 中与 \mathbf{w} 距离最小的码字 \mathbf{y} 所对应的授权用户一定是参与违法复制的合谋者之一, 其译码算法的时间复杂度为 $O(nM)$.
- 如果 C 是 t -IPPC, 通过检查 C 的所有大小 $\leq t$ 的子码可以得到 \mathbf{w} 的限定父代码集 $\mathcal{P}_t(\mathbf{w})$. 由 (2.2) 可以确定 $\mathcal{P}_t(\mathbf{w})$ 中所有的 C' 之间必有非空交集且这个交集也是 C_0 的子集, 因此能够追踪到至少一个参与违法复制的合谋攻击者, 此追踪算法的时间复杂度是 $O(nM^t)$.
- 如果 C 是 t -SFPC, 由 (2.3) 可知, 对于任何一个人数不超过 t 且未参与合谋攻击的授权用户集 C_1 , 都有 $C_1 \notin \mathcal{P}_t(\mathbf{w})$, 这意味着 C_1 中的用户不会被 C_0 中的合谋攻击用户所构陷. 另外, 如果一个人数不超过 t 的授权用户集 C_2 能够产生指纹 \mathbf{w} , 即 $\mathbf{w} \in \text{desc}(C_2)$ 或 $C_2 \in \mathcal{P}_t(\mathbf{w})$, 则说明 C_2 中一定包含参与违法复制的合谋攻击用户, 但是合谋攻击用户的身份无法确定.
- 如果 C 是 t -FPC, 由 (2.4) 可知, 对于任意一位未参与合谋攻击的授权用户, 其数字指纹不可能被 C_0 中的合谋攻击者用户构造出来, 这就能够保证其不被合谋攻击者用户所构陷, 但是版权所有者无法根据从盗版数字媒体中提取出的 \mathbf{w} 来追踪任何一个合谋攻击者用户.

当参与违法复制的用户人数不超过 t 时, 上述说明展示了 t -TAC、 t -IPPC、 t -SFPC 和 t -FPC 在数字指纹识别中的追踪能力和译码算法时间. 需要特别注意的是, 这里 t -TAC 和 t -IPPC 的译码算法都是根据这两类码本身的定义而设计的, 虽然具有一般适用性, 但是当用户人数 M 非常大时, 它们的译码算法并不高效, 尤其是 t -IPPC. 很多已有的研究表明, 我们可以利用性质较好的纠错码 (error correcting code) 和代数几何结构来构造 TAC 和 IPPC, 虽然构造出来的 TAC 和 IPPC 其码率一般都比较小, 但是它们的译码算法却能得到极大的改进.

- Silverberg 等^[81] 利用 Reed-Solomon 码构造了一类特殊的 TAC, 并基于列表译码 (list decoding) 的思想 (参见文献 [52]) 为此类 TAC 开发了更高效的追踪算法, 将时间复杂度降低到 $O(\text{Ploy}(\log M))$.
- Barg 和 Kabatiansky^[5] 利用一个码率较大的 IPPC 和一个最小 Hamming 距离较大的线性码 (linear code) 串联构造出了一类新的 IPPC, 其码率虽然变小, 但是能在 $\text{Ploy}(n)$ 时间内追踪到至少一个合谋用户.
- Trung 和 Martirosyan^[90] 对 IPPC 进行多次递归构造, 并证明得到的 IPPC 其译码算法的时间复杂度为 $O(M)$.

更多关于 TAC 和 IPPC 译码算法方面的改进结果, 可参见文献 [7].

2.1.2 针对密钥加密的安全集合系

第 1.4.2 小节介绍了针对会话密钥加密的 Stinson-Wei 模型. 本小节介绍基于此模型的可用于设计可追踪密钥分配方案的几类组合安全集合系.

定义 2.2 (广播加密中的安全集合系) 设 \mathcal{B} 是一个 (w, M, n) 集合系, $d \leq w$ 是一个正整数, 对于任意 $B_0 \subseteq \mathcal{B}$, $|B_0| \leq t$ 和任意 $T \subseteq \bigcup_{B \in B_0} B$ 且 $|T| = w$,

- (i) 称 \mathcal{B} 为 t -可追踪分配方案 (traceability scheme, TS), 记作 t -TS(w, M, n), 如果存在 $B' \in B_0$ 使

得对于任意 $B \in \mathcal{B} \setminus \mathcal{B}_0$ 都有

$$|B \cap T| < |B' \cap T|; \tag{2.5}$$

(ii) 称 \mathcal{B} 为 t - 父代可识别集合系 (identifiable-parent-property set system, IPPS, 或者 parent-identifying set system), 记作 t -IPPS(w, M, n), 若

$$\bigcap_{B' \in \mathcal{P}_t(T)} B' \neq \emptyset, \tag{2.6}$$

其中 $\mathcal{P}_t(T)$ 为 T 在 \mathcal{B} 中的所有可能的限定父代集 (bounded parent sets) 的集合

$$\mathcal{P}_t(T) := \left\{ B' \subseteq \mathcal{B} : T \subseteq \bigcup_{B \in B'} B, |B'| \leq t \right\}; \tag{2.7}$$

(iii) 称 \mathcal{B} 为 $(s, t; d)$ - 并交限定族 (union-intersection-bounded family, UIBF), 记作 $(s, t; d)$ -UIBF(w, M, n), 若对于任意 $s \leq t$ 和任意 $s + t$ 个不同的区组 $A_1, \dots, A_s, B_1, \dots, B_t \in \mathcal{B}$ 都有

$$\left| \left(\bigcup_{1 \leq i \leq s} A_i \right) \cap \left(\bigcup_{1 \leq j \leq t} B_j \right) \right| < d; \tag{2.8}$$

(iv) 称 \mathcal{B} 为 t - 无覆盖族 (cover-free family, CFF), 记作 t -CFF(w, M, n), 若对于任意 $F \in \mathcal{B} \setminus \mathcal{B}_0$ 都有

$$F \not\subseteq \bigcup_{B \in \mathcal{B}_0} B. \tag{2.9}$$

例 2.2 令

$$\mathcal{B} = \{B_1 = \{1, 2, 3\}, B_2 = \{1, 4, 5\}, B_3 = \{3, 4, 6\}, B_4 = \{7, 8, 9\}\}$$

为一个 $(3, 4, 9)$ 集合系. 根据定义 2.2, 可以证明:

(1) \mathcal{B} 是一个 2 -CFF($3, 4, 9$). 这是因为任何一个区组都不会被其他两个区组的并集所覆盖. 例如, $B_3 = \{3, 4, 6\} \not\subseteq B_1 \cup B_2 = \{1, 2, 3, 4, 5\}$.

(2) \mathcal{B} 是一个 $(2, 2; 3)$ -UIBF($3, 4, 9$). 例如,

$$|(B_1 \cup B_2) \cap (B_3 \cup B_4)| = |(\{1, 2, 3\} \cup \{1, 4, 5\}) \cap (\{3, 4, 6\} \cup \{7, 8, 9\})| = |\{3, 4\}| = 2 < 3.$$

(3) \mathcal{B} 不是一个 2 -IPPS($3, 4, 9$). 这是因为对于 $\mathcal{B}_0 = \{B_1, B_2\} \subseteq \mathcal{B}$ 和 $T = \{1, 3, 4\} \subseteq (B_1 \cup B_2)$, 可以得到 T 在 \mathcal{B} 中所有可能的限定父代集的集合 $\mathcal{P}_2(T) = \{\mathcal{B}_0, \{B_1, B_3\}, \{B_2, B_3\}\}$, 但是 $\mathcal{B}_0 \cap \{B_1, B_3\} \cap \{B_2, B_3\} = \emptyset$.

(4) \mathcal{B} 不是一个 2 -TS($3, 4, 9$). 这是因为对于上述 (3) 中的 \mathcal{B}_0 和 T , 有 $|B_1 \cap T| = |B_2 \cap T| = |B_3 \cap T|$, 不满足 2 -TS 的条件.

由定义 2.2 可以推出 TS、IPPS、UIBF 和 CFF 之间的关系, 见如下定理 2.2 和图 4, 其中, $A \Rightarrow B$ 表示, 若一个集合系 \mathcal{B} 具有 A 的性质, 则 \mathcal{B} 也具有 B 的性质. 其具体证明和相关例子可参见文献 [45, 49].

$$t\text{-TS} \Rightarrow t\text{-IPPS} \Rightarrow (t, t; w)\text{-UIBF} \Rightarrow t\text{-CFF}$$

图 4 定义 2.2 中的组合安全集合系之间的关系

定理 2.2^[45, 49] 设 \mathcal{B} 是一个 (w, M, n) 集合系, 则有如下结论:

- (1) 如果 \mathcal{B} 是 t -TS, 则 \mathcal{B} 也是 t -IPPS; 反之, 则不一定成立.
- (2) 如果 \mathcal{B} 是 t -IPPS, 则 \mathcal{B} 也是 $(t, t; w)$ -UIBF; 反之, 则不一定成立.
- (3) 如果 \mathcal{B} 是 $(t, t; w)$ -UIBF, 则 \mathcal{B} 也是 t -CFF; 反之, 则不一定成立.

相应地, 则有如下推论:

推论 2.2 给定 t, w 和 n . 令 $M_{\text{TS}}(t; w, n)$ 、 $M_{\text{IPPS}}(t; w, n)$ 、 $M_{\text{UIBF}}(t, t; w; w, n)$ 和 $M_{\text{CFF}}(t; w, n)$ 分别为 t -TS、 t -IPPS、 $(t, t; w)$ -UIBF 和 t -CFF 中的最大区组个数, 则

$$M_{\text{TS}}(t; w, n) \leq M_{\text{IPPS}}(t; w, n) \leq M_{\text{UIBF}}(t, t; w; w, n) \leq M_{\text{CFF}}(t; w, n).$$

当将定义 2.2 中这几类组合安全集合系实际应用于广播加密中的会话密钥加密时, 一个 (w, M, n) 集合系 \mathcal{B} 中的 M 个区组代表被分配给 M 个授权用户的个人解密密钥. 假设有不超过 t 个授权用户 (记为集合 \mathcal{B}_0) 合谋攻击, 则他们可以根据所拥有的基本密钥制造出海盗版解密密钥 $T \subseteq \bigcup_{B \in \mathcal{B}_0} B$. 基于 (n, w) 阈值秘密共享方案的原理, 盗版 T 可以用来解密会话密钥, 从而合谋攻击者能够非法获取原始数据, 因而损害数据发行商的版权利益. 若盗版 T 被抓获, 我们希望能够利用集合系 \mathcal{B} 和 T 来追踪 \mathcal{B}_0 中的部分或者全部成员. 显然, 盗版 T 只包含了合谋者集合 \mathcal{B}_0 的部分信息, 如果要求追踪到所有的合谋攻击者, 那么 \mathcal{B} 需要满足苛刻的条件, 相应的集合系包含的区组个数 (对应于可容纳的授权用户数) 会非常少, 因而不实用. 这里将考虑从盗版 T 追踪到真实合谋攻击者集合 \mathcal{B}_0 中的至少一个合谋攻击者的情形. 根据定义 2.2, 有如下结论:

- 如果 \mathcal{B} 是 t -TS, 发行商通过比较盗版解密密钥 T 与每一个授权用户 $B \in \mathcal{B}$ 的个人密钥中相同的基本密钥的数目, 即 $|T \cap B|$ 的大小, 可以确定与盗版 T 的相同基本密钥数目最多的授权用户一定是真实的合谋攻击者之一, 其相应的追踪 / 译码算法的时间复杂度为 $O(nM)$.

- 如果 \mathcal{B} 是 t -IPPS, 通过检查 \mathcal{B} 的所有大小不超过 t 的子集可以找出盗版 T 的所有可能的限定父代集 \mathcal{B}' , 即 $\mathcal{P}_t(T)$. 由 (2.6) 可确定 $\mathcal{P}_t(T)$ 中所有 \mathcal{B}' 之间必有非空交集且这个交集也是 \mathcal{B}_0 的子集, 因此能够追踪到至少一个参与制造盗版 T 的合谋攻击者, 此追踪算法的时间复杂度是 $O(nM^t)$.

- 如果 \mathcal{B} 是 $(s, t; d)$ -UIBF, 由 (2.8) 可知, 对于任何一个人数不超过 s 且未参与合谋攻击的授权用户集 \mathcal{B}_1 , \mathcal{B}_1 中的任何用户都不能被其他由至多 t 个授权用户构成的合谋攻击团体 \mathcal{B}_0 所构陷, 因为这两个不相交的授权用户集不可能产生相同的盗版 T . 当 $s = 1$ 时, $(s, t; d)$ -UIBF 即为 t -CFF. 另外, 如果一个人数不超过 s 的授权用户集 \mathcal{B}_1 能够产生盗版 T , 则说明其中一定有参与违法制造盗版的合谋攻击用户, 但是合谋攻击用户的身份仍无法具体确定.

- 如果 \mathcal{B} 是 t -CFF, 由 (2.9) 可知, 对于任意一位未参与合谋攻击的授权用户, 其个人解密密钥不可能被 \mathcal{B}_0 中的授权用户合谋制造出来, 这就能够保证其不会被合谋攻击用户所构陷, 但是发行商无法根据盗版 T 来追踪任何一位参与违法复制的合谋攻击用户.

上面几类集合系都是在 (n, w) 阈值秘密共享方案的基础上提出的, 除此之外, 也有基于对指定用户集合授权的密钥分配方案设计等的研究, 有兴趣的读者可参见文献 [94].

2.2 多媒体指纹识别中组合安全码

本小节主要介绍多媒体指纹识别中的几类组合安全码及其追踪/译码算法思想, 包括防诬陷码、强可分离码、可分离码和签名码这四种具有完全可追踪能力的组合安全码, 以及强多媒体父代可识别码、多媒体父代可识别码这两种具有不完全可追踪能力的组合安全码.

首先, 在定义 2.3 中给出多媒体指纹识别中的几类组合安全码的数学定义. 其中, 强可分离码和可分离码的定义可分别参见文献 [56, 定义 3.3] 和 [17, 定义 1.1], 签名码的定义可参见文献 [34], 强多媒体父代可识别码和多媒体父代可识别码的定义可分别参见文献 [57, 定义 3.5] 和 [15, 定义 4.1].

定义 2.3 (多媒体指纹识别中的组合安全码) 设 \mathcal{C} 是一个 (n, M, q) 码, 对于任意 $\mathcal{C}_0 \subseteq \mathcal{C}$, $|\mathcal{C}_0| \leq t$, 记 $\mathcal{S} = \text{desc}(\mathcal{C}_0)$, 则

(i) 称 \mathcal{C} 为 \bar{t} -强可分离码 (strongly separable code, SSC), 记作 \bar{t} -SSC (n, M, q) , 如果

$$\bigcap_{\mathcal{C}' \in \mathcal{P}(\mathcal{S})} \mathcal{C}' = \mathcal{C}_0; \tag{2.10}$$

(ii) 称 \mathcal{C} 为 \bar{t} -可分离码 (separable code, SC), 记作 \bar{t} -SC (n, M, q) , 如果对于 \mathcal{C} 中任意不同于 \mathcal{C}_0 的子集 \mathcal{C}_1 , $|\mathcal{C}_1| \leq t$, 都有

$$\text{desc}(\mathcal{C}_1) \neq \text{desc}(\mathcal{C}_0); \tag{2.11}$$

(iii) 称 \mathcal{C} 为 \bar{t} -签名码 (signature code), 如果 $q = 2$, 而且对于 \mathcal{C} 中任意不同于 \mathcal{C}_0 的子集 \mathcal{C}_1 , $|\mathcal{C}_1| \leq t$, 以及任意满足 $\sum_{\mathbf{c}_j \in \mathcal{C}_0} \lambda_j = \sum_{\mathbf{c}_k \in \mathcal{C}_1} \lambda'_k = 1$ 的正实数 λ_j 和 λ'_k , 都有

$$\sum_{\mathbf{c}_j \in \mathcal{C}_0} \lambda_j \mathbf{c}_j \neq \sum_{\mathbf{c}_k \in \mathcal{C}_1} \lambda'_k \mathbf{c}_k; \tag{2.12}$$

(iv) 称 \mathcal{C} 为多媒体指纹识别中的 t -强父代可识别码 (strongly multimedia identifiable-parent-property code, SMIPPC), 记作 t -SMIPPC (n, M, q) , 如果

$$\bigcap_{\mathcal{C}' \in \mathcal{P}(\mathcal{S})} \mathcal{C}' \neq \emptyset; \tag{2.13}$$

(v) 称 \mathcal{C} 为多媒体指纹识别中的 t -父代可识别码 (multimedia identifiable-parent-property code, MIPPC), 记作 t -MIPPC (n, M, q) , 如果

$$\bigcap_{\mathcal{C}' \in \mathcal{P}_t(\mathcal{S})} \mathcal{C}' \neq \emptyset. \tag{2.14}$$

例 2.3 根据定义 2.3, 不难验证:

(1) 例 2.1 中的码 \mathcal{C} 是一个 $\bar{2}$ -SC, 也是一个 2-MIPPC, 这是因为对于 \mathcal{C} 的任意一个大小不超过 2 的子集 \mathcal{C}_0 , 其限定父代码集满足 $\mathcal{P}_t(\text{desc}(\mathcal{C}_0)) = \{\mathcal{C}_0\}$, 因此满足 (2.11) 和 (2.14).

(2) 例 2.1 中的码 \mathcal{C} 是一个 $\bar{2}$ -SSC, 也是一个 2-SMIPPC, 这是因为对于 \mathcal{C} 的任意一个大小不超过 2 的子集 \mathcal{C}_0 , 其父代码集也满足 $\mathcal{P}(\text{desc}(\mathcal{C}_0)) = \{\mathcal{C}_0\}$, 因此满足 (2.10) 和 (2.13).

(3) 令

$$\mathcal{C} = \begin{matrix} & \mathbf{c}_1 & \mathbf{c}_2 & \mathbf{c}_3 & \mathbf{c}_4 \\ \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \end{matrix}$$

则 \mathcal{C} 是一个长度为 3 的 $\bar{3}$ -签名码. 例如, 对于 \mathcal{C} 的两个不同的子集 $\mathcal{C}_0 = \{c_1, c_2, c_3\}$ 和 $\mathcal{C}_1 = \{c_1, c_2\}$, 假设公式 (2.12) 不成立, 即存在正实数 λ_i ($1 \leq i \leq 5$) 使得 $\lambda_1 c_1 + \lambda_2 c_2 + \lambda_3 c_3 = \lambda_4 c_1 + \lambda_5 c_2$, 其中 $\lambda_1 + \lambda_2 + \lambda_3 = \lambda_4 + \lambda_5 = 1$. 观察 \mathcal{C}_0 和 \mathcal{C}_1 的前两行会发现 $\lambda_1 + \lambda_3 = \lambda_4$ 且 $\lambda_2 + \lambda_3 = \lambda_5$, 这时可得到 $\lambda_3 = 0$, 与 $\lambda_3 > 0$ 矛盾. 事实上, 除了定义 2.3, 也可以利用定理 3.5 从几何的角度来证明 \mathcal{C} 是 $\bar{3}$ -签名码 (感兴趣的读者可以尝试).

但是 \mathcal{C} 不是 $\bar{3}$ -SC, 也不是 $\bar{3}$ -SSC, 这是因为对于上述 \mathcal{C}_0 和 \mathcal{C}_1 , $\text{desc}(\mathcal{C}_0) = \{0, 1\} \times \{0, 1\} \times \{0\} = \text{desc}(\mathcal{C}_1)$, 即 $\mathcal{C}_1 \in \mathcal{P}_3(\text{desc}(\mathcal{C}_0)) \subseteq \mathcal{P}(\text{desc}(\mathcal{C}_0))$, 但是 $\mathcal{C}_0 \neq \mathcal{C}_1$ 且 $\mathcal{C}_0 \not\subseteq \mathcal{C}_1$.

根据定义 2.1 和 2.3, 可以证明 FPC、SSC、SC、签名码、SMIPPC 和 MIPPC 之间有如下关系.

定理 2.3 设 \mathcal{C} 是一个 (n, M, q) 码, 则有如下结论:

- (1) (参见文献 [56, 引理 3.5, 例 3.6]) 如果 \mathcal{C} 是 t -FPC, 则 \mathcal{C} 也是 \bar{t} -SSC; 反之, 则不一定成立.
- (2) (参见文献 [56, 引理 4.2, 定理 4.4]) 如果 \mathcal{C} 是 \bar{t} -SSC, 则 \mathcal{C} 也是 \bar{t} -SC; 当 $n = t = 2$ 时, 反之也成立.
- (3) 如果 \mathcal{C} 是 \bar{t} -SSC, 则 \mathcal{C} 也是 t -SMIPPC; 反之, 则不一定成立.
- (4) (参见文献 [37, 引理 3.3.1]) 当 $q = 2$ 时, 如果 \mathcal{C} 是 \bar{t} -SC, 则 \mathcal{C} 也是 \bar{t} -签名码; 当 $q = t = 2$ 时, 反之也成立.
- (5) (参见文献 [15]) 如果 \mathcal{C} 是 \bar{t} -SC, 则 \mathcal{C} 也是 t -MIPPC; 当 $t = 2$ 时, 反之也成立.
- (6) (参见文献 [57, 引理 3.6]) 如果 \mathcal{C} 是 t -SMIPPC, 则 \mathcal{C} 也是 t -MIPPC; 反之, 则不一定成立.

根据定理 2.3, 对于一般的 $t \geq 2$, 我们在图 3 中同样简要刻画了 FPC、SSC、SC、SMIPPC、签名码和 MIPPC 之间的关系. 因此, 关于它们的最大渐进码率, 则有如下结论:

推论 2.3 设 $R_{\text{SSC}}(q, \bar{t})$ 、 $R_{\text{SC}}(q, \bar{t})$ 、 $R_{\text{SMIPPC}}(q, t)$ 和 $R_{\text{MIPPC}}(q, t)$ 分别为 q 元 \bar{t} -SSC、 \bar{t} -SC、 t -SMIPPC 和 t -MIPPC 的最大渐进码率, 则

$$R_{\text{FPC}}(q, t) \leq R_{\text{SSC}}(q, \bar{t}) \leq R_{\text{SC}}(q, \bar{t}) \leq R_{\text{MIPPC}}(q, t),$$

且

$$R_{\text{SSC}}(q, \bar{t}) \leq R_{\text{SMIPPC}}(q, t) \leq R_{\text{MIPPC}}(q, t).$$

根据第 1.5 小节介绍, 将一个 $(n, M, 2)$ 组合安全码 \mathcal{C} 实际应用于 M 个授权用户构成的多媒体指纹识别系统时, \mathcal{C} 中的一个码字被用于构造一个授权用户的多媒体指纹. 同样地, 假设有不超过 t 个授权用户参与合谋攻击 (记为集合 \mathcal{C}_0). 当发行商发现流通的违法复制品时, 可从中提取出盗版指纹并计算得到 $\mathbf{r} \in \mathbb{R}^n$ (参见第 1.5 小节). 设 \mathbf{r} 离散化后的结果为 $\mathcal{S} \subseteq \{0, 1\}^n$, 则 $\mathcal{S} = \text{desc}(\mathcal{C}_0)$. 因为 $\mathcal{P}_t(\mathcal{S})$ 中每一个 \mathcal{C}' 是 \mathcal{S} 的一个限定父代码, $\mathcal{P}(\mathcal{S})$ 中每一个 \mathcal{C}' 是 \mathcal{S} 的一个父代码, 所以 $\mathcal{C}_0 \in \mathcal{P}_t(\mathcal{S})$ 且 $\mathcal{C}_0 \in \mathcal{P}(\mathcal{S})$. 发行商的目的是利用码 \mathcal{C} 和得到的 \mathbf{r} 或 \mathcal{S} 来追踪 \mathcal{C}_0 中的成员. 根据定义 2.1 和 2.3, 有如下结论:

- 如果 \mathcal{C} 是 t -FPC, 由 (2.4) 可知, 对于任意 $\mathbf{c}' \in \mathcal{C} \setminus \mathcal{C}_0$, 都有 $\mathbf{c}' \notin \text{desc}(\mathcal{C}_0) = \mathcal{S}$. 因此, 可以检查码 \mathcal{C} 中的每一个码字并找出所有满足 $\mathbf{c} \in \mathcal{S}$ 的码字 \mathbf{c} , 这些码字构成的集合即为 \mathcal{C}_0 , 该追踪算法的时间复杂度为 $O(nM)$.

- 如果 \mathcal{C} 是 \bar{t} -SSC, 则追踪 \mathcal{C}_0 的过程分两步完成 (参见文献 [56]):

首先, 由 (1.1) 可知, 对于任意 $\mathbf{c}' \in \mathcal{C}$, 如果 $\mathbf{c}' \notin \text{desc}(\mathcal{C}_0) = \mathcal{S}$, 则对于 \mathcal{S} 的任意一个父代码 \mathcal{C}' 都有 $\mathbf{c}' \notin \mathcal{C}'$. 因此, 与 t -FPC 的追踪算法思想一样, 我们可以检查码 \mathcal{C} 的每一个码字并找出所有满足 $\mathbf{c} \in \mathcal{S}$

的码字 c , 不同的是, 这里找出来的码字所构成的集合不一定是 C_0 , 而是一个包含 C_0 的更大的集合. 为方便起见, 记这个集合为 C_L . 根据 C_L 中码字的选择方式, 不难证明 $C_0 \subseteq C_L$ 且 $C_L = \bigcup_{C' \in \mathcal{P}(S)} C'$. 进一步地, $\text{desc}(C_L) = S$.

接着, 需要从 C_L 中找出 C_0 . 由 (2.10) 可以证明, 对于任意 $c \in \bigcap_{C' \in \mathcal{P}(S)} C' = C_0$, 都有 $\text{desc}(C_L \setminus \{c\}) \neq S$. 特别地, 对于任意 $c \in C_0$, 一定存在某个 $i \in [n]$ 使得 $c(i) = 1$, 而对于任意 $x \in C_L \setminus \{c\}$ 都有 $x(i) = 0$; 或者一定存在某个 $k \in [n]$ 使得 $c(k) = 0$, 而对于任意 $x \in C_L \setminus \{c\}$ 都有 $x(k) = 1$. 因此, 通过检查 C_L 中的每一个码字并从中找出所有满足 $\text{desc}(C_L \setminus \{c\}) \neq S$ 的码字 c , 可以证明这些最终找出来的码字所构成的集合即为 C_0 , 该追踪算法的时间复杂度也为 $O(nM)$.

- 如果 C 是 \bar{t} -SC, 由 (2.11) 可知, S 的限定父代码只可能是 C_0 本身. 因此, 可以检查 C 的所有大小不超过 t 的子码, 比较这些子码的后代码和 S , 即可找出 S 的唯一限定父代码 C_0 , 追踪算法的时间复杂度为 $O(nM^t)$.

- 如果 C 是 \bar{t} -签名码, 根据第 1.5 小节介绍, $r = \sum_{c_j \in C_0} \lambda_j c_j$, 其中 λ_j 为正实数, 由 C_0 中的合谋攻击者决定, 且满足 $\sum_{c_j \in C_0} \lambda_j = 1$. 由 (2.12) 知, C 中能产生 r 且大小不超过 t 的子集只有 C_0 , 因此, 可以检查 C 的所有大小不超过 t 的子码, 比较这些子码和 r , 即可找出 C_0 , 追踪算法的时间复杂度为 $O(nM^t)$.

- 如果 C 是 t -SMIPPC, 类似于 \bar{t} -SSC 的追踪算法思想, 首先, 检查码 C 的每一个码字并找出所有满足 $c \in S$ 的码字 c , 得到一个码字集合 C_L 且 $\text{desc}(C_L) = S$; 其次, 根据 (2.13), 检查 C_L 中的每一个码字并从中找出所有满足 $\text{desc}(C_L \setminus \{c\}) \neq S$ 的码字 c , 可以证明这些码字对应的授权用户一定是参与违法复制的合谋攻击用户. 不同于 \bar{t} -SSC 的是, 这里从 C_L 中找出的码字集合不一定是 C_0 , 而是 C_0 的一个非空子集, 该追踪算法的时间复杂度为 $O(nM)$.

- 如果 C 是 t -MIPPC, 类似于 \bar{t} -SC 的追踪思想, 通过检查 C 的所有大小不超过 t 的子集, 找出 S 的所有限定父代码, 由 (2.14) 可以确定这些限定父代码有非空交集且这个交集是 C_0 的子集, 因此能够追踪到至少一个合谋攻击者, 该追踪算法的时间复杂度为 $O(nM^t)$.

表 1 总结上述二元 FPC、SSC、SC、签名码、SMIPPC 和 MIPPC 这几类组合安全码在多媒体指纹识别中的追踪能力和相应的 (一般情形下的) 追踪/译码算法时间复杂度.

2.3 公开问题

到目前为止, 对上述已经提出的组合安全码和集合系的研究都在努力地回答着问题 1.1, 即在寻找着其相应的应用场景中可允许的最大用户容量 (如最大渐近码率等) 与追踪/译码算法的时间复杂度之间的平衡. 通常情形下, 如果一类组合安全码的最大渐近码率相对较大, 则其对应的追踪算法的时间复杂度也会相对较大. 鉴于此, 我们提出如下问题:

公开问题 2.1 在广播加密应用场景中, 是否存在新的组合安全码或集合系可以更好地诠释最大渐近码率 (即可容纳的最大用户数目) 与追踪/译码算法之间的平衡关系? 如果有, 该如何定义它, 以及它与已经存在的组合安全码之间又有着怎样的联系?

表 1 多媒体指纹识别中的组合安全码

$(n, M, 2)$ 组合安全码	t -FPC	\bar{t} -SSC	\bar{t} -SC	\bar{t} -签名码	t -SMIPPC	t -MIPPC
能够追踪到的合谋攻击用户人数	全部	全部	全部	全部	≥ 1	≥ 1
译码算法的时间复杂度	$O(nM)$	$O(nM)$	$O(nM^t)$	$O(nM^t)$	$O(nM)$	$O(nM^t)$

公开问题 2.2 在多媒体指纹识别的应用中, 是否存在其他具有完全可追踪性能的组合安全码使其既具有签名码码率大的优点又具有防诬陷码和强可分离码译码快的优点?

公开问题 2.3 是否存在更多的特殊组合安全码或集合系, 使得基于这些码或集合系的追踪/译码算法的时间复杂度为 $O(\text{Poly}(\log M))$, 其中 M 为授权用户人数?

公开问题 2.4 通常在考虑组合安全码或集合系时, 人们假设合谋攻击团体的大小不会超过某个常数阈值 t . 但在现实中, 合谋攻击团体的大小有可能随着码长 n 或基本密钥个数 w 变化, 如 $\log n$ 或 $\log w$ 的多项式函数. 这种情形下, 一种可能的解决思路是利用概率方法来分析因果关系 (参见文献 [69]) 从而追踪合谋攻击者. 沿着这条思路进行研究, 将会得到什么样的结果?

3 组合安全码最大渐近码率的上界

本节主要介绍组合安全码最大渐近码率 (或最大码字数) 的上界的研究进展, 其中第 3.1 和 3.2 小节将分别介绍极值集合论和极值图论中的方法在研究和改进组合安全码最大渐近码率的上界方面所发挥的重要作用, 并列与组合安全码最大渐近码率上界相关的部分公开问题. 第 3.3 小节将从几何角度阐述在有噪情形的多媒体指纹识别中不存在具有完全可追踪能力的组合安全码的理由.

3.1 极值集合论方法

以 Sperner 定理和 Erdős-Ko-Rado 定理等为代表的极值集合论方法从 20 世纪发展至今, 在现代组合学和信息科学等领域有着丰富的应用 (参见文献 [59]). 本小节主要介绍用极值集合论的方法来研究广播加密中的组合安全集合系的大小的界定问题. 特别地, 本文以定义 2.2 中的 TS 和 CFF 为对象来展开讨论.

Erdős 等 [35, 36] 提出了特有子集 (own-subset) 的概念. 假设 \mathcal{B} 是一个 (w, M, n) 集合系, 对于任意 $B \in \mathcal{B}$ 和任意 $B_0 \subseteq B$, 称 B_0 是 B 的一个 $|B_0|$ -特有子集, 如果对于任意 $B' \in \mathcal{B} \setminus \{B\}$, 都有 $B_0 \not\subseteq B'$, 也就是说, B_0 只能被 B 所覆盖, 而不能被其他的任何区组所覆盖. Erdős 等 [36] 对 CFF 中任何一个区组所包含的特有子集个数进行了以下刻画:

引理 3.1 [36] 设 \mathcal{B} 是一个 t -CFF(w, M, n), 且 B 是 \mathcal{B} 中的任意一个区组, 则 B 所包含的 $\lfloor w/t \rfloor$ -特有子集的个数至少是 $\binom{w-1}{\lfloor w/t \rfloor - 1}$.

利用引理 3.1, Erdős 等 [36] 给出了 CFF 的大小的一个上界. 后来, Gu 和 Miao [48] 利用算两次 (double counting) 的方法对其进行了改进, 得到如下结果.

定理 3.1 [48] 对于任意 $n \geq w \geq 2$ 和 $t \geq 2$, 则有

$$M_{\text{CFF}}(t; w, n) \leq \frac{\binom{n}{\lfloor w/t \rfloor} - \binom{w-1}{\lfloor w/t \rfloor}}{\binom{w-1}{\lfloor w/t \rfloor - 1}}.$$

在关于 TS 的研究中, Stinson 和 Wei [85] 证明了, 一个 t -TS(w, M, n) 一定是一个 t -CFF(w, M, n), 并利用此关系和 CFF 的大小的已知上界得到了 TS 的大小的一般上界. 后来, Collins [22] 利用对 TS 中特有子集的研究对其进行了改进, 然而得到的上界也并非紧的. 最近, Gu 和 Miao [48] 发现了一个 TS 与 CFF 之间的非常有趣的关系, 并结合特有子集的结构, 得到了 TS 的大小的新的上界, 并证明了新的上界在很多情形下是可以达到的 (即最优的).

引理 3.2 [48] 一个 t -TS(w, M, n) 一定是一个 t^2 -CFF(w, M, n).

事实上, 引理 3.2 首次证明了组合安全集合系之间安全参数从 t 到 t^2 的非线性关系. 引理 3.2 的严格证明可参见文献 [48]. 其证明思路是利用反证法, 即假设一个集合系是 t -TS(w, M, n) 但不是 t^2 -CFF(w, M, n); 然后利用鸽巢原理 (pigeonhole principle) 等组合计数的技巧, 推导出与定义和假设的矛盾. 结合引理 3.1 和 3.2, 则有如下结论:

引理 3.3 [48] 设 \mathcal{B} 是一个 t -TS(w, M, n), 且 B 是 \mathcal{B} 中的任意一个区组. 则 B 所包含的 $\lfloor w/t^2 \rfloor$ -特有子集的个数至少是 $\binom{w-1}{\lfloor w/t^2 \rfloor - 1}$.

基于引理 3.3, 利用算两次的方法可以得到如下结果:

定理 3.2 [48] 对于任意 $n \geq w \geq 2$ 和 $t \geq 2$, 有

$$M_{\text{TS}}(t; w, n) \leq \frac{\binom{n}{\lfloor w/t^2 \rfloor} - \binom{w-1}{\lfloor w/t^2 \rfloor}}{\binom{w-1}{\lfloor w/t^2 \rfloor - 1}}.$$

定理 3.2 中的新的上界在某些情形下是紧的. 事实上, 依据太阳花 (sunflower) 和组合设计 (参见第 4.2 小节) 等构造得到的 TS 可以达到定理 3.2 中的上界, 即存在最优的 TS 的无穷类 (参见文献 [48]).

值得一提的是, 上述利用特有子集的证明思路也可以用于对 IPPS 和 UIBF 的区组个数上界的研究 (详见文献 [22, 46, 49]), 类似的研究思路在 FPC 的最大码率的研究中也有应用 (参见文献 [79]). 表 2 给出 TS、IPPS、UIBF 和 CFF 的最大区组个数的一般情形的渐近上下界, 其中 TS 和 CFF 的一般情形的渐近下界是通过构造性的方法得到的, 而 IPPS 和 UIBF 的一般情形的渐近下界是通过概率方法得到的.

针对表 2, 自然地提出下面问题:

公开问题 3.1 当 $(\lfloor (t/2 + 1)^2 \rfloor - 1) \nmid w$ 时, 能否确定 t -IPPS 的最优阶?

公开问题 3.2 当 $(s + t - 1) \nmid d$ 时, 能否确定 $(s, t; d)$ -UIBF 的最优阶?

关于公开问题 3.1, 在表 2 中的一般渐近上下界的基础上, 也有对于小参数情形 $t = 2$ 时的改进, 如文献 [46, 93] 中利用图论的方法对 2-IPPS 渐近上界的改进和文献 [50] 中利用加法数论 (additive number theory) 的方法对 2-IPPS 渐近下界的改进等.

类比于引理 3.2 中 TS 和 CFF 的关系, 我们也想知道定义 2.1 中的 TAC 和 FPC 是否也有类似的安全参数从 t 到 t^2 的非线性关系, 即有如下公开问题:

公开问题 3.3 一个 t -TAC(n, M, q) 一定是一个 t^2 -FPC(n, M, q) 吗?

公开问题 3.4 对于一个 t -TAC(n, M, q), 是否有 $M = O(q^{\lfloor n/t^2 \rfloor})$ 的渐近上界成立?

对于公开问题 3.4, 对于特殊情形 $t = 2, 3$, 文献 [9, 77] 给出了肯定的回答; 但是对于 $t \geq 4$ 的一般情形, 仍未解决.

3.2 极值图论方法

本小节以签名码为例阐述极值图论中的方法在研究组合安全码最大码字个数的上界方面的应用. 首先解释如何将一个 \bar{t} -签名码转化成一个围长至少为 $2t + 2$ 的二部图, 再利用具有该性质的二部图

表 2 安全集合系 TS、IPPS、UIBF 和 CFF 的渐近上下界 [45]

(w, M, n) 组合安全集合系	t -TS	t -IPPS	$(s, t; d)$ -UIBF	t -CFF
渐近上界 ($n \rightarrow \infty$)	$O(n^{\lfloor w/t^2 \rfloor})$	$O(n^{\lfloor \frac{w}{\lfloor (t/2+1)^2 \rfloor - 1} \rfloor})$	$O(n^{\lfloor \frac{d}{s+t-1} \rfloor})$	$O(n^{\lfloor w/t \rfloor})$
渐近下界 ($n \rightarrow \infty$)	$\Omega(n^{\lfloor w/t^2 \rfloor})$	$\Omega(n^{\lfloor \frac{w}{\lfloor (t/2+1)^2 \rfloor - 1} \rfloor})$	$\Omega(n^{\lfloor \frac{d}{s+t-1} \rfloor})$	$\Omega(n^{\lfloor w/t \rfloor})$

的结果得到 \bar{t} -签名码最大码字个数的上界.

二部图 (bipartite graph) 通常定义为三元组 $G = (V_1, V_2, E)$, 其中, V_1 和 V_2 是点集, E 是边集, E 中任意一条边分别关联 V_1 中的一个点和 V_2 中的一个点. $G = (V_1, V_2, E)$ 称为完全二部图, 如果 V_1 中任意一个点和 V_2 中所有的点都相邻, 若 $|V_1| = s, |V_2| = m$, 则完全二部图也记作 $K_{s,m}$. 对于由二部图 G 中 k 个不同的点构成的一个序列 $v_1 v_2 \cdots v_k v_1$, 若 $\{v_i, v_{i+1}\} \in E$ 且 $\{v_k, v_1\} \in E$, 则称该序列为 G 中一个长度为 k 的圈 (cycle), 记为 C_k . 如果 G 不含长度为 k 的圈, 则称 G 是不含 C_k (C_k -free) 的. 二部图 G 中最短圈的长度称为 G 的围长 (girth).

设 \mathcal{C} 是一个 $(n, M, 2)$ 码. 将标记码长坐标的集合 $[n]$ 划分成两个不交的子集 I_1 和 I_2 , 其中, $|I_1| = n_1, |I_2| = n_2$ 且 $n_1 + n_2 = n$. 对于 \mathcal{C} 中的任意一个码字 c , 记 $c|_{I_1}$ 为删除 c 中坐标落于 I_2 的元素后得到的新码字, $c|_{I_2}$ 为删除 c 中坐标落于 I_1 的元素后得到的新码字. 令 $V_1 = \{0, 1\}^{n_1}$ 和 $V_2 = \{0, 1\}^{n_2}$ 为二部图 $G = (V_1, V_2, E)$ 的点集. 显然, $c|_{I_1} \in V_1$ 且 $c|_{I_2} \in V_2$. 根据码 \mathcal{C} 以及坐标的划分 I_1 和 I_2 , 构造二部图 G 的边集如下: 对于 G 中任意两个点 $u \in V_1$ 和 $v \in V_2$, 若 \mathcal{C} 中存在一个码字 c 使得 $c|_{I_1} = u$ 且 $c|_{I_2} = v$, 则连接点 u 和 v , 即 $\{u, v\} \in E$. 根据构造不难发现, $|E| = M$. Fan 等^[39] 观察得到了以下的结论:

引理 3.4^[39] 设 \mathcal{C} 是一个长度为 n 的 \bar{t} -签名码, 则对于 $[n]$ 的任意一个划分 $[n] = I_1 \cup I_2, I_1 \cap I_2 = \emptyset$, 基于码 \mathcal{C} 和该划分得到的二部图 $G = (V_1, V_2, E)$ 的围长至少是 $2t + 2$.

根据引理 3.4, 研究 \bar{t} -签名码的码字个数的上界可以转化为研究不含 C_{2k} ($k \leq t$) 的二部图的边数的上界, 而二部图中的 C_4 也是 $K_{2,2}$. 以下定理给出了不含偶圈或 $K_{2,2}$ 的二部图的已知结果:

定理 3.3^[67, 72] 设 t, p, s 和 m 均为正整数且 $t \geq 2, p \geq s - 1, G = (V_1, V_2, E)$ 是一个二部图, 其中 $|V_1| = a, |V_2| = b$.

(1) 如果 G 不含 C_{2t} , 则

$$|E| \leq \begin{cases} (2t - 3)(b^{\frac{1}{2}} a^{\frac{t+2}{2t}} + a + b), & t \equiv 0 \pmod{2}, \\ (2t - 3)((ab)^{\frac{t+1}{2t}} + a + b), & t \equiv 1 \pmod{2}. \end{cases}$$

(2) 如果 G 不含 $K_{s,m}$, 则

$$|E| \leq \frac{m-1}{\binom{p}{s-1}} \binom{a}{s} + b \frac{(p+1)(s-1)}{s}.$$

设 $M_{\text{SigC}}(n, \bar{t})$ 为 n 长 \bar{t} -签名码的最大码字个数, 由定理 3.3 和引理 3.4 可得如下定理:

定理 3.4^[39] 设 $n, t \geq 2$ 均为正整数. 则

(1) 当 $t \geq 3$ 时,

$$M_{\text{SigC}}(n, \bar{t}) \leq \begin{cases} (2t - 3)(2^{\frac{(t+2)n}{2t+2} + 1} + 2^{\frac{tn}{2t+2}}), & t \equiv 0 \pmod{2}, \\ (2t - 3)(2^{\frac{(t+1)n}{2t}} + 2^{\lceil \frac{n}{2} \rceil} + 2^{\lfloor \frac{n}{2} \rfloor}), & t \equiv 1 \pmod{2}; \end{cases} \quad (3.1)$$

(2) 当 $t = 2$ 时,

$$M_{\text{SigC}}(n, \bar{2}) \leq 2^{\lfloor 2n/3 \rfloor} + 2^{\lceil n/3 \rceil} \frac{2^{\lceil n/3 \rceil} - 1}{2}. \quad (3.2)$$

定理 3.4 得到的关于签名码的一般上界在某些特定参数下是最优的 (参见文献 [39]). 特别注意的是, 关于不含 C_{2t} 的二部图的边数上界的结果在部分参数下有新的改进, 例如, 当 $t \geq 3$ 为奇数时,

定理 3.3 中的上界可以被改进为 $|E| \leq (ab)^{\frac{t+1}{2t}} + \max\{a, b\}$ (参见文献 [61]). 这也意味着 (3.1) 中关于 $M_{\text{SigC}}(n, t)$ 上界的结果也可以在系数方面得到改进. 这里不再一一列举, 感兴趣的读者可参见文献 [39, 61, 63] 等. 事实上, 上述将码转化成不含偶圈的二部图来界定最大码字个数的方法也被应用在其他组合安全码的研究中, 例如, D'yachkov 等 [28] 用此方法研究了可分离码, Cheng 等 [15] 利用该方法得到了多媒体父代可识别码最大码字个数的一般上界.

基于本小节的讨论, 可以考虑如下问题:

公开问题 3.5 能否找到更多关于签名码的直接构造, 使得构造出来的签名码的码字个数达到定理 3.4 中的上界? 如果不能, 如何通过其他方法, 如极值集合论和编码理论等, 来改进目前已知的上界?

3.3 基于几何方法的码的不存在性

本小节从几何角度阐述在有噪情形的多媒体指纹识别中不存在具有完全可追踪能力的组合安全码. 首先简单介绍相关的几何学术语, 并从几何学的角度给出多媒体指纹识别中的组合安全码的一个新的刻画, 再解释有噪情形下组合安全码不存在的原因.

为方便起见, 称由 $\{0, 1\}^n$ 中所有的点组成的几何结构为 n 维超方体 (hypercube). 设 $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^n$ 是 n 维 Euclid 空间中的点, $\delta > 0$ 是一个实数, 称集合 $\{\mathbf{z}' \in \mathbb{R}^n : \|\mathbf{z}' - \mathbf{z}\| < \delta\}$ 中所有的点构成的几何结构为以点 \mathbf{z} 为中心、 δ 为半径的开球 (不包含球面). 设 $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_s\} \subseteq \{0, 1\}^n$ 为 n 维超方体中的一个点集, 由 \mathcal{C} 中的点构成的称为开凸包 (open convex polytope) 的几何结构, 可表示为

$$\text{Cov}(\mathcal{C}) := \left\{ \sum_{i=1}^s \lambda_i \mathbf{c}_i \in \mathbb{R}^n : \sum_{i=1}^s \lambda_i = 1, \lambda_i > 0, \forall 1 \leq i \leq s \right\}. \quad (3.3)$$

因此, 可以从几何的角度刻画多媒体指纹识别中的签名码如下:

定理 3.5 [39] $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\} \subseteq \{0, 1\}^n$ 是一个 t -签名码, 当且仅当 \mathcal{C} 中任意两个不同的子集 \mathcal{C}_0 和 \mathcal{C}_1 , $2 \leq |\mathcal{C}_0|, |\mathcal{C}_1| \leq t$, 都满足 $\text{Cov}(\mathcal{C}_0) \cap \text{Cov}(\mathcal{C}_1) = \emptyset$.

第 3.2 小节介绍了如何利用极值图论的方法去研究签名码最大码字个数的上界, Egorova 等 [34] 利用 Goppa 码 [73] 的校验矩阵给出了签名码的一个直接构造, 得到了签名码的最大码字个数的一个下界, 但是目前已知的上下界之间还有一定的差距 (参见文献 [39]). 根据定理 3.5, 可以考虑以下问题:

公开问题 3.6 能否从几何学角度找到相关研究方法改进签名码最大码字个数的上下界? 特别地, 在 n 维超方体中, 最多能找到多少个顶点, 使得其中任意不超过 t 个顶点所构成的开凸包两两不相交?

对于有噪情形的多媒体指纹识别, 基于有损的海盗版计算得到的结果为 $\tilde{\mathbf{r}} = \sum_{j \in S} \lambda_j \mathbf{c}_j + \mathbf{e}$, 这里不妨设 $\|\mathbf{e}\| < \delta$, 则 $\tilde{\mathbf{r}}$ 可看成是以点 $\sum_{j \in S} \lambda_j \mathbf{c}_j$ 为中心、 δ 为半径的开球中的一个点. 为了能够追踪到合谋攻击用户集 S , 具有完全可追踪能力的组合安全码 $\mathcal{C} \subseteq \{0, 1\}^n$ 必须要满足这样的性质:

(P) 对于 \mathcal{C} 的任意两个不同且大小不超过 t 的子集 \mathcal{C}_1 和 \mathcal{C}_2 , 以及任意满足 $\sum_{\mathbf{c}_j \in \mathcal{C}_1} \lambda_j = \sum_{\mathbf{c}_k \in \mathcal{C}_2} \lambda'_k = 1$ 的正实数 λ_j, λ'_k 和任意 $\mathbf{e}, \mathbf{e}' \in \mathbb{R}^n, \|\mathbf{e}\|, \|\mathbf{e}'\| < \delta$, 都有

$$\sum_{\mathbf{c}_j \in \mathcal{C}_1} \lambda_j \mathbf{c}_j + \mathbf{e} \neq \sum_{\mathbf{c}_k \in \mathcal{C}_2} \lambda'_k \mathbf{c}_k + \mathbf{e}'.$$

为了从几何的角度描述上述性质, Fan 等 [39] 定义了 \mathbb{R}^n 中任意两个点集 V_1 和 V_2 之间的距离为 $d(V_1, V_2) := \inf\{\|\mathbf{z} - \mathbf{z}'\| : \mathbf{z} \in V_1, \mathbf{z}' \in V_2\}$, 基于此定义, 他们发现:

定理 3.6 ^[39] 性质 (P) 等价于, \mathcal{C} 中任意两个不同且大小不超过 t 的子集 \mathcal{C}_1 和 \mathcal{C}_2 都满足

$$d(\text{Cov}(\mathcal{C}_1), \text{Cov}(\mathcal{C}_2)) \geq 2\delta. \quad (3.4)$$

然而, 当 $\mathcal{C}_1 \cap \mathcal{C}_2 \neq \emptyset$ 时, (3.4) 不可能成立. 因此, 在有噪情形的多媒体指纹识别中, 不存在具有完全可追踪能力的组合安全码. 基于定理 3.6 以及版权保护中对具有不完全可追踪能力的组合安全码的研究思路, 可以考虑如下问题:

公开问题 3.7 在有噪情形下, 多媒体指纹识别中是否存在不完全可追踪的组合安全码? 如果存在, 该如何定义这样的码? 其最大码字个数或最大渐近码率是多少? 怎么具体构造这样的码?

4 组合安全码最大渐近码率的下界

本节分别从概率方法、组合构造等角度介绍研究组合安全码最大渐近码率下界的方法, 每小节将以部分安全码的研究为例进行介绍, 其中, 第 4.3 小节的内容与朱烈老师在哈希 (hash) 族方面的工作有着紧密的联系.

4.1 基于概率方法的存在性

概率方法是一种可以用来证明某一类数学对象的存在性的强有力的方法 (参见文献 [3]), 其在数学、信息科学和人工智能等领域都有着极为普遍的应用. 在 20 世纪, 国际著名的数学家 Erdős 将概率方法引入到组合学的研究中, 从而大大促进了现代组合学的进步和发展, 感兴趣的读者可参见文献 [59]. 毫不例外, 在对于版权保护中的组合安全码的已有研究中, 多类安全码和集合系的存在性结果也都是利用概率方法给出的, 如父代可识别集合系 ^[46]、并交限定族 ^[49]、可分离码 ^[8, 43, 47]、多媒体父代可识别码 ^[46] 和强多媒体父代可识别码 ^[57] 等. 本小节以父代可识别集合系为例来介绍概率方法在版权保护中的应用.

Gu 等 ^[46] 将不同应用下的具有父代可识别性质的码和集合系进行了统一, 给出了如下定义. 其中, 他们采用了术语 f -信道 (f -channel) 来表示一种因果关系 (cause-effect relationship), 该信道有至多 t 个输入 (cause) 和一些可能的输出 (effect) 集合.

定义 4.1 设 \mathcal{C} 是 Q^n 或 2^Q 的一个子集. 称 \mathcal{C} 是一个在 f -信道下的父代可识别方案 (parent-identifying scheme), 如果对于任意 $\mathcal{C}_0 \subseteq \mathcal{C}$ (满足 $|\mathcal{C}_0| \leq t$) 和任意 $d \in f(\mathcal{C}_0)$, 都有

$$\bigcap_{\mathcal{P} \subseteq \mathcal{C}: |\mathcal{P}| \leq t, d \in f(\mathcal{P})} \mathcal{P} \neq \emptyset.$$

上述定义中的 f -信道可以根据不同的应用场景进行具体设定, 从而一些已知的码和集合系都可看成是特殊的父代可识别方案, 如父代可识别码 (定义 2.1)、父代可识别集合系 (定义 2.2)、多媒体父代可识别码 (定义 2.3), 以及群试理论中的单人可追踪族 (single-user-tracing family) ^[1, 24] 等.

Gu 等 ^[46] 试图将父代可识别方案与极值组合学中的 Turán 类型问题建立联系. Turán 类型问题主要研究组合对象在不允许某些子结构出现时的最大量化问题. Barg 等 ^[4] 提出了最小构型的概念.

定义 4.2 给定集合 \mathcal{C} 以及它的子集族 $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_m\}$, 其中 $\mathcal{F}_i \subseteq \mathcal{C}$ 且 $|\mathcal{F}_i| \leq t$. 称 \mathcal{F} 为一个构型 (configuration), 如果 \mathcal{F} 中元素的交集为空集, 即 $\bigcap_{1 \leq i \leq m} \mathcal{F}_i = \emptyset$. 进一步地, 称 \mathcal{F} 为一个最小构型 (minimal configuration), 如果去掉 \mathcal{F} 中任何一个元素后, 其交集不再为空集.

针对父代可识别方案的要求, Gu 等 ^[46] 定义最小禁止构型如下:

定义 4.3 设 \mathcal{C} 是在 f -信道下的一个父代可识别方案. \mathcal{C} 中的一个最小禁止构型 (minimal forbidden configuration) 是指一个最小构型 $\mathcal{F} = \{\mathcal{F}_1, \dots, \mathcal{F}_m\}$, 满足性质 $f(\mathcal{F}_1) \cap \dots \cap f(\mathcal{F}_m) \neq \emptyset$.

下面的对父代可识别分配方案的等价刻画已在文献 [46] 中被证明:

引理 4.1 [46] 一个集合 \mathcal{C} 是在 f -信道下的父代可识别方案当且仅当 \mathcal{C} 中不存在最小禁止构型 \mathcal{F} 使得 $|\bigcup_{1 \leq i \leq m} \mathcal{F}_i| \leq \lfloor (t/2 + 1)^2 \rfloor$.

根据引理 4.1 和概率方法中的删去法 (deletion method, 或 probabilistic expurgation method), Gu 等 [46] 证明了 IPPS 的存在性下界:

定理 4.1 [46] 对于给定的 w 和 $t \geq 2$, 当 n 足够大时, 存在 t -IPPS(w, M, n), 其中

$$M = \Omega(n^{\frac{w}{\lfloor (t/2+1)^2 - 1 \rfloor}}).$$

定理 4.1 的证明思路是, 首先随机选取一些区组构成一个随机集合, 然后对该随机集合中的最小禁止构型的个数进行估计, 再利用删去法在每一个最小禁止构型中去掉一个区组来摧毁所有的最小禁止构型, 最后得到 IPPS 的存在性结果 (详细证明可参见文献 [46]). 上述一般的 IPPS 的下界, 在某些特定情形下也可以被改进, 如文献 [50, 78] 等.

相对于定理 4.1 的存在性结果, 自然地有如下问题:

公开问题 4.1 是否存在具体的 IPPS 构造方法, 使得构造得到的 IPPS 也可以达到定理 4.1 中的下界?

公开问题 4.2 定理 4.1 中的下界是否是阶最优的? 如果不是, 怎样改进?

4.2 基于纠错码和组合设计的构造

在 20 世纪 40 年代, 美国数学家 Hamming 提出了可在信息传输中进行错误检测与纠正的重要数学工具—纠错码. 组合设计理论 (combinatorial design theory) 起源于 18 世纪, 在试验设计、统计科学和异常检测等领域有着丰富的应用. 在当今现代信息科学的发展中, 纠错码理论和组合设计理论仍起着重要的不可替代的作用. 本小节将介绍如何利用纠错码和组合设计方法来构造组合安全码和集合系. 设 \mathcal{C} 是一个 (n, M, q) 码, 其最小 Hamming 距离定义为 $d(\mathcal{C}) := \min_{\mathbf{c}, \mathbf{c}' \in \mathcal{C}: \mathbf{c} \neq \mathbf{c}'} d(\mathbf{c}, \mathbf{c}')$.

Staddon 等 [82] 给出了如下的关系:

定理 4.2 [82] 设 \mathcal{C} 是一个 (n, M, q) 码, 其最小 Hamming 距离为 $d(\mathcal{C})$, 则

(1) \mathcal{C} 是一个 t -FPC, 如果 $d(\mathcal{C}) > (1 - \frac{1}{t})n$;

(2) \mathcal{C} 是一个 t -TAC 和一个 t -IPPC, 如果 $d(\mathcal{C}) > (1 - \frac{1}{t^2})n$.

在定义 2.1 中, IPPC 是比 TAC 的要求更弱的一类码, 而定理 4.2 中的由纠错码来构造 IPPC 和 TAC 的充分条件却是一样的. 自然地, 有如下问题:

公开问题 4.3 定理 4.2 中的充分条件是否也是必要条件? 特别地, 是否存在一个 (n, M, q) 码 \mathcal{C} , 其最小 Hamming 距离 $d(\mathcal{C}) \leq (1 - \frac{1}{t^2})n$, 且 \mathcal{C} 是一个 t -IPPC?

考虑码 \mathcal{C} 是特殊的最大距离可分 (maximum distance separable, MDS) 码或者是 Reed-Solomon 码 [73] 时, Fernandez 等 [40] 以及 Jin 和 Blaum [58] 证明了当 $d(\mathcal{C}) \leq (1 - \frac{1}{t^2})n$ 时, 在某些特殊参数下, 码 \mathcal{C} 不再满足 IPPC 或者 TAC 的条件, 即不再是 IPPC 或 TAC. 但是对于很多其他参数的情形仍然是未解决的公开问题.

类似地, 组合设计可以用来构造安全码和集合系. 首先, 给出 τ -设计 (τ -design) 的定义 (参见文献 [20]). 一个 τ - (n, w, λ) 设计是一个 (w, M, n) 集合系 \mathcal{B} 且 $[n]$ 的任意 τ -子集都恰好被包含在 λ 个

区组中. 易知, 一个 τ -(n, w, λ) 设计包含 $M = \lambda \binom{n}{\tau} / \binom{w}{\tau}$ 个区组 (参见文献 [20]). Stinson 和 Wei^[85] 证明了如下关系:

定理 4.3^[85] 如果存在一个 τ -($n, w, 1$) 设计, 则

(1) 一定存在一个 t -FPC($n, M, 2$), 其中 $M = \binom{n}{\tau} / \binom{w}{\tau}$ 且 $t = \lfloor (w-1)/(\tau-1) \rfloor$;

(2) 一定存在一个 t -TS(w, M, n), 其中 $M = \binom{n}{\tau} / \binom{w}{\tau}$ 且 $t = \lfloor \sqrt{(w-1)/(\tau-1)} \rfloor$.

利用定理 4.3 中的构造方法, 可以得到能够达到第 3.1 小节中介绍的 TS 码字数上界的最优 TS 的无穷类, 详细可参见文献 [45, 48].

在 SC 的构造方面, Cheng 等^[16] 利用 \bar{t} -SC(n_1, M, q) 和 \bar{t} -SC(n_2, q, q') 通过级联构造 (concatenated construction)^[17] 得到 \bar{t} -SC($n_1 n_2, M, q'$), 其中 $2 \leq q' < q$. 这个构造方法促使他们去研究具有小参数的 \bar{t} -SC(n, M, q), 如 $(t, n) = (2, 2), (2, 3)$. 他们证明了定义在 $[q]$ 上的 $\bar{2}$ -SC($2, M, q$), 记为 \mathcal{C} , 等价于一种特殊的组合设计, 即 $(q, W, 1)$ 广义填充 (generalized packing), 记为 $([q], \mathcal{B})$, 其中 $\mathcal{B} = \{B_1, B_2, \dots, B_q\}$, $B_i = \{x \in [q] : (i, x) \in \mathcal{C}\}$, $W = \{|B_i| : B_i \in \mathcal{B}\}$, 且 $[q]$ 的任意 2-子集都至多被包含在 \mathcal{B} 的一个区组中. 根据构造不难看出, $M = \sum_{i=1}^q |B_i|$. 通过简单的最优化方法, 他们证明了 $M \leq \lfloor q(1 + \sqrt{4q-3})/2 \rfloor$, 然后通过射影平面 (projective plane) 得到了达到此上界的最优 $\bar{2}$ -SC($2, M, q$) 的无穷类. 而后, Cheng 等^[14] 利用图论和射影平面得到了更多的结果, 并揭示了此问题与极值图论中 Zarankiewicz 数 (参见文献 [10]) 的联系. 通过差矩阵 (difference matrix), Cheng 等^[16] 也构造出了所有的最优 $\bar{2}$ -SC($3, M, q$).

除了上述利用 τ -设计、广义填充和差矩阵来构造组合安全码和集合系以外, 在相关文献中, 也有利用正交阵列 (orthogonal array) 等组合结构来构造组合安全码的研究, 如文献 [13, 85, 86] 等.

公开问题 4.4 对 $(t, n) = (2, 2)$, 能否具体构造出所有的最优 SC? 对 $(t, n) = (2, \geq 4)$ 和 $(t, n) = (> 2, \geq t)$, 怎样刻画和构造最优 SC?

4.3 与哈希族的联系

本小节主要介绍版权保护中的组合安全码与可分和完美哈希族之间的关系, 并举例说明如何利用哈希族证明组合安全码的存在性.

一个 (M, q) -哈希函数 (hash function) 是指将集合 A 映射到集合 B 的函数, 其中, $|A| = M$, $|B| = q$ 且 $M \geq q$. 一个 (M, q) -哈希族 (hash family) \mathcal{H} 是指由有限个 (M, q) -哈希函数构成的集合, 当 $|\mathcal{H}| = n$ 时, \mathcal{H} 称为 (n, M, q) -哈希族. 如果对于 A 的任意 t 个两两不相交的子集 X_1, X_2, \dots, X_t , $|X_j| \leq w_j$, $1 \leq j \leq t$, 都存在 $h \in \mathcal{H}$ 使得 $h(X_1), h(X_2), \dots, h(X_t)$ 也是两两不相交的, 则称 (n, M, q) -哈希族 \mathcal{H} 是可分的 (separating), 记为 $(n, M, q, \{w_1, w_2, \dots, w_t\})$ -可分哈希族. 当 $w_1 = w_2 = \dots = w_t = 1$ 时, $(n, M, q, \{w_1, w_2, \dots, w_t\})$ -可分哈希族也称为完美的 (perfect), 记为 (n, M, q, t) -完美哈希族. 可分哈希族和完美哈希族是组合学中的研究热点, 这两个结构与极值图论中的稀疏超图有着紧密的联系 (参见文献 [44, 76]), 在密码学和信息科学中也有诸多应用 (参见文献 [68, 84, 92]).

为了跟组合安全码建立联系, 先将一个 (n, M, q) -哈希族 \mathcal{H} 转化成一个 q -元的 $n \times M$ 矩阵, 集合 A 中的元素对应为矩阵的列, \mathcal{H} 中的哈希函数对应为矩阵的行, 当将该矩阵中每个列向量看成一个 n 长的码字时, \mathcal{H} 对应到一个 (n, M, q) 码. 基于这样的转化, 可分哈希族和完美哈希族与组合安全码之间的关系得到大量的关注和研究. 例如, 根据定义 2.1, 不难看出 $(n, M, q, \{t, 1\})$ -可分哈希族和 t -FPC(n, M, q) 实际上是等价的 (参见文献 [84]), 而以下结论则揭示了完美哈希族和父代可识别码之间的关系:

定理 4.4^[82] 若 \mathcal{H} 是一个 $(n, M, q, \lfloor (t+2)^2/4 \rfloor)$ -完美哈希族, 则由 \mathcal{H} 可得到一个 t -IPPC(n, M, q).

由定理 4.4 可知, 如果能构造一个 $(n, M, q, \lfloor (t+2)^2/4 \rfloor)$ -完美哈希族, 则能得到一个 t -IPPC (n, M, q) , 这为研究 IPPC 的构造或最大渐近码率的下界提供了一种思路. 值得一提的是, 在可分和完美哈希族的构造等方面, 朱烈老师及其合作者有着突出的贡献 (参见文献 [86, 87]).

其实定理 4.4 中的完美哈希族可以弱化为偏哈希族 (partially hash family). 一个 (n, M, q) -哈希族 \mathcal{H} 称为偏哈希族, 记为 (n, M, q, k, u) -偏哈希族, 如果对于 A 的任意一个 u -子集 D 及 D 的任意一个 k -子集 X , 都存在 $h \in \mathcal{H}$, 使得对于任意 $x \in X$, 均满足 $h(\{x\}) \cap h(D \setminus \{x\}) = \emptyset$.

定理 4.5 [7] 若 \mathcal{H} 是一个 $(n, M, q, t, \lfloor (t+2)^2/4 \rfloor)$ -偏哈希族, 则由 \mathcal{H} 可得到一个 t -IPPC (n, M, q) .

与完美哈希族的情形不同, 关于偏哈希族的研究很少. 为得到更多的 IPPC, 一方面, 可以沿着朱烈老师及其合作者的思路构造更多的完美哈希族; 另一方面, 也应该考虑如下问题:

公开问题 4.5 当给定 M, q, k 和 u 时, (n, M, q, k, u) -偏哈希族中 n 的最小值是多少? 怎样构造 n 达到最小值时的偏哈希族?

5 与组合安全码相关的其他问题

本节将介绍版权保护中的组合安全码与其他研究领域之间的关联, 其中, 第 5.1 小节阐述组合安全码与群试理论中两个著名的测试矩阵之间的关系, 第 5.2 小节介绍组合安全码在多用户通信系统中的应用.

5.1 组合安全码与群试理论

群试理论是 Dorfman [26] 于 1943 年提出的一种关于大规模样本检测方法的理论, 其基本思想是, 对待测样本进行分组, 对同一组内的混合样本进行合并检测, 从而达到减少检测次数和节约检测成本的目的. 之后群试理论被广泛研究并应用于工程和生物等领域, 如血液检测、化学品泄漏检测和网络故障诊断等. 值得一提的是, 这种合并检测的方法在新冠病毒检测中发挥着重要作用 (参见文献 [66]), 在实际检测中已经被中国、以色列和新加坡等国家采用.

对样本进行分组检测时, 如果一个组的测试结果是阳性 (positive), 则认为该组中存在至少一个阳性样本; 若结果是阴性 (negative), 则说明该组中所有样本都是阴性. 关于样本的分组策略, 一般来说有适应性 (adaptive) 和非适应性 (non-adaptive) 这两种类别. 在适应性分组测试中, 样本的分组方法由上一个组的测试结果决定, 而非适应性分组测试则是为所有的样本设计好分组之后再并行检测. 如果阳性样本服从某个概率分布, 则称关于该模型下的分组测试为概率群试 (probabilistic group testing); 如果阳性样本的个数不超过某个特定阈值 t , 则称该模型下的分组测试为组合群试 (combinatorial group testing). 根据文献 [31, 32], 一个包含 M 个样本和 n 个分组的非适应性组合群试方案可以通过一个 $n \times M$ 的二元矩阵表示, 矩阵的行和列分别代表分组和样本. 为方便起见, 用 $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\} \subseteq \{0, 1\}^n$ 来表示测试矩阵中的列向量, 若第 j 个样本在第 i 组中, 则 $c_j(i) = 1$; 否则 $c_j(i) = 0$. 所有 n 个组的测试结果可以表示成一个二元向量 $\mathbf{r} = (r(1), r(2), \dots, r(n))^T \in \{0, 1\}^n$, 其中, $r(i) = 1$ 表示第 i 个组的测试结果为阳性, $r(i) = 0$ 表示第 i 个组的测试结果为阴性. 设 $S \subseteq [M]$ 表示阳性样本的集合, 则

$$\mathbf{r} = \bigvee_{j \in S} \mathbf{c}_j, \quad (5.1)$$

其中“ \vee ”为向量之间的布尔和 (Boolean sum). 非适应性组合群试理论的研究内容即为设计和研究分组测试方案 $\mathcal{C} = \{c_1, c_2, \dots, c_M\} \subseteq \{0, 1\}^n$, 使得根据测试结果 r 和分组方案 \mathcal{C} , 当 $|S| \leq t$ 时, 可以快速找出 S 中的阳性样本. 虽然 (5.1) 中的结果和版权保护中的攻击模型两者之间的数学表达有所不同, 但是关于非适应性组合群试中测试矩阵的研究和版权保护中组合安全码的研究非常相似.

可分离矩阵 (separable matrix) 和析取矩阵 (disjunct matrix) 是非适应性组合群试中两类经典的测试矩阵, 关于它们的研究起源于 Kautz 和 Singleton^[60] 为文件修复而提出的叠加码. 1982 年, Erdős 等^[35] 在集合系中给出了与可分离矩阵和析取矩阵分别对应的两种组合结构: 无并族 (union-free family) 和无覆盖族, 并对这两类集合族的概念作了一般化的推广 (参见文献 [36]), 更多关于这两类经典的组合结构的研究可参见文献 [23, 30, 41, 42, 55, 74]. Cheng 和 Miao^[17] 研究多媒体指纹识别中的组合安全码时发现了多媒体指纹识别和非适应性组合群试之间的联系, 提出了可分离码的概念, 并给出了防诬陷码在多媒体指纹识别中的快速追踪/译码算法. 同时, 他们还证明了可分离矩阵和析取矩阵可作为组合安全码应用于版权保护.

当阳性样本的个数不超过 t 时, \bar{t} -可分离矩阵和 t -析取矩阵均可以识别出所有的阳性样本, 但是二者在码率和译码算法方面各有优缺点. t -析取矩阵的译码算法比 \bar{t} -可分离矩阵的更高效, 但是 \bar{t} -可分离矩阵的码率更大, 这意味着当待测样本数量相同时, 相比于 t -析取矩阵, \bar{t} -可分离矩阵所需要的检测次数更少 (参见文献 [31, 32]). 最近, Fan 等^[38] 借鉴多媒体指纹识别中的强可分离码的追踪算法思想, 在非适应性组合群试理论中提出了强可分离矩阵 (strongly separable matrix) 的概念, 证明了强可分离矩阵的性质介于可分离矩阵和析取矩阵之间, 但是能够提供的译码算法与析取矩阵的一样高效. 由于强可分离矩阵是新提出的概念, 目前关于它的研究结果还很少, 感兴趣的读者可参见文献 [38, 71]. 在具有不完全可识别能力的测试矩阵或集合族方面, 单人可追踪族^[1, 24] 和多人可追踪族 (multiple-user-tracing family)^[2, 62] 是常被研究的两类组合结构, 可用于识别至少一个或多个阳性样本, 其概念和译码算法与版权保护中的父代可识别码类似.

5.2 组合安全码与多用户通信

多用户通信, 也称为“多对一”通信, 常见于卫星通信、移动通信和无线通信等通信系统. 关于多用户通信系统的数学理论研究起源于 Shannon 在 1961 年发表的文献 [80], 该文献引入了接入两个用户的通信系统的数学模型, 而后一般化的多用户接入的模型 (见图 5) 得到了广泛研究 (参见文献 [21, 29, 53]).

在多用户通信系统中, 多址接入信道 (multiple access channel, MAC) 作为共享的信息传输媒介连接多个用户与一个信息接收方进行通信. 不同于传统的“一对一”和“一对多” (广播) 通信, “多对一”通信的特点是多个用户发送的信息在信道中传输时会相互干扰. 信息接收方收到信息之后, 需要重建出用户发送的信息. 当每个用户所要发送的信息只有一个时, 信息的重建问题也就是用户身份的识别

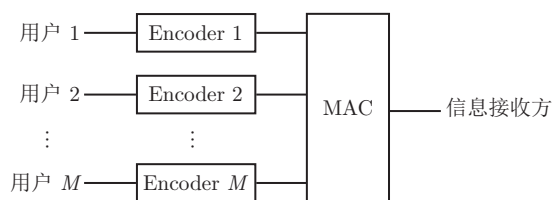


图 5 多用户通信的信道模型

问题. 例如, 在一个接入了 M 个移动用户但只能提供有限个信道 (个数远小于 M) 的通信系统中, 用户 j 活跃时可以使用编码器 (encoder) j 进行编码, 然后通过无线信道向中央控制单元发送一个信息 $\mathbf{c}_j \in \{0, 1\}^n$, 中央控制单元再从无线信道的输出信息 \mathbf{r} 中检测出所有活跃用户并为其分配专用信道. 信道不同, 输出信息 \mathbf{r} 也不同. 为了识别出所有活跃用户, 需要根据无线信道的类型来设计发送的信息, 即 $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\} \subseteq \{0, 1\}^n$. 该问题的研究和版权保护中组合安全码的研究也极为相似, 这里不同的多址接入信道可以对应为版权保护中不同的攻击模型 (参见文献 [33]). 特别地, Egorova 等^[34]指出多媒体指纹识别中的线性攻击和一类二元加权信道的数学模型是等价的, 并提出了既可以应用于二元加权信道又可以应用于多媒体指纹识别的签名码 (参见定义 2.3). 基于他们的工作, Fan^[37]对二元加权信道和多媒体指纹识别展开了系统研究, 证明了防诬陷码、可分离码和强可分离码均是特殊的签名码 (参见定理 2.3), 因此都可以作为签名码应用于二元加权信道. 更多关于多用户通信和版权保护的共同研究可参见文献 [28].

6 结论

本文介绍了版权保护中对抗合谋攻击的数学理论及其最新研究进展, 也介绍了版权保护与群试理论和多用户通信领域相关组合问题的联系. 其实, 版权保护中用到的数学思想也能应用到压缩感知 (compressed sensing)、差分隐私 (differential privacy) 和机器学习 (machine learning) 等领域的研究中, 限于篇幅, 这里无法展开介绍, 有兴趣的读者可参见文献 [12, 25, 27, 64, 70, 91] 等. 我们期待有更多的研究者进入这一引人入胜的研究领域, 解决遗留的公开问题, 建立新的理论体系, 推动此领域及相关领域更快地发展.

致谢 作者非常感谢审稿人提出的宝贵建议. 本文所讨论的可追踪组合安全码、防诬陷组合安全码以及相关的集合系, 与朱烈老师曾经做出突出贡献的无覆盖族、哈希族等组合问题^[86, 87]有着不可分割的联系. 谨以此文献给朱烈老师八十年华诞, 恭祝朱烈老师生日快乐, 健康长寿!

参考文献

- 1 Alon N, Asodi V. Tracing a single user. *European J Combin*, 2006, 27: 1227–1234
- 2 Alon N, Asodi V. Tracing many users with almost no rate penalty. *IEEE Trans Inform Theory*, 2007, 53: 437–439
- 3 Alon N, Spencer J H. *The Probabilistic Method*, 3rd ed. New York: John Wiley & Sons, 2008
- 4 Barg A, Cohen G, Encheva S, et al. A hypergraph approach to the identifying parent property: The case of multiple parents. *SIAM J Discrete Math*, 2001, 14: 423–431
- 5 Barg A, Kabatiansky G A. A class of I.P.P. codes with efficient identification. *J Complexity*, 2004, 20: 137–147
- 6 Blackburn S R. Frameproof codes. *SIAM J Discrete Math*, 2003, 16: 499–510
- 7 Blackburn S R. Combinatorial schemes for protecting digital content. In: *London Mathematical Society Lecture Note Series. Surveys in Combinatorics*, vol. 307. Cambridge: Cambridge University Press, 2003, 43–78
- 8 Blackburn S R. Probabilistic existence results for separable codes. *IEEE Trans Inform Theory*, 2015, 61: 5822–5827
- 9 Blackburn S R, Etzion T, Ng S L. Traceability codes. *J Combin Theory Ser A*, 2010, 117: 1049–1057
- 10 Bollobás B. *Extremal Graph Theory*. New York: Dover, 2004
- 11 Boneh D, Shaw J. Collusion-secure fingerprinting for digital data. *IEEE Trans Inform Theory*, 1998, 44: 1897–1905
- 12 Candes E J, Romberg J, Tao T. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans Inform Theory*, 2006, 52: 489–509
- 13 Chee Y M, Zhang X. Improved constructions of frameproof codes. *IEEE Trans Inform Theory*, 2012, 58: 5449–5453
- 14 Cheng M, Fu H L, Jiang J, et al. New bounds on $\bar{2}$ -separable codes of length 2. *Des Codes Cryptogr*, 2015, 74: 31–40
- 15 Cheng M, Fu H L, Jiang J, et al. Codes with the identifiable parent property for multimedia fingerprinting. *Des Codes Cryptogr*, 2017, 83: 71–82
- 16 Cheng M, Ji L, Miao Y. Separable codes. *IEEE Trans Inform Theory*, 2012, 58: 1791–1803

- 17 Cheng M, Miao Y. On anti-collusion codes and detection algorithms for multimedia fingerprinting. *IEEE Trans Inform Theory*, 2011, 57: 4843–4851
- 18 Chor B, Fiat A, Naor M. Tracing traitors. In: *Advances in Cryptology—CRYPTO’94*. Lecture Notes in Computer Science, vol. 839. Berlin: Springer-Verlag, 1994, 480–491
- 19 Chor B, Fiat A, Naor M, et al. Tracing traitors. *IEEE Trans Inform Theory*, 2000, 46: 893–910
- 20 Colbourn C J, Dinitz J. *Handbook of Combinatorial Designs*, 2nd ed. Boca Raton: CRC Press, 2007
- 21 Colbourn C J, Heller J, Viterbi A. A new coding technique for asynchronous multiple access communication. *IEEE Trans Comm Technol*, 1971, 19: 849–855
- 22 Collins M J. Upper bounds for parent-identifying set systems. *Des Codes Cryptogr*, 2009, 51: 167–173
- 23 Coppersmith D, Shearer J B. New bounds for union-free families of sets. *Electron J Combin*, 1998, 5: R39
- 24 Csürös M, Ruzinkó M. Single-user tracing and disjointly superimposed codes. *IEEE Trans Inform Theory*, 2005, 51: 1606–1611
- 25 Donoho D L. Compressed sensing. *IEEE Trans Inform Theory*, 2006, 52: 1289–1306
- 26 Dorfman R. The detection of defective members of large populations. *Ann of Math Stud*, 1943, 14: 436–440
- 27 Dwork C, Roth A. The algorithmic foundations of differential privacy. *Found Trends Theoret Comput Sci*, 2014, 9: 211–407
- 28 D’yachkov A G, Polyanskii N, Shchukin V Y, et al. Separable codes for the symmetric multiple-access channel. *IEEE Trans Inform Theory*, 2019, 65: 3738–3750
- 29 D’yachkov A G, Rykov V V. A coding model for a multiple-access adder channel. *Probl Peredachi Inf*, 1981, 17: 26–38
- 30 D’yachkov A G, Rykov V V. Bounds on the length of disjunctive codes. *Probl Peredachi Inf*, 1982, 18: 7–13
- 31 Du D, Hwang F K. *Combinatorial Group Testing and its Applications*, 2nd ed. Singapore: World Scientific, 2000
- 32 Du D, Hwang F K. *Pooling Designs and Nonadaptive Group Testing*. Singapore: World Scientific, 2006
- 33 Egorova E E, Fernandez M, Kabatiansky G A, et al. Signature codes for A -channel and collusion-secure multimedia fingerprinting codes. In: *Proceedings of IEEE International Symposium on Information Theory (ISIT)*. Barcelona: IEEE, 2016, 3043–3047
- 34 Egorova E E, Fernandez M, Kabatiansky G A, et al. Signature codes for weighted noisy adder channel, multimedia fingerprinting and compressed sensing. *Des Codes Cryptogr*, 2019, 87: 455–462
- 35 Erdős P, Frankl P, Füredi Z. Families of finite sets in which no set is covered by the union of two others. *J Combin Theory Ser A*, 1982, 33: 158–166
- 36 Erdős P, Frankl P, Füredi Z. Families of finite sets in which no set is covered by the union of r others. *Israel J Math*, 1985, 51: 79–89
- 37 Fan J. Signature coding schemes for multiple access channels with applications to multimedia fingerprinting and group testing. PhD Thesis. Tsukuba: University of Tsukuba, 2020
- 38 Fan J, Fu H L, Gu Y, et al. Strongly separable matrices for nonadaptive combinatorial group testing. *Discrete Appl Math*, 2021, 291: 180–187
- 39 Fan J, Gu Y, Hachimori M, et al. Signature codes for weighted binary adder channel and multimedia fingerprinting. *IEEE Trans Inform Theory*, 2021, 67: 200–216
- 40 Fernandez M, Cotrina J, Soriano M, et al. A note about the identifier parent property in Reed-Solomon codes. *Comput Security*, 2010, 29: 628–635
- 41 Frankl P, Füredi Z. Union-free hypergraphs and probability theory. *European J Combin*, 1984, 5: 127–131
- 42 Füredi Z. On r -cover-free families. *J Combin Theory Ser A*, 1996, 73: 172–173
- 43 Gao F, Ge G. New bounds on separable codes for multimedia fingerprinting. *IEEE Trans Inform Theory*, 2014, 60: 5257–5262
- 44 Ge G, Shangguan C, Wang X. Some intriguing upper bounds for separating hash families. *Sci China Math*, 2019, 62: 269–282
- 45 Gu Y. New progress on combinatorial schemes for broadcast encryption and codes for multimedia fingerprinting. PhD Thesis. Tsukuba: University of Tsukuba, 2018
- 46 Gu Y, Cheng M, Kabatiansky G A, et al. Probabilistic existence results for parent-identifying schemes. *IEEE Trans Inform Theory*, 2019, 65: 6160–6170
- 47 Gu Y, Fan J, Miao Y. Improved bounds for separable codes and B_2 codes. *IEEE Commun Lett*, 2020, 24: 15–19
- 48 Gu Y, Miao Y. Bounds on traceability schemes. *IEEE Trans Inform Theory*, 2018, 64: 3450–3460
- 49 Gu Y, Miao Y. Union-intersection-bounded families and their applications. *Discrete Appl Math*, 2019, 266: 346–354
- 50 Gu Y, Satake S. On 2-parent-identifying set systems of block size 4. *Des Codes Cryptogr*, 2020, 88: 2067–2076
- 51 Gu Y, Vorobyev I, Miao Y. Secure codes with list decoding. In: *Proceedings of IEEE International Symposium on*

- Information Theory (ISIT). Espoo: IEEE, 2022, 2350–2355
- 52 Guruswami V, Sudan M. Improved decoding of Reed-Solomon and algebraic-geometry codes. *IEEE Trans Inform Theory*, 1999, 45: 1757–1767
- 53 Györi S. Signature coding over multiple access OR channel. In: *Proceedings of IEEE Workshop on Information Theory*. Paris: IEEE, 2003, 115–118
- 54 Hollmann H D L, van Lint J H, Linnartz J P, et al. On codes with the identifiable parent property. *J Combin Theory Ser A*, 1998, 82: 121–133
- 55 Hwang F K, Sós V T. Nonadaptive hypergeometric group testing. *Studia Sci Math Hungar*, 1987, 22: 257–263
- 56 Jiang J, Cheng M, Miao Y. Strongly separable codes. *Des Codes Cryptogr*, 2016, 79: 303–318
- 57 Jiang J, Gu Y, Cheng M. Multimedia IPP codes with efficient tracing. *Des Codes Cryptogr*, 2020, 88: 851–866
- 58 Jin H, Blaum M. Combinatorial properties for traceability codes using error correcting codes. *IEEE Trans Inform Theory*, 2007, 53: 804–808
- 59 Jukna S. *Extremal Combinatorics: With Applications in Computer Science*, 2nd ed. Heidelberg: Springer, 2011
- 60 Kautz W H, Singleton R C. Nonrandom binary superimposed codes. *IEEE Trans Inform Theory*, 1964, 10: 363–377
- 61 Keevash P, Sudakov B, Verstraëte J. On a conjecture of Erdős and Simonovits: Even cycles. *Combinatorica*, 2013, 33: 699–732
- 62 Laczay B, Ruzinkó M. Multiple user tracing codes. In: *Proceedings of IEEE International Symposium on Information Theory (ISIT)*. Seattle: IEEE, 2006, 1900–1904
- 63 Li B, Ning B. Exact bipartite Turán numbers of large even cycles. *J Graph Theory*, 2021, 97: 642–656
- 64 Li S, Gao F, Ge G, et al. Deterministic construction of compressed sensing matrices via algebraic curves. *IEEE Trans Inform Theory*, 2012, 58: 5035–5041
- 65 Liu K J R, Trappe W, Wang Z J, et al. *Multimedia Fingerprinting Forensics for Traitor Tracing*. New York: Hindawi, 2005
- 66 Mutesa L, Ndishimye P, Butera Y, et al. A pooled testing strategy for identifying SARS-CoV-2 at low prevalence. *Nature*, 2021, 589: 276–280
- 67 Naor A, Verstraëte J. A note on bipartite graphs without $2k$ -cycles. *Combin Probab Comput*, 2005, 14: 845–849
- 68 Newman I, Wigderson A. Lower bounds on formula size of Boolean functions using hypergraph entropy. *SIAM J Discrete Math*, 1995, 8: 536–542
- 69 Pearl J. *Causality: Models, Reasoning, and Inference*, 2nd ed. New York: Cambridge University Press, 2009
- 70 Polyanskii N. On learning sparse vectors from mixture of responses. *Adv Neural Inform Process Syst*, 2021, in press
- 71 Qian B, Wang X, Ge G. Improved lower bounds for strongly separable matrices and related combinatorial structures. *arXiv:2110.07381*, 2021
- 72 Roman S. A problem of Zarankiewicz. *J Combin Theory Ser A*, 1975, 18: 187–198
- 73 Roth R. *Introduction to Coding Theory*. Cambridge: Cambridge University Press, 2012
- 74 Ruzinkó M. On the upper bound of the size of the r -cover-free families. *J Combin Theory Ser A*, 1994, 66: 302–310
- 75 Shamir A. How to share a secret. *Commun ACM*, 1979, 22: 612–613
- 76 Shangguan C, Ge G. Separating hash families: A Johnson-type bound and new constructions. *SIAM J Discrete Math*, 2016, 30: 2243–2264
- 77 Shangguan C, Ma J, Ge G. New upper bounds for parent-identifying codes and traceability codes. *Des Codes Cryptogr*, 2018, 86: 1727–1737
- 78 Shangguan C, Tamo I. Sparse hypergraphs with applications to coding theory. *SIAM J Discrete Math*, 2020, 34: 1493–1504
- 79 Shangguan C, Wang X, Ge G, et al. New bounds for frameproof codes. *IEEE Trans Inform Theory*, 2017, 63: 7247–7252
- 80 Shannon C E. Two-way communication channels. In: *Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1. Berkeley: University of California Press, 1961, 611–644
- 81 Silverberg A, Staddon J, Walker J. Efficient traitor tracing algorithms using list decoding. In: *Advances in Cryptology—ASIACRYPT’01*. Lecture Notes in Computer Science. Berlin: Springer, 2001, 175–192
- 82 Staddon J N, Stinson D R, Wei R. Combinatorial properties of frameproof and traceability codes. *IEEE Trans Inform Theory*, 2001, 47: 1042–1049
- 83 Stinson D R, Paterson M. *Cryptography: Theory and Practice*. Boca Raton: CRC Press, 2018
- 84 Stinson D R, van Trung T, Wei R. Secure frameproof codes, key distribution patterns, group testing algorithms and related structures. *J Statist Plann Inference*, 2000, 86: 595–617
- 85 Stinson D R, Wei R. Combinatorial properties and constructions of traceability schemes and frameproof codes. *SIAM J Discrete Math*, 1998, 11: 41–53

- 86 Stinson D R, Wei R, Zhu L. New constructions for perfect hash families and related structures using combinatorial designs and codes. *J Combin Des*, 2000, 8: 189–200
- 87 Stinson D R, Wei R, Zhu L. Some new bounds for cover-free families. *J Combin Theory Ser A*, 2000, 90: 224–234
- 88 Trappe W, Wu M, Wang Z J, et al. Anti-collusion fingerprinting for multimedia. *IEEE Trans Signal Process*, 2003, 51: 1069–1087
- 89 van Trung T, Martirosyan S. On a class of traceability codes. *Des Codes Cryptogr*, 2004, 31: 125–132
- 90 van Trung T, Martirosyan S. New constructions for IPP codes. *Des Codes Cryptogr*, 2005, 35: 227–239
- 91 Vadhan S. The complexity of differential privacy. In: Lindell Y, ed. *Tutorials on the Foundations of Cryptography. Information Security and Cryptography*. Cham: Springer, 2017, 347–450
- 92 Walker II R A, Colbourn C J. Perfect hash families: Constructions and existence. *J Math Cryptol*, 2007, 1: 125–150
- 93 Wang X. Improved upper bounds for parent-identifying set systems and separable codes. *Des Codes Cryptogr*, 2021, 89: 91–104
- 94 Wei R. Traceability schemes, frameproof codes, key distribution patterns and related topics: A combinatorial approach. PhD Thesis. Nebraska: University of Nebraska, 1998

Combinatorial secure codes for copyright protection and related problems

Jinping Fan, Yujie Gu & Ying Miao

Abstract The rapid development of modern science and technology not only provides convenience for data communication, but also poses a tremendous threat to the copyright of digital content. This paper focuses on the mathematical theory of traitor-tracing for copyright protection and its latest progress. First, for applications in different scenarios, such as broadcast encryption and multimedia fingerprinting, we introduce unified models to characterize the (key/fingerprint) distribution schemes with traceability property and frameproof property, respectively. Next, we review several classes of combinatorial secure codes with the traceability property and frameproof property, and the combinatorial methods used to investigate the bounds on the maximum code size and the explicit constructions, as well as the latest results and the open problems. The relationships between the combinatorial problems in copyright protection and the related problems in group testing and multiple access communication are discussed as well.

Keywords copyright protection, combinatorial secure code, traitor-tracing scheme, frameproof allocation scheme, set system, group testing theory, multi-user communication

MSC(2020) 05B20, 05D05, 05D40, 68P30, 94A60, 94B25

doi: 10.1360/SSM-2022-0079