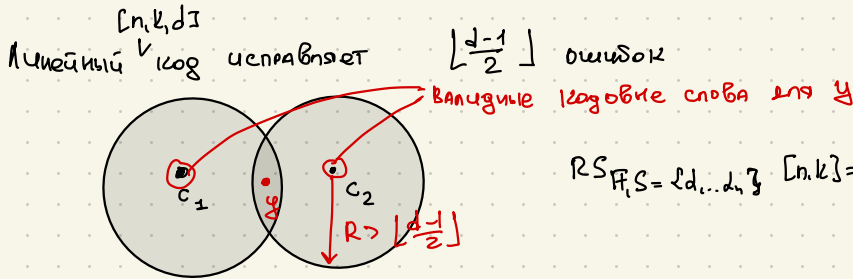


Лекция 89

Списочное декодирование Приложение теории кодирования в вычислительной биологии

I. Списочное декодирование



$$RS_{\mathbb{F}, S} = \{d_1 \dots d_n\} \quad [n, k] = \{(p(d_1) \dots p(d_n)) \in \mathbb{F}^n \mid p \in \mathbb{F}[x], \deg p \leq k-1\}$$

ЗАДАЧА ПОИСКА ДВУХ МНОГОЧЛЕНОВ: Положим $p_1(x), p_2(x) \in \mathbb{F}[x]$,

$\deg p_1(x) = \deg p_2(x) = k-1$; $n \geq 4k$ — чётное; $d_1 \dots d_n$ — различны.

$\exists T \subset \{1 \dots n\}$, $|T| = \frac{n}{2}$. Далее, пусть нам дан $y \in \mathbb{F}^n$

$$y_i = \begin{cases} p_1(d_i) & , i \in T \\ p_2(d_i) & , i \in \{1 \dots n\} \setminus T \end{cases}$$

ЗАДАЧА состоит в вычислении $p_1(x), p_2(x)$ по данным парам $\{(d_i, y_i)\}_{i=1}^n$

Связь с декодированием RS: $(y_i)_{i=1}^n$ — полученное слово, полагая что

$p_1(x)$ — исходное сообщение. Однако, $p_1(x)$ не совпадает с y на $\frac{n}{2}$

значениях; $\frac{n}{2} > \frac{d-1}{2} = \lfloor \frac{n-k+1-1}{2} \rfloor = \lfloor \frac{n-k}{2} \rfloor$. Поэтому обычные

Алг-мы декодирования когда RS не подходит.

Идея списочного декодирования $(y_i - p_1(d_i))(y_i - p_2(d_i)) = 0 \quad \forall i$

Положим, мн-м $Q(x, y) = (y - p_1(x))(y - p_2(x)) = y^2 - \underbrace{(p_1(x) + p_2(x))}_{B(x)} y + \underbrace{p_1(x) \cdot p_2(x)}_{C(x)}$

$$= y^2 - B(x)y + C(x), \quad \text{где } \deg B(x) = k-1, \quad \deg C(x) = 2(k-1)$$

$Q(d_i, y_i) = 0 \quad \forall i \leq n.$

Алгоритм

Шаг 1

Составить систему лнн. ур-ий из $Q(d_i, y_i) = 0$;
 $\{b_0, \dots, b_{k-1}, c_0, \dots, c_{2(k-1)}\}$ - неизвестные; и n ур-ий. Т.к. $n \geq 4k$,
коэф-ты $B(x)$ коэф-ты $C(x)$

система будет иметь решение. Получим $B(x), C(x)$ в явном виде

Сложность: $O(n^w)$, $2 \leq w \leq 3$

Шаг 2

Факторизуем мнн $Q(x, y) = (y - f_1(x))(y - f_2(x))$.
Верны $f_1(x), f_2(x)$.

Сложность: факторизация мн-ов от двух переменных степени d
на \mathbb{F} : $O(d^5 \cdot \lg |F|)$

Корректность

На шаге I мы всегда отыщем какие-либо $B(x), C(x)$, \exists
решение $B(x) = p_1(x) + p_2(x)$, $C(x) = p_1(x) \cdot p_2(x)$.

Докажем, что на IIм шаге мы получим корректные $p_1(x), p_2(x)$

Лемма

$\forall Q(x, y)$, полученного на шаге I, справедливо
 $(y - p_1(x)) \mid Q(x, y)$ и $(y - p_2(x)) \mid Q(x, y)$

◁ Док-м утверждение для $p_1(x)$

Заметим, что $Q(x, y)$ - унитарный (от y) мн-н. Для того, чтобы показать,

что $(y - \beta)$ делит Q , достаточно показать, что β - корень $Q \Rightarrow$

$Q(\beta) = 0$. Чтобы показать, $y - p_1(x) \mid Q(x, y)$, покажем, что

$$Q(x, p_1(x)) = 0$$

$$\neq R(x) := Q(x, p_1(x)), \deg R(x) \leq 2(k-1)$$

$$\stackrel{||}{=} p_1(x)^2 - (p_1(x) + p_2(x)) \cdot p_1(x) + p_1(x) p_2(x)$$

Заметим, что $\exists \frac{n}{2} \geq 2k$ d_i , т.ч. $p_1(d_i) = y_i$. Для таких d_i , справедливо:

$$R(d_i) = Q(d_i, p_1(d_i)) = Q(d_i, y_i) = 0 \Rightarrow \text{мы нашли } \frac{n}{2} \geq 2k \text{ корней мн-на}$$

$$\text{степени } \leq 2(k-1) \Rightarrow R(x) \equiv 0.$$

II ЗАДАЧА: ГРУППОВОЕ ТЕСТИРОВАНИЕ (group testing)

Пример Даны N человек, $s \ll N$ из которых заражены.
ЗАДАЧА: выявить зараженных за min. число тестов

Тривиальное решение сделать N тестов может быть дорогим.

МОЖНО ли смешать образцы крови так, чтобы можно было выявить больных/здоровых за меньшее число тестов

ИДЕЯ

- Ассоциируем с каждым человеком кодовое слово $c \in RS_{\mathbb{F}_s}[n, k]$
- имеем $n \cdot |F|$ кодов/тестов, сформированных в матрицу $|F| \times n$, где столбцы пронумерованы $\{e_i\}$.
- образцы крови человека, ассоциированного с кодовым словом c , помещаются в коды $(c[i], c_i)$
- $d = n - k + 1$ - min. расстояние RS $\Rightarrow \forall$ два c_i, c_j отличаются на как min d позиций.

Пример $RS_{\mathbb{F}_5}$, $n=5, k=3, d=3, s \leq 2$; $|RS_{\mathbb{F}_5}| = 5^3 = 125 \Rightarrow$ макс. 125 человек

$$c^* = (c_0^*, c_1^*, c_2^*, c_3^*, c_4^*) \\ = (2, 0, 2, 1, 4) -$$

- зараженный №1

$$c^* = (1, 2, 0, 3, 3) -$$

зараженный №2

$$c = (1, 0, 2, 4, 3) -$$

здоровый

\mathbb{F}_5	c_0	c_1	c_2	c_3	c_4
0	0	5	10	15	20
1	1	6	11	16	21
2	2	7	12	17	22
3	3	8	13	18	23
4	4	9	14	19	24

Клетка - код (всего 25)

Если $s=1 \Rightarrow \exists$ 3 негативных теста для \neq здорового;

Если $s=2 \Rightarrow \exists$ 1 негативный тест — || —

Замечание

Можно показать, что для $1 \leq s \ll N$, такое групповое тестирование работает корректно, если число тестов удовлетворяет

$$T = O\left(s^2 \cdot \left(\frac{\log N}{\log s}\right)^2\right)$$