

УДК 621.391 : 519.72

© 2020 г. Е.Е. Егорова¹, М. Фернандес², Г.А. Кабатянский¹, И. Мяо³**СУЩЕСТВОВАНИЕ И КОНСТРУКЦИИ МУЛЬТИМЕДИЙНЫХ КОДОВ,
СПОСОБНЫХ НАХОДИТЬ ПОЛНУЮ КОАЛИЦИЮ ПРИ АТАКЕ
УСРЕДНЕНИЯ И ШУМЕ**

Как было недавно показано в [1], не существует мультимедийных кодов цифровых водяных знаков, способных полностью восстановить коалицию недобросовестных пользователей в условиях общей линейной атаки и целенаправленного шума. Мы покажем, что такие коды существуют, если сузить класс атак до атаки усреднения. Возникающая математическая задача близка к задаче построения сигнатурных кодов для двоичного суммирующего канала с шумом.

Ключевые слова: мультимедийный код цифровых отпечатков пальцев, канал множественного доступа, целенаправленный шум, сигнатурный код, атака на основе сговора, коды без перекрытия, дизъюнктивные коды.

DOI: 10.31857/S0555292320040087

§ 1. Введение

Математическая постановка задачи защиты цифрового контента от нелегального копирования и перераспределения возникла в конце прошлого века, см. [2–4]. Наибольшее внимание привлекла постановка задачи, известная как коды цифровых отпечатков пальцев и существующая в различных вариациях, см. [5–8]. Соответствующие модели являются дискретными, и первая непрерывная модель, мотивированная приложениями к защите мультимедийного контента (изображения, музыка), появилась в работах [9, 10] под названием мультимедийные коды отпечатков пальцев. То, что эти коды тесно связаны с сигнатурными кодами для соответствующих каналов множественного доступа, в частности, для А-канала [11], неявно появилось в работе [12], а позже – в явной форме в [13]. Затем эта связь была распространена на сигнатурные коды для взвешенного двоичного суммирующего канала [14]. Этот подход был дальше развит в [1], где было доказано, в частности, что мультимедийные коды отпечатков пальцев не способны полностью восстановить коалицию недобросовестных пользователей в условиях общей линейной атаки и целенаправленного шума. В данной статье мы покажем, что ситуация более оптимистична, если несколько ограничить атаки коалиций, а именно при атаке усреднения существуют коды, которые находят всех членов коалиции даже в присутствии целенаправленного шума. Отметим, что далее мы будем употреблять термин “цифровые водяные знаки” вместо “отпечатки пальцев”.

¹ Работа выполнена при частичной финансовой поддержке Российского фонда фундаментальных исследований (номера проектов 20-51-50007 и 20-07-00652).

² Работа выполнена при частичной финансовой поддержке гранта правительства Испании TCORISEBLOCK (номер гранта PID2019-110224RB-I00, проект MINECO/FEDER) и гранта правительства Каталонии 2017-SGR-782.

³ Работа выполнена при частичной финансовой поддержке Японского общества содействия развитию науки (JSPS) в рамках научного проекта JPJSBP120204802.

§ 2. Мультимедийные коды цифровых водяных знаков

Рассмотрим математическую модель защиты мультимедийного контента от нелегального перераспределения. Мультимедийный контент представляется как N -мерный вектор \mathbf{x} над полем \mathbb{R} вещественных чисел. Имеется дистрибьютор, который видоизменяет \mathbf{x} специальным образом для каждого пользователя так, чтобы если коалиция недобросовестных пользователей подделает \mathbf{x} , то он может найти всех членов коалиции (complete traceability; см. [1]). Для этого дистрибьютор выбирает m попарно ортогональных векторов $\mathbf{f}_1, \dots, \mathbf{f}_m$ в \mathbb{R}^N одинаковой длины (для простоты – длины 1), которые известны только ему – это важное предположение, которое будет использоваться в дальнейшем. Затем дистрибьютор формирует так называемые *цифровые водяные знаки* (ЦВЗ) как линейные комбинации этих векторов с двоичными коэффициентами (известен и другой вариант “модуляции” ЦВЗ, когда используются коэффициенты из $\{-1, +1\}$).

ЦВЗ \mathbf{w}_j , предназначенный j -му пользователю, имеет вид

$$\mathbf{w}_j = \sum_{i=1}^m h_{ij} \mathbf{f}_i, \quad (1)$$

где $h_{ij} \in \{0, 1\}$. Вложение ЦВЗ осуществляется аддитивно, т.е. дистрибьютор выдает j -му пользователю вектор

$$\mathbf{y}_j = \mathbf{x} + \mathbf{w}_j \quad (2)$$

как копию \mathbf{x} , где предполагается, что длина вектора \mathbf{x} много больше длины ЦВЗ \mathbf{w}_j (т.е. $\|\mathbf{x}\| \gg \|\mathbf{w}_j\|$, где $\|\cdot\|$ здесь и далее обозначает евклидову норму вектора), чтобы копия \mathbf{y}_j мало отличалась от оригинала \mathbf{x} .

Пусть имеется M пользователей и среди них коалиция $A \subset [M]$ недобросовестных пользователей. *Линейная атака* состоит в том, что коалиция A генерирует поддельную копию \mathbf{y} как линейную комбинацию имеющихся у нее копий \mathbf{y}_j с вещественными коэффициентами $\lambda_1, \dots, \lambda_M$, такими что $\sum_{j=1}^M \lambda_j = 1$, $\lambda_j > 0$ для всех $j \in A$ и $\lambda_j = 0$ для $j \notin A$, т.е.

$$\mathbf{y} = \sum_{j=1}^M \lambda_j \mathbf{y}_j = \sum_{a \in A} \lambda_a \mathbf{y}_a, \quad (3)$$

Так как $\sum_{j=1}^M \lambda_j = 1$, то $\mathbf{y} = \mathbf{x} + \sum_{j=1}^M \lambda_j \mathbf{w}_j$, а так как все $\lambda_j \geq 0$, то в силу неравенства треугольника (для нормы) имеем

$$\|\mathbf{y} - \mathbf{x}\| = \left\| \sum_{j=1}^M \lambda_j \mathbf{w}_j \right\| \leq \sum_{j=1}^M \lambda_j \|\mathbf{w}_j\| \leq \max_j \|\mathbf{w}_j\| \ll \|\mathbf{x}\|, \quad (4)$$

и следовательно, \mathbf{y} является достаточно хорошей копией оригинала \mathbf{x} .

Так как дистрибьютор знает значение \mathbf{x} , то чтобы определить, что \mathbf{y} – это нелегальная копия \mathbf{x} , и найти всех участников коалиции, создавших \mathbf{y} , он вычисляет скалярные произведения

$$s_k = (\mathbf{y} - \mathbf{x}, \mathbf{f}_k) = \left(\sum_{j=1}^M \lambda_j \sum_{i=1}^m h_{ij} \mathbf{f}_i, \mathbf{f}_k \right) = \sum_{j=1}^M \lambda_j h_{kj} \quad (5)$$

и формирует вектор-синдром $\mathbf{S} = \mathbf{S}(\Lambda) = (s_1, \dots, s_m)$, где $\Lambda := (\lambda_1, \dots, \lambda_M)$. Отметим, что для вектора Λ его носителем $\text{supp}(\Lambda) := \{j : \lambda_j \neq 0\}$ является коалиция A .

Введем векторы $\mathbf{h}_1, \dots, \mathbf{h}_M$, где $\mathbf{h}_j = (h_{1j}, \dots, h_{mj})$. Тогда (5) можно переписать в виде

$$\mathbf{S} = \sum_{j=1}^M \lambda_j \mathbf{h}_j = \sum_{a \in A} \lambda_a \mathbf{h}_a. \quad (6)$$

Это уравнение, в свою очередь, можно записать как матричное уравнение

$$\mathbf{S} = H\Lambda^T, \quad (7)$$

где H – матрица размера $m \times M$, составленная из векторов-столбцов $\mathbf{h}_1, \dots, \mathbf{h}_M$.

Так как векторы $\mathbf{f}_1, \dots, \mathbf{f}_m$ ортонормированные, а векторы ЦВЗ $\mathbf{w}_1, \dots, \mathbf{w}_M$ выражаются в базисе $\mathbf{f}_1, \dots, \mathbf{f}_m$ с двоичными коэффициентами, а именно

$$\mathbf{w}_j = \sum_{i=1}^m h_{ij} \mathbf{f}_i, \quad \text{где } h_{ij} \in \{0, 1\},$$

то множества $\mathcal{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$ и $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\}$ изометричны. При этом векторы $\mathbf{f}_1, \dots, \mathbf{f}_m$ известны дистрибьютору (и только ему), поэтому для него равносильно, работать ли с множеством ЦВЗ \mathcal{W} или с множеством \mathcal{H} соответствующих двоичных векторов. Поэтому далее мы будем оба множества называть мультимедийным кодом, а если по синдрому \mathbf{S} можно однозначно найти носитель $\text{supp}(\Lambda)$, т.е. коалицию A , то будем называть их *t-мультимедийным кодом со свойством полного поиска коалиций* (*t*-МППК-кодом).

Определение 1. Двоичный код $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\} \subset \{0, 1\}^m$ называется *t-МППК-кодом*, если для любых двух вещественных векторов Λ и Λ' , таких что

$$\sum_{j=1}^M \lambda_j = \sum_{j=1}^M \lambda'_j = 1 \quad \text{и} \quad |\text{supp}(\Lambda)|, |\text{supp}(\Lambda')| \leq t,$$

из $H\Lambda^T = H\Lambda'^T$ следует, что $\text{supp}(\Lambda) = \text{supp}(\Lambda')$.

Замечание 1. Заметим, что в данном определении мы не воспользовались ограничением, что все λ_j неотрицательны.

Замечание 2. Всюду далее мы предполагаем, что значение t фиксировано.

Среди всех линейных атак особо выделяется *атака усреднения*, для которой $\lambda_j = |A|^{-1}$ при $j \in A$ и $\lambda_j = 0$ в противном случае. Ранее, начиная с первых работ по этой тематике (см. [9, 10]), считалось, что “атака усреднения является наиболее справедливой для участников коалиции, чтобы избежать обнаружения” [12]. Поэтому все работы до [13] ограничивались рассмотрением только атаки усреднения. В этой статье мы покажем, что атака усреднения намного слабее общей линейной атаки, по крайней мере, в случае целенаправленного шума.

Напомним, что двоичный код $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\} \subset \{0, 1\}^m$ называется *t-сигнатурным кодом* для двоичного суммирующего канала (см. [15]), если для любых двух различных кодовых подмножеств A и B , где оба подмножества мощности не более t , справедливо неравенство

$$\sum_{j \in A} \mathbf{h}_j \neq \sum_{j \in B} \mathbf{h}_j. \quad (8)$$

Код, для которого неравенство (8) выполнено для любых двух различных кодовых подмножеств A и B одинаковой мощности не более t , будем называть $(=, t)$ -сигнатурным кодом для суммирующего канала. Двоичный код, обладающий свойством нахождения полной коалиции при атаке усреднения, является $(=, t)$ -сигнатурным кодом для суммирующего канала. Обратное неверно ни для $(=, t)$ -, ни для t -сигнатурного кода для двоичного суммирующего канала.

Обозначим через $M(m, t)$ максимально возможную мощность t -МППК-кода длины m , а через $R(m, t) := m^{-1} \log_2 M(m, t)$ – соответствующую кодовую скорость. Тогда из сделанного выше замечания об атаке усреднения и обычных соображений мощности вытекает следующий аналог границы Хэмминга:

$$C_{M(m,t)}^t \leq (t+1)^m. \quad (9)$$

Отметим, что из известной верхней границы для скорости сигнатурных кодов для суммирующего канала [16] следует асимптотически в два раза лучшая, чем (9), граница на скорость кода

$$\limsup_m R(m, t) \leq \frac{\log_2 t}{2t} (1 + o(1)). \quad (10)$$

С другой стороны, в [14] была доказана следующая нижняя граница:

$$M(m, t) \geq 2^{\lfloor m/t \rfloor}, \quad (11)$$

из которой следует, что $R(m, t) \geq t^{-1}(1 + o(1))$.

Повторим кратко основные аргументы из [14]. Прежде всего отметим, что если среди векторов $\mathbf{h}_1, \dots, \mathbf{h}_M$ любые $2t$ линейно независимы над \mathbb{R} , то определение 1 выполнено, и более того, дистрибьютор может найти не только те j , которым соответствуют $\lambda_j \neq 0$, но и соответствующие значения λ_j (и кроме того, ограничение $\sum_{j=1}^m \lambda_j = 1$ не требуется). Примером такого множества двоичных векторов являются столбцы проверочной матрицы двоичного кода, исправляющего t ошибок, так как любые $2t$ ее столбцов линейно независимы над полем из двух элементов, а следовательно, и над \mathbb{R} . Теперь, чтобы получить (11), остается взять в качестве $\mathbf{h}_1, \dots, \mathbf{h}_M$ столбцы проверочной $(m \times M)$ -матрицы неприводимого кода Гоппы (см. [17]) длины $M = 2^\ell$ и избыточности $m = \ell t$, исправляющего t ошибок.

§ 3. Мультимедийные коды цифровых водяных знаков в присутствии шума

Отметим, что ранее в литературе уже рассматривались модели мультимедийных кодов цифровых водяных знаков в присутствии шума: вероятностная модель шума (см. [18]) и модель целенаправленного шума, где нет ограничения на норму (длину) вектора ошибки, а есть только ограничение на вес Хэмминга ошибки [14].

В этой статье, следуя [1], мы рассматриваем модель, когда коалиция A не только создает ложную копию

$$\mathbf{y} = \mathbf{x} + \sum_{j \in A} \lambda_j \mathbf{w}_j \in \mathbb{R}^N$$

в соответствии с моделью линейной атаки, но еще и целенаправленно добавляет вектор шума $\mathbf{e} \in \mathbb{R}^N$, такой что $\|\mathbf{e}\| \leq \delta$, где $\|\cdot\|$ – евклидова норма на \mathbb{R}^N . В результате коалиция A формирует и перераспределяет копию

$$\hat{\mathbf{y}} = \mathbf{x} + \sum_{j \in A} \lambda_j \mathbf{w}_j + \mathbf{e}. \quad (12)$$

Множество ЦВЗ $\mathcal{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\} \subset \mathbb{R}^N$ будем называть (t, δ) -мультимедийным кодом ЦВЗ со свойством полного поиска коалиций, устойчивым к δ -шуму, если по любой ложной копии $\hat{\mathbf{y}} = \sum_{j \in A} \lambda_j \mathbf{y}_j + \mathbf{e}$ можно однозначно найти коалицию A .

Это условие равносильно тому, что для любых двух различных коалиций A и B , $|A|, |B| \leq t$, и любых двух вещественных векторов Λ и Λ' , таких что

$$\sum_{j \in A} \lambda_j = \sum_{j \in B} \lambda'_j = 1,$$

справедливо неравенство

$$\left\| \sum_{j \in A} \lambda_j \mathbf{w}_j - \sum_{j \in B} \lambda'_j \mathbf{w}_j \right\| > 2\delta. \quad (13)$$

Довольно очевидно, что условие (13) слишком сильное, и таких кодов не существует, что и было показано в [1]. Приведем дополнительные к [1] аргументы, почему таких кодов нет. Положим $A = \{1, 2, \dots, t\}$ и $B = \{1, 2, \dots, t-1, t+1\}$. Обозначим $w = \max_{j \in [t+1]} \|\mathbf{w}_j\|$ и выберем $\lambda_t = \lambda'_{t+1} = \lambda$ так, чтобы $0 < \lambda < \min\{\delta w^{-1}, 1\}$. Выберем положительные $\lambda_j = \lambda'_j$ для $j = 1, \dots, t-1$ такими, что $\lambda + \sum_{j=1}^{t-1} \lambda_j = 1$. Тогда

$$\sum_{j \in A} \lambda_j \mathbf{w}_j + \mathbf{e} = \sum_{j=1}^{t-1} \lambda_j \mathbf{w}_j = \sum_{j \in B} \lambda_j \mathbf{w}_j + \mathbf{e}',$$

где $\mathbf{e} = -\lambda \mathbf{w}_t$, $\mathbf{e}' = -\lambda \mathbf{w}_{t+1}$ и $\|\mathbf{e}\|, \|\mathbf{e}'\| \leq \delta$, и неравенство (13) не выполнено.

Этот анализ условия (13) показывает, что если какое-то λ_j отлично от нуля, но достаточно мало, то от него можно “избавиться”, введя взамен другое ненулевое λ_i , что и не позволяет найти коалицию целиком. Однако у атаки усреднения все λ_j равны и достаточно отделены от нуля, что позволяет надеяться, что если ограничить класс линейных атак только атакой усреднения, то существуют коды, находящие коалицию целиком; см. следующее

Определение 2. Множество $\mathcal{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$ будем называть (t, δ) -мультимедийным кодом со свойством полного поиска коалиций, устойчивым к атаке усреднения и δ -шуму, если для любых двух различных подмножеств кода $A, B \subset \mathcal{W}$, таких что $|A|, |B| \leq t$, справедливо неравенство

$$\left\| \frac{1}{|A|} \sum_{j \in A} \mathbf{w}_j - \frac{1}{|B|} \sum_{j \in B} \mathbf{w}_j \right\| > 2\delta. \quad (14)$$

Ниже нам будет удобнее рассматривать не код \mathcal{W} , а изометричный ему двоичный код \mathcal{H} длины m , для которого условие (14) переписется в виде

$$\left\| \frac{1}{|A|} \sum_{j \in A} \mathbf{h}_j - \frac{1}{|B|} \sum_{j \in B} \mathbf{h}_j \right\| > 2\delta. \quad (15)$$

Обозначим через $\mathcal{M}(m, t, \delta)$ максимальную мощность (t, δ) -мультимедийного кода со свойством полного поиска коалиций, устойчивого к атаке усреднения и δ -шуму. Определим, как обычно, соответствующую максимальную скорость

$$\mathcal{R}(m, t, \delta) := m^{-1} \log_2 \mathcal{M}(m, t, \delta).$$

Основным результатом статьи является доказательство существования соответствующих (t, δ) -мультимедийных кодов со скоростью, отделенной от нуля.

Теорема 1. Для любых фиксированных t и δ

$$\liminf_m \mathcal{R}(m, t, \delta) \geq \frac{\gamma_t \log_2 e}{t(1 + \gamma_t \log_2 e)} > \frac{1}{t(1 + e(\log_2 e)^{-1})} > \frac{0,346}{t}, \quad (16)$$

где $\gamma_t = (1 - t^{-1})^{t-1}$.

Разобьем построение таких кодов на две подзадачи: первая, когда мощность коалиции заранее известна, вторая – построение кодов, которые позволяют найти мощность коалиции по любой ложной копии. Начнем с первой подзадачи.

Если мощность коалиции известна, то возникающая задача – это по существу задача о сигнатурном коде для двоичного суммирующего канала с шумом (ДСКШ). А именно, дадим

Определение 3. Двоичный код $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$ называется $(=, t, \Delta)$ -сигнатурным кодом для ДСКШ, если для любых двух различных подмножеств кода $A, B \subset \mathcal{C}$, таких что $|A| = |B| \leq t$, справедливо неравенство

$$\left\| \sum_{\mathbf{c} \in A} \mathbf{c} - \sum_{\mathbf{c}' \in B} \mathbf{c}' \right\| > 2\Delta. \quad (17)$$

Очевидно, что $(=, t, \Delta)$ -сигнатурный код для ДСКШ позволяет однозначно найти всю коалицию при атаке усреднения и δ -шуме, где $\delta = \Delta/t$, если мощность коалиции заранее известна и не превышает t . Заметим также, что в определении 3 условие $|A| = |B| \leq t$ можно без ограничения общности заменить на условие $|A| = |B| = t$.

Воспользуемся конструкцией, предложенной в [19] для построения сигнатурных кодов для двоичного суммирующего по модулю 2 канала, и перероткнутой в [20] как коды, исправляющие ошибки в канале и синдроме.

Начнем с кода $\widehat{\mathcal{H}} = \{\widehat{\mathbf{h}}_1, \dots, \widehat{\mathbf{h}}_M\} \subset \{0, 1\}^m$ длины m , слова которого – это столбцы проверочной $(m \times M)$ -матрицы линейного двоичного кода V , исправляющего t ошибок. Закодируем слова $\widehat{\mathbf{h}}_1, \dots, \widehat{\mathbf{h}}_M$ двоичным кодом U длины n с t информационными символами и минимальным кодовым расстоянием d (в метрике Хэмминга). Обозначим полученные векторы через $\mathbf{h}_1, \dots, \mathbf{h}_M$ и зададим код $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\}$.

Лемма 1. Код $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\} \subset \{0, 1\}^n$ является $(=, t, \Delta)$ -сигнатурным кодом для ДСКШ с $\Delta = \sqrt{d}/2$.

Доказательство. Рассмотрим два произвольных различных подмножества A, B кода \mathcal{H} одинаковой мощности не более t и соответствующие им векторы

$$\mathbf{h}^{(A)} = \sum_{\mathbf{h} \in A} \mathbf{h} \quad \text{и} \quad \mathbf{h}^{(B)} = \sum_{\mathbf{h} \in B} \mathbf{h}.$$

Будем обозначать через

$$\mathbf{a} \bmod 2 = (a_1 \bmod 2, \dots, a_n \bmod 2) \in \{0, 1\}^n$$

двоичную проекцию произвольного целочисленного вектора $\mathbf{a} = (a_1, \dots, a_n)$. Так как различные суммы по модулю 2 из t и менее векторов $\widehat{\mathbf{h}}_j$ различны (и отличны от нуля), то это же справедливо и для сумм по модулю 2 векторов \mathbf{h}_j . Следовательно, $\mathbf{h}^{(A)} \bmod 2 \neq \mathbf{h}^{(B)} \bmod 2$. Так как векторы $\mathbf{h}^{(A)} \bmod 2$ и $\mathbf{h}^{(B)} \bmod 2$ принадлежат коду U (в силу линейности кода), то

$$d_H(\mathbf{h}^{(A)} \bmod 2, \mathbf{h}^{(B)} \bmod 2) \geq d.$$

Теперь остается применить неравенство

$$d_E(\mathbf{a}, \mathbf{b}) \geq \sqrt{d_H(\mathbf{a} \bmod 2, \mathbf{b} \bmod 2)}, \quad (18)$$

справедливое для любых двух целочисленных векторов $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, где через d_E обозначено евклидово расстояние, а через d_H – расстояние Хэмминга.

Таким образом, код $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\}$ позволяет найти полностью коалицию недобросовестных пользователей в условиях атаки усреднения и целенаправленного шума длины не более $\delta = (2t)^{-1}\sqrt{d}$, если число таких пользователей заранее известно и не превосходит t . ▲

§ 4. Коды, определяющие мощность коалиции

Определение 4. Будем называть двоичный код C *кодом, определяющим мощность коалиции вплоть до t в условиях δ -шума*, если для любых двух его подмножеств A, B различной мощности не более t справедливо неравенство

$$\left\| \frac{1}{|A|} \sum_{c \in A} \mathbf{c} - \frac{1}{|B|} \sum_{c' \in B} \mathbf{c}' \right\| > 2\delta. \quad (19)$$

С помощью такого кода дистрибьютор сформирует итоговый код, приписывая в качестве “хвостов” к словам кода \mathcal{H} , построенного выше, слова кода, определяющего мощность. По “хвостам” дистрибьютор найдет мощность коалиции, а затем по коду \mathcal{H} – и саму коалицию. Параметры получаемых кодов мы оценим в самом конце статьи.

Введем, как нам представляется, новое понятие в комбинаторной теории кодирования.

Определение 5. Двоичный код называется *слабым t -дизъюнктивным кодом*, если для любых t' различных кодовых слов, $1 \leq t' \leq t$, существует координата, в которой ровно одно слово равно 1.

Иначе говоря, для любого кодового подмножества A , такого что $1 \leq |A| \leq t$, существует координата i , такая что $\pi_i(A) = 1$, где $\pi_i(A) := |\{\mathbf{a} \in A : a_i = 1\}|$.

Отметим очевидную связь с дизъюнктивными кодами (см. [21, 22]). Напомним, что код называется t -дизъюнктивным кодом, если для любого кодового подмножества A мощности не более t и любого кодового слова $\mathbf{b} \notin A$ существует координата i , такая что $a_i = 0$ для всех $\mathbf{a} \in A$ и $b_i = 1$. Очевидно, что $(t-1)$ -дизъюнктивный код является слабым t -дизъюнктивным кодом, так как для любых t различных кодовых слов существует t соответствующих координат i , в проекциях на которые встречаются все t слов веса Хэмминга 1, а для свойства слабой дизъюнктивности достаточно всего одного такого слова и одной такой координаты, но для любых t' слов, $1 \leq t' \leq t$.

Дизъюнктивные коды были переоткрыты в экстремальной комбинаторике под названием семейства множеств без t -перекрытий [23, 24]. Напомним, что семейство множеств $\mathcal{F} = \{F_1, \dots, F_M\}$ называется семейством без t -перекрытий, если ни одно множество семейства не покрывается объединением t других множеств этого семейства.

Обобщим понятие слабого t -дизъюнктивного кода, введя, по существу, соответствующее расстояние.

Определение 6. Двоичный код называется *слабым (t, T) -дизъюнктивным кодом*, если для любого кодового подмножества A , такого что $2 \leq |A| \leq t$, существует не менее T координат i , таких что $\pi_i(A) = 1$.

В качестве примера заметим, что слабый $(2, 1)$ -дизъюнктивный код – это любой двоичный код, состоящий из различных слов, а слабый $(2, T)$ -дизъюнктивный код – это двоичный код с минимальным кодовым расстоянием (в метрике Хэмминга) не менее T .

Обозначим через $F(t, T; L)$ максимальное число слов в слабом (t, T) -дизъюнктивном коде длины L , а через $R_F(t, T; L) = L^{-1} \log_2 F(t, T; L)$ – соответствующую максимальную скорость. Докажем следующий аналог границы Варшамова–Гилберта для слабых (t, T) -дизъюнктивных кодов.

Теорема 2. Для фиксированных t и T

$$\liminf_L R_F(t, T; L) \geq \frac{1}{t} \left(1 - \frac{1}{t}\right)^{t-1} \log_2 e > \frac{\log_2 e}{et}. \quad (20)$$

Доказательство. Рассмотрим случайный двоичный код C мощности M и длины L , в котором координаты кодовых слов выбираются независимо с вероятностью $p = t^{-1}$ того, что координата равна 1. Возьмем произвольное множество A из t' слов кода C , $2 \leq t' \leq t$, и посчитаем вероятность \mathcal{P} того, что для i -й координаты $\pi_i(A) = 1$. Очевидно, что

$$\mathcal{P} = \mathcal{P}(t') = t' p (1 - p)^{t'-1}. \quad (21)$$

Тогда для вероятности p_{bad} того, что у A нет искомым T координат, справедливо

$$p_{\text{bad}} = p_{\text{bad}}(t') = \sum_{i=0}^{T-1} C_L^i \mathcal{P}^i (1 - \mathcal{P})^{L-i} = (1 - \mathcal{P})^L \sum_{i=0}^{T-1} C_L^i \left(\frac{\mathcal{P}}{1 - \mathcal{P}}\right)^i. \quad (22)$$

В силу того, что $\mathcal{P} \leq 1/2$ и

$$\sum_{i=0}^{T-1} C_L^i < L^T$$

при $L > 1$, получаем, что $p_{\text{bad}} = p_{\text{bad}}(t') \leq (1 - \mathcal{P})^L L^T$.

Так как имеется всего $C_M^{t'} < M^{t'}$ различных t' -подмножеств, то для вероятности P_{bad} того, что существует хотя бы одно “плохое” t' -подмножество A , т.е. такое, у которого нет T искомым координат, справедливо неравенство $P_{\text{bad}} < M^{t'} p_{\text{bad}}$. Потребуем, чтобы

$$P_{\text{bad}}(t') < M^{t'} p_{\text{bad}}(t') \leq (2t)^{-1}, \quad (23)$$

для чего достаточно, чтобы выполнялось неравенство

$$M = M(t') \leq (2tL^T (1 - \mathcal{P}(t'))^L)^{-1/t'}. \quad (24)$$

Оценим теперь скорость кода $R(t') := \frac{1}{L} \log_2 M(t')$:

$$\begin{aligned} R(t') &\geq -\frac{1}{t'} (\log_2(1 - \mathcal{P}(t')) + o(1)) \geq \frac{1}{t'} \mathcal{P}(t') \log_2 e + o(1) = \\ &= t^{-1} (1 - t^{-1})^{t'-1} \log_2 e + o(1), \end{aligned} \quad (25)$$

где второе неравенство следует из оценки $\ln(1 + x) \leq x$. Так как правая часть (25) монотонно убывает по t' , то выберем $t' = t$ и итоговую скорость кода

$$R := \frac{1}{t} \left(1 - \frac{1}{t}\right)^{t-1} \log_2 e + o(1) > \frac{\log_2 e}{et} + o(1).$$

Тогда для случайного кода с данной скоростью вероятность того, что при некотором t' не выполнено требуемое условие – для любого t' -подмножества кода существует как минимум T искомым координат – не превосходит $(2t')^{-1}$ (см. (23)). Так как таких условий не более t , то следовательно, с вероятностью не меньше $1/2$ случайный код является слабым (t, T) -дизъюнктивным кодом. \blacktriangle

Перейдем теперь к доказательству того, что слабый (t, T) -дизъюнктивный код позволяет найти мощность коалиции в присутствии шума и, более того, приведем алгоритм вычисления мощности коалиции.

Рассмотрим следующие *непересекающиеся* отрезки на числовой оси:

$$S_k := \left[\frac{1}{k} - \frac{1}{2t^2}, \frac{1}{k} + \frac{1}{2t^2} \right], \quad k = 1, 2, \dots, t.$$

Пусть C – слабый (t, T) -дизъюнктивный код, $A \subset C$ – некоторая коалиция мощности не более t и $\hat{\mathbf{y}} = \mathbf{x} + |A|^{-1} \mathbf{w}^{(A)} + \mathbf{e}$ – ложная копия \mathbf{x} , распространяемая этой коалицией (напомним, что $\mathbf{w}^{(A)} := \sum_{\mathbf{w} \in A} \mathbf{w}$). Так как дистрибьютор знает \mathbf{x} , то он может вычислить $\mathbf{z} := \hat{\mathbf{y}} - \mathbf{x}$.

Алгоритм нахождения мощности коалиции

1. Вычислить $q_k := |\{i : z_i \in S_k\}|$.
2. Положить мощность коалиции равной максимальному k , такому что $q_k \geq T/2$, $k \leq t$.

Лемма 2. Для любого слабого (t, T) -дизъюнктивного кода алгоритм правильно находит мощность коалиции, если длина шума

$$\|\mathbf{e}\| \leq \delta = \frac{\sqrt{T}}{2\sqrt{2}} t^{-2}.$$

Доказательство. Из определения кода и того, что $\mathbf{z} := \hat{\mathbf{y}} - \mathbf{x} = |A|^{-1} \mathbf{w}^{(A)} + \mathbf{e}$, следует, что как минимум $T/2$ координат z_i попадут в отрезок $S_{|A|}$. Действительно, из определения слабого (t, T) -дизъюнктивного кода следует, что по меньшей мере T координат вектора $|A|^{-1} \mathbf{w}^{(A)}$ равны $1/|A|$, и следовательно, если более $T/2$ этих координат вектора \mathbf{z} не принадлежат отрезку $S_{|A|}$, то для квадрата длины вектора шума справедливо неравенство

$$\|\mathbf{e}\|^2 > \frac{T}{2} \left(\frac{1}{2t^2} \right)^2 = \delta^2,$$

что противоречит предположению $\|\mathbf{e}\| \leq \delta$.

Предположим, что есть такое k , что $|A| < k \leq t$ и $q_k \geq T/2$. Пусть координата z_i попала в отрезок S_k . Если $w_i^{(A)} \geq 1$, то

$$|e_i| \geq \frac{w_i^{(A)}}{|A|} - \left(\frac{1}{k} + \frac{1}{2t^2} \right) \geq \frac{1}{|A|} - \left(\frac{1}{k} + \frac{1}{2t^2} \right) \geq \frac{1}{t(t-1)} - \frac{1}{2t^2} > \frac{1}{2t^2}.$$

Если же $w_i^{(A)} = 0$, то

$$|e_i| \geq t^{-1} - (2t^2)^{-1} > \frac{1}{t^2} - \frac{1}{2t^2} > \frac{1}{2t^2}.$$

Общее число таких координат, которые могли бы привести к неправильному определению $|A|$, равно q_k . Тогда $\|\mathbf{e}\|^2 > q_k (2t^2)^{-2}$, а так как по предположению леммы

$\|e\|^2 \leq \frac{T}{2}(2t^2)^{-2}$, то $q_k < T/2$, и следовательно, алгоритм не может выдать k как свой выход, что и требовалось доказать. ▲

§ 5. Выбор кодов

Естественно выбрать параметры d и T таким образом, чтобы обеспечиваемый леммами 1 и 2 уровень шума δ совпадал. Таким образом,

$$\delta = \frac{\sqrt{d}}{2t} = \frac{\sqrt{T}}{2\sqrt{2}t^2},$$

или, что то же самое, $T = 2dt^2$. Всюду далее и t , и δ – константы.

Напомним, что мы строим код следующим образом. Мы начинаем с t -МППК-кода мощности $M = 2^t$ и длины $m = \ell t$. Затем мы удлиняем слова этого кода, рассматривая их как информационные последовательности для линейного (n, m) -кода V с минимальным кодовым расстоянием d . Известная в этом случае избыточность $r_V = n - m$ кода V имеет порядок $\frac{d}{2} \log_2 n$ и достигается, если взять соответствующие двоичные коды БЧХ или Гоппы. Заметим, что это удлинение не влияет на асимптотику итоговой длины кода. Следующее удлинение (конкатенация) происходит с помощью слабого (t, T) -дизъюнктивного кода длины L . Ясно, что дистрибутору следует выбрать длины кодов так, чтобы мощности кодов были (примерно) равными, т.е.

$$t^{-1}n \approx t^{-1}L\gamma_t \log_2 e, \quad \text{где } \gamma_t = (1 - t^{-1})^{t-1}.$$

Итоговая длина кода \tilde{m} равна $n + L$, и следовательно, для скорости $\mathcal{R}(\tilde{m}, t, \delta)$ наилучшего (t, δ) -мультимедийного кода со свойством полного поиска коалиций, устойчивого к атаке усреднения и δ -шуму, справедлива следующая асимптотическая оценка:

$$\mathcal{R}(\tilde{m}, t, \delta) \geq \frac{\gamma_t \log_2 e}{t(1 + \gamma_t \log_2 e)} + o(1) > \frac{1}{t(1 + e(\log_2 e)^{-1})} + o(1) > \frac{0,346}{t} + o(1),$$

что и завершает доказательство теоремы 1. ▲

Таким образом, мы доказали существование мультимедийных кодов со скоростью не меньше $0,346t^{-1}$, способных находить целиком коалицию из не более чем t недобросовестных пользователей, которые применяют атаку усреднения и целенаправленный шум ограниченной евклидовой длины. Отметим, что главный член асимптотики скорости кода (см. выше) зависит от t , но не зависит от длины шума δ .

§ 6. Заключение

В заключение следует отметить, что рассмотренные в этой статье задачи близки к задаче о евклидовых “дизъюнктивных” кодах, см. [25, 26], которую можно получить, если неравенство (14) в определении 2 заменить на

$$\left\| \sum_{j \in A} \mathbf{w}_j - \sum_{j \in B} \mathbf{w}_j \right\| > 2\delta. \quad (26)$$

Другое отличие состоит в том, что в задаче из [25, 26] в качестве \mathbf{w}_j рассматривались произвольные векторы евклидова пространства, а в данной статье – только двоичные.

Авторы считают своим приятным долгом выразить благодарность И.В. Воробьеву за полезные обсуждения.

СПИСОК ЛИТЕРАТУРЫ

1. *Fan J., Gu Y., Hachimori M., Miao Y.* Signature Codes for Weighted Binary Adder Channel and Multimedia Fingerprinting // IEEE Trans. Inform. Theory. 2020 (to appear).
2. *Wagner N.R.* Fingerprinting // Proc. 1983 IEEE Symp. on Security and Privacy. Oakland, CA, USA. April 25–27, 1983. P. 18–22.
3. *Blakley G.R., Meadows C., Purdy G.B.* Fingerprinting Long Forgiving Messages // Advances in Cryptology—CRYPTO'85 (Proc. Conf. on the Theory and Application of Cryptographic Techniques. Santa Barbara, CA, USA. August 18–22, 1985). Lect. Notes Comp. Sci. V. 218. Berlin: Springer, 1986. P. 180–189.
4. *Chor B., Fiat A., Naor M.* Tracing Traitors // Advances in Cryptology—CRYPTO'94 (Proc. 14th Annu. Int. Cryptology Conf. Santa Barbara, CA, USA. August 21–25, 1994). Lect. Notes Comp. Sci. V. 839. Berlin: Springer, 1994. P. 257–270.
5. *Hollmann H.D.L., van Lint J.H., Linnartz J.-P., Tolhuizen L.M.G.M.* On Codes with the Identifiable Parent Property // J. Combin. Theory Ser. A. 1998. V. 82. № 2. P. 121–133.
6. *Boneh D., Shaw J.* Collusion-Secure Fingerprinting for Digital Data // IEEE Trans. Inform. Theory. 1998. V. 44. № 5. P. 1897–1905.
7. *Barg A., Blakley G.R., Kabatiansky G.A.* Digital Fingerprinting Codes: Problem Statements, Constructions, Identification of Traitors // IEEE Trans. Inform. Theory. 2003. V. 49. № 4. P. 852–865.
8. *Tardos G.* Optimal Probabilistic Fingerprint Codes // Proc. 35th Annu. ACM Symp. on Theory of Computing (STOC'03). San Diego, CA, USA. June 9–11, 2003. P. 116–125.
9. *Trappe W., Wu M., Wang Z.J., Liu K.J.R.* Anti-Collusion Fingerprinting for Multimedia // IEEE Trans. Signal Process. 2003. V. 51. № 4. P. 1069–1087.
10. *Liu K.J.R., Trappe W., Wang Z.J., Wu M., Zhao H.* Multimedia Fingerprinting Forensics for Traitor Tracing. Cairo, Egypt: Hindawi, 2005.
11. *Chang S.C., Wolf J.K.* On the T -User M -Frequency Noiseless Multiple-Access Channel with and without Intensity Information // IEEE Trans. Inform. Theory. 1981. V. 27. № 1. P. 41–48.
12. *Cheng M., Miao Y.* On Anti-Collusion Codes and Detection Algorithms for Multimedia Fingerprinting // IEEE Trans. Inform. Theory. 2011. V. 57. № 7. P. 4843–4851.
13. *Egorova E., Fernandez M., Kabatiansky G., Lee M.H.* Signature Codes for the A-Channel and Collusion-Secure Multimedia Fingerprinting Codes // Proc. 2016 IEEE Int. Symp. on Information Theory (ISIT'2016). Barcelona, Spain. July 10–15, 2016. P. 3043–3047.
14. *Egorova E., Fernandez M., Kabatiansky G., Lee M.H.* Signature Codes for Weighted Noisy Adder Channel, Multimedia Fingerprinting and Compressed Sensing // Des. Codes Cryptogr. 2019. V. 87. № 2–3. P. 455–462.
15. *Györfi L., Györi S., Laczay B., Ruszinkó M.* Lectures on Multiple Access Channels. Book draft, 2005. Available at http://www.szit.bme.hu/~gyori/AFOSR_05/book.pdf.
16. *D'yakov A.G.* On a Search Model of False Coins // Topics in Information Theory (Proc. 2nd Colloq. on Information Theory. Keszthely, Hungary. August 25–30, 1975). Colloq. Math. Soc. János Bolyai. V. 16. Amsterdam: North Holland, 1977. P. 163–170.
17. *Мак-Вильямс Ф.Дж., Слоэн Н.Дж.А.* Теория кодов, исправляющих ошибки. М.: Связь, 1979.
18. *Kabatiansky G., Fernandez M., Egorova E.* Multimedia Fingerprinting Codes Resistant against Colluders and Noise // Proc. 8th IEEE Int. Workshop on Information Forensics and Security (WIFS'2016). Abu Dhabi, UAE. December 4–7, 2016. P. 1–5.
19. *Ericson T., Levenshtein V.I.* Superimposed Codes in the Hamming Space // IEEE Trans. Inform. Theory. 1994. V. 40. № 6. P. 1882–1893.
20. *Влэдуч С.Г., Кабатянский Г.А., Ломаков В.В.* Об исправлении ошибок при искажениях в канале и синдроме // Пробл. передачи информ. 2015. Т. 51. № 2. С. 50–56.

21. *Kautz W.H., Singleton R.C.* Nonrandom Binary Superimposed Codes // IEEE Trans. Inform. Theory. 1964. V. 10. № 4. P. 363–377.
22. *Дьячков А.Г., Рыков В.В.* Границы длины дизъюнктивных кодов // Пробл. передачи информ. 1982. Т. 18. № 3. С. 7–13.
23. *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. V. 33. № 2. P. 158–166.
24. *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. V. 51. № 1–2. P. 79–89.
25. *Ericson T., Györfi L.* Superimposed Codes in \mathbb{R}^n // IEEE Trans. Inform. Theory. 1988. V. 34. № 4. P. 877–880.
26. *Füredi Z., Ruszinkó M.* An Improved Upper Bound of the Rate of Euclidean Superimposed Codes // IEEE Trans. Inform. Theory. 1999. V. 45. № 2. P. 799–802.

Егорова Елена Евгеньевна
 Сколковский институт науки и технологий (Сколтех)
 egorovahelene@gmail.com
Фернандес Марсель
 Политехнический университет Каталонии,
 Барселона, Испания
 marcelf@entel.upc.edu
Кабатянский Григорий Анатольевич
 Сколковский институт науки и технологий (Сколтех)
 g.kabatyansky@skoltech.ru
Мяо Ин
 Университет Цукубы, Цукуба, префектура Ибараки, Япония
 miao@sk.tsukuba.ac.jp

Поступила в редакцию
 23.10.2020
 После доработки
 24.11.2020
 Принята к публикации
 24.11.2020