

Теория кодирования

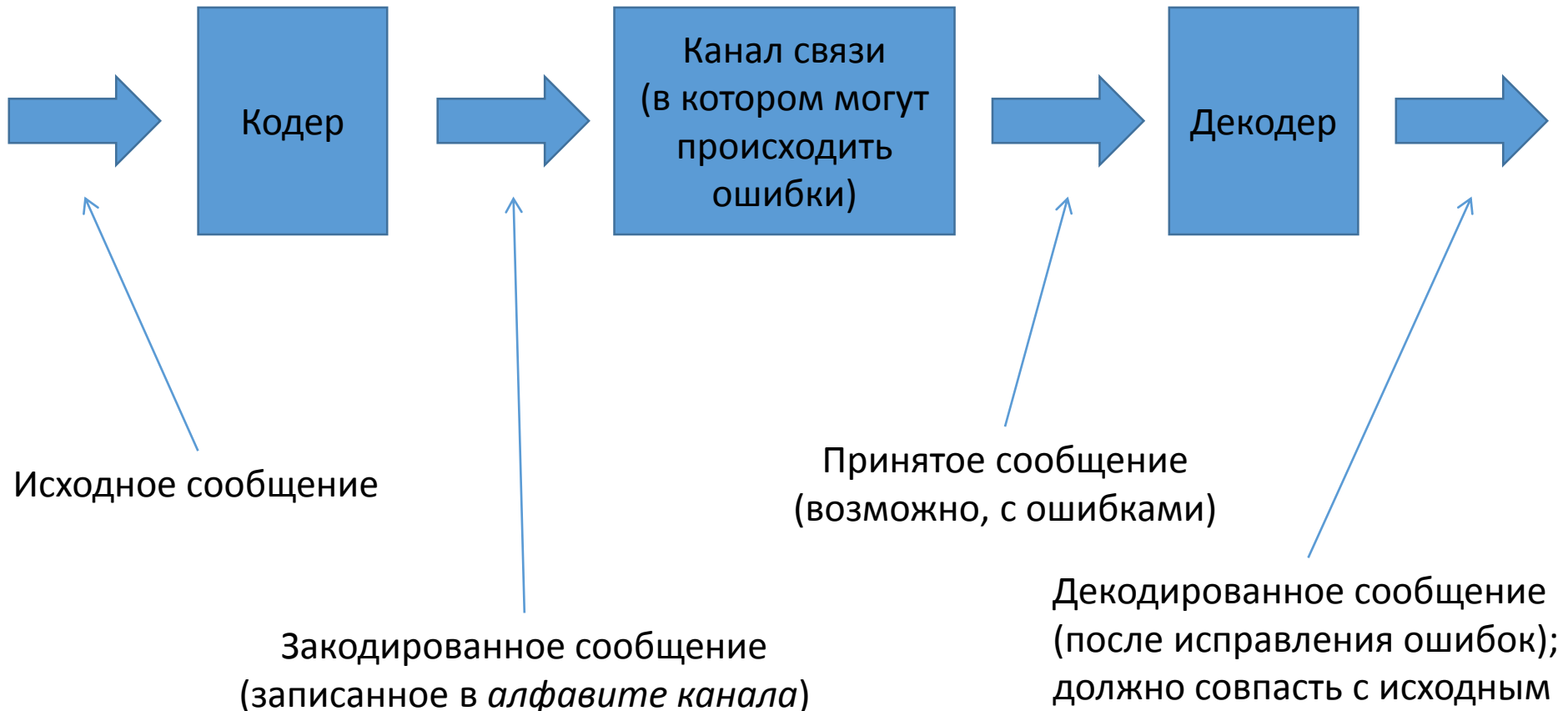
МФТИ, осень 2013

Александр Дайняк

www.dainiak.com

Коды, исправляющие ошибки

Основная модель канала связи:



Типы ошибок

- Ошибки замещения: муха → мука
 - Симметричные
 - Несимметричные
- Ошибки стирания: муха → му?а
- Ошибки выпадения: муха → уха
- Ошибки вставки: мука → мурка
- Комбинации перечисленных типов

Типы ошибок

Всегда задаются ограничения на «ненадёжность» канала, например:

- верхняя оценка числа ошибок на одно сообщение (детерминированные ограничения)
- вероятность возникновения ошибки на один символ сообщения (вероятностные ограничения)

Чаще всего алфавит канала двоичный: $\{0,1\}$

Коды

Пусть \mathbb{A}_q — алфавит канала, $|\mathbb{A}_q| = q$.

q -ичным кодом называется любое подмножество
$$C \subseteq \mathbb{A}_q^n$$

n — длина кода (длина кодовых слов)

$|C|$ — мощность кода (число кодовых слов)

Чаще всего рассматривают двоичные коды, т.е. когда $q = 2$ и $\mathbb{A}_q = \{0,1\}$.

Для произвольного двоичного слова \mathbf{a} будем через $\|\mathbf{a}\|$ обозначать вес слова, т.е. величину

$$\#\{i \mid a_i \neq 0\}$$

Обнаружение/исправление ошибок

Пусть \mathbf{a} и \mathbf{b} — слова в алфавите канала.

Обозначим через $\tilde{d}(\mathbf{a}, \mathbf{b})$ минимальное число ошибок, в результате которых \mathbf{a} может перейти в \mathbf{b} .

Способ кодирования позволяет *обнаруживать k ошибок*, если для любых различных кодовых сообщений \mathbf{a}' и \mathbf{a}'' при передаче в канал \mathbf{a}' на выходе не может получиться \mathbf{a}'' (если в канале произошло не более k ошибок).

Иначе говоря, $\tilde{d}(\mathbf{a}', \mathbf{a}'') > k$.

Обнаружение/исправление ошибок

Способ кодирования позволяет *исправлять k ошибок*, если при передаче в канал различных кодовых сообщений \mathbf{a}' и \mathbf{a}'' на выходе из канала будут получаться различные сообщения (при условии, что с каждым отдельным сообщением в канале происходит не более k ошибок).

Формально:

$$\nexists \mathbf{a}', \mathbf{a}'' \in C, \mathbf{a}: (\mathbf{a}' \neq \mathbf{a}'' \wedge \tilde{d}(\mathbf{a}', \mathbf{a}) \leq k \wedge \tilde{d}(\mathbf{a}'', \mathbf{a}) \leq k)$$

Метрика

Особенно удобно, когда \tilde{d} является *метрикой*:

- $\forall \mathbf{a}, \mathbf{b} \quad \tilde{d}(\mathbf{a}, \mathbf{b}) = \tilde{d}(\mathbf{b}, \mathbf{a})$
- $\forall \mathbf{a} \neq \mathbf{b} \quad \tilde{d}(\mathbf{a}, \mathbf{b}) > 0$
- $\forall \mathbf{a} \quad \tilde{d}(\mathbf{a}, \mathbf{a}) = 0$
- $\forall \mathbf{a}, \mathbf{b}, \mathbf{c} \quad \tilde{d}(\mathbf{a}, \mathbf{b}) \leq \tilde{d}(\mathbf{a}, \mathbf{c}) + \tilde{d}(\mathbf{c}, \mathbf{b})$

Так бывает не всегда. Например, если в канале есть только ошибки вставки и никаких других, то при $\mathbf{a} \neq \mathbf{b}$ по крайней мере одна из двух величин $\tilde{d}(\mathbf{a}, \mathbf{b})$, $\tilde{d}(\mathbf{b}, \mathbf{a})$ вовсе не определена.

Метрика Хемминга

Если рассматриваются слова одной и той же длины, а в канале возможны только ошибки типа замещения (любые), то $\tilde{d}(\mathbf{a}, \mathbf{b}) = d_X(\mathbf{a}, \mathbf{b})$, где

$$d_X(\mathbf{a}, \mathbf{b}) := \#\{i \mid a_i \neq b_i\}$$

Функционал d_X — метрика Хемминга,

$d_X(\mathbf{a}, \mathbf{b})$ — расстояние Хемминга между \mathbf{a} и \mathbf{b}

Метрика Левенштейна

Если в канале происходят ошибки выпадения/вставки, то канал описывается *метрикой Левенштейна*:

$d_L(\mathbf{a}, \mathbf{b}) := \min \# \text{ выпадений и вставок, переводящих } \mathbf{a} \text{ в } \mathbf{b}$

Например:

- $d_L(\langle \mathbf{aba} \rangle, \langle \mathbf{aa} \rangle) = 1$
- $d_L(\langle \mathbf{abbaba} \rangle, \langle \mathbf{abaab} \rangle) = 3$

Кодовое расстояние

Пусть $\tilde{d}(\cdot, \cdot)$ — метрика и C — код.

Кодовым расстоянием кода C называется величина

$$\tilde{d}(C) := \min_{\substack{a \neq b \\ a, b \in C}} \tilde{d}(a, b)$$

Кодовое расстояние определяет устойчивость к ошибкам:

- C обнаруживает t ошибок $\Leftrightarrow \tilde{d}(C) > t$
- C исправляет t ошибок $\Leftrightarrow \tilde{d}(C) > 2t$

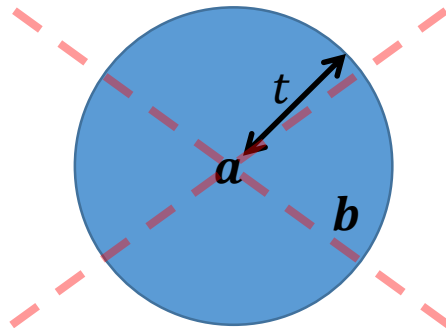
Геометрическая интерпретация

Шар радиуса r с центром в \mathbf{a} — это множество

$$S_r(\mathbf{a}) := \{\mathbf{b} \mid \tilde{d}(\mathbf{a}, \mathbf{b}) \leq r\}$$

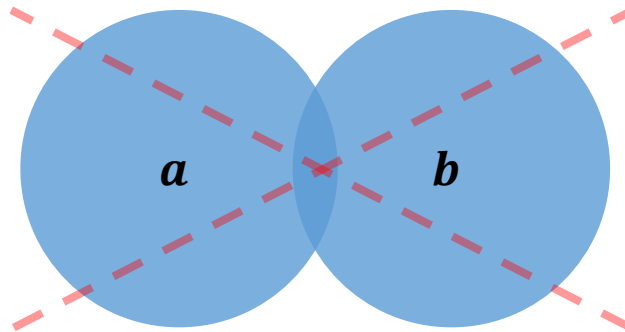
Если в канал передавалось \mathbf{a} , то на выходе из канала может быть любое слово $\mathbf{b} \in S_t(\mathbf{a})$.

Значит, код обнаруживает t ошибок т. и т.т., когда никакое кодовое слово не попадает в шар радиуса t с центром в другом кодовом слове:



Геометрическая интерпретация

Код исправляет t ошибок т. и т.т., когда при передаче в канал различных кодовых слов на выходе получаются различные слова, то есть когда шары радиуса t с центрами в кодовых словах не пересекаются:



Основные задачи теории кодов, исправляющих ошибки

Основная задача: строить коды, для которых

- число кодовых слов как можно больше,
- кодовое расстояние как можно больше,
- длина кодовых слов как можно меньше.

Задачи, связанные с ресурсами:

- Процессы кодирования и декодирования (исправление ошибок) должны быть возможно менее трудоёмкими по количеству операций и по памяти

Основные задачи теории кодов, исправляющих ошибки

Основная задача: строить коды, для которых

- число кодовых слов как можно больше,
- кодовое расстояние как можно больше,
- длина кодовых слов как можно меньше.

Геометрически, это **задача об упаковке**

- возможно большего числа шаров,
- возможно большего радиуса,
- в пространстве возможно меньшей размерности.

Коды Варшамова—Тененгольца

- Пример кодов, исправляющих ошибки выпадения/вставки
- Простой алгоритм исправления ошибок

Коды Варшамова—Тененгольца

Код Варшамова—Тененгольца длины n :

$$C := \left\{ a_1 a_2 \dots a_n \mid \sum_{i=1}^n i a_i \equiv 0 \pmod{(n+1)} \right\}$$

Для мощности кода справедлива формула (без доказательства)

$$|C| = \frac{1}{2^{n+1}} \sum_{\substack{d|(n+1) \\ d \text{ нечётно}}} \phi(d) 2^{(n+1)/d}$$

Асимптотически это максимально возможная мощность кода, исправляющего одну ошибку выпадения/вставки символа.

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Пусть C — код В.—Т. длины n , и пусть $\mathbf{a} \in C$.

Пусть в канал передали $\mathbf{a} = a_1 \dots a_n$, и на выходе получили слово

$$\mathbf{a}' := a'_1 \dots a'_{n-1} = a_1 \dots a_{k-1} a_{k+1} \dots a_n$$

(символ a_k выпал).

Наша задача: по \mathbf{a}' восстановить \mathbf{a} .

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Восстановить \mathbf{a} — не то же самое, что восстановить пару (k, a_k) .

Например, если $\mathbf{a}' = 1001$, то $\mathbf{a} = 10001$, но мы не узнаем, какой именно из нулей выпал.

Положим

$$\begin{aligned} n_0 &:= \#\{i > k \mid a_i = 0\} \\ n_1 &:= \#\{i > k \mid a_i = 1\} \end{aligned}$$

Заметим, что если $a_k = 0$, то \mathbf{a} можно восстановить по \mathbf{a}' , если известно n_1 .

Аналогично, если $a_k = 1$, то \mathbf{a} можно восстановить по \mathbf{a}' , если известно n_0 .

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Рассмотрим суммы

$$S := \sum_{i=1}^n i a_i \quad \text{и} \quad S' := \sum_{i=1}^{n-1} i a'_i$$

Заметим, что

$$\begin{aligned} S - S' &= \sum_{i=1}^n i a_i - \left(\sum_{i=1}^{k-1} i a_i + \sum_{i=k}^{n-1} i a_{i+1} \right) = \sum_{i=k}^n i a_i - \sum_{i=k+1}^n (i-1) a_i \\ &= k a_k + \sum_{i=k+1}^n a_i \end{aligned}$$

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Получаем

$$S' = S - \left(ka_k + \sum_{i=k+1}^n a_i \right) = S - ka_k - n_1$$

Так как $S \equiv 0 \pmod{(n+1)}$, то

$$S' \equiv -n_1 - ka_k \pmod{(n+1)}$$

Если $a_k = 0$, то $-S' \equiv n_1$.

Если $a_k = 1$, то

$$-S' \equiv n_1 + k = (n - k - n_0) + k = n - n_0$$

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Итак,

- если $a_k = 0$, то $(-S') \bmod (n + 1) = n_1$,
- если $a_k = 1$, то $(-S') \bmod (n + 1) = n - n_0$.

Осталось определить, чему равно a_k .

Заметим, что $\|\mathbf{a}'\| \geq n_1$,
 $\|\mathbf{a}'\| \leq (n - 1) - n_0$

Отсюда $n_1 \leq \|\mathbf{a}'\| < n - n_0$.

То есть, если $(-S') \bmod (n + 1) \leq \|\mathbf{a}'\|$, то это n_1 , а в противном случае это $n - n_0$.

Коды Варшамова—Тененгольца: исправление одной ошибки выпадения

Итоговый алгоритм восстановления \mathbf{a} по \mathbf{a}' :

- Вычисляем величину

$$T := \left(- \sum_{i=1}^{n-1} i a'_i \right) \bmod (n + 1)$$

- Если $T \leq \|\mathbf{a}'\|$, то в слово \mathbf{a}' вставляем перед T -й с конца единицей символ 0.
- Если $T > \|\mathbf{a}'\|$, то в слово \mathbf{a}' вставляем перед $(n - T)$ -м с конца нулём символ 1.

Коды Варшамова—Тененгольца: исправление одной ошибки вставки

Теперь рассмотрим задачу, когда \mathbf{a}' получено из \mathbf{a} вставкой символа:

$$\mathbf{a}' = \dots a_k x a_{k+1} \dots$$

(Если $k = 0$, то $\mathbf{a}' = x\mathbf{a}$; если $k = n$, то $\mathbf{a}' = \mathbf{a}x$)

Тогда $S' = S + (k + 1)x + \sum_{i>k} a_i$, и значит

$$S' \equiv (k + 1)x + \sum_{i>k} a_i = (k + 1)x + n_1$$

Положим $T := S' \bmod (n + 1)$.

Коды Варшамова—Тененгольца: исправление одной ошибки вставки

$$\begin{aligned} \mathbf{a}' &= \dots a_k x a_{k+1} \dots \\ S' &\equiv (k+1)x + n_1 \end{aligned}$$

Положим $T := S' \bmod (n+1)$.

Есть два случая, когда $T = 0$:

- $(k+1)x + n_1 = 0$. Тогда $x = 0$ и $a_{k+1} = \dots = a_n = 0$.
- $(k+1)x + n_1 = n+1$. Тогда $x = 1$ и $a_{k+1} = \dots = a_n = 1$.

В обоих случаях \mathbf{a} получается из \mathbf{a}' удалением последнего символа.

Коды Варшамова—Тененгольца: исправление одной ошибки вставки

$$\begin{aligned} \mathbf{a}' &= \dots a_k x a_{k+1} \dots \\ T &\equiv (k + 1)x + n_1 \end{aligned}$$

Теперь рассмотрим случай, когда $T = \|\mathbf{a}'\| > 0$.

Это возможно только в одном из двух случаев:

- $a_1 = \dots = a_k = x = 0$
- $a_1 = \dots = a_k = x = 1$

В любом случае, если $T = \|\mathbf{a}'\|$, то \mathbf{a} получается из \mathbf{a}' удалением первого символа.

Коды Варшамова—Тененгольца: исправление одной ошибки вставки

$$\begin{aligned} \mathbf{a}' &= \dots a_k x a_{k+1} \dots \\ T &\equiv (k + 1)x + n_1 \end{aligned}$$

Остался случай $0 < T \neq \|\mathbf{a}'\|$.

- Если $x = 0$, то $T = n_1 < \|\mathbf{a}'\|$.
- Если $x = 1$, то $T = k + 1 + n_1 > \|\mathbf{a}'\|$.

При этом оказывается, что

$$T = k + 1 + (n - k - n_0) = n + 1 - n_0.$$

В обоих случаях нужная для восстановления \mathbf{a} информация у нас есть.

Коды Варшамова—Тененгольца

Возможные обобщения кодов В.—Т.:

- Произвольный фиксированный модуль $l > n$:

$$\{a_1 a_2 \dots a_n \mid \sum_{i=1}^n i a_i \equiv 0 \pmod{l}\}$$

- Дополнительные соотношения, например:

$$\{a_1 a_2 \dots a_n \mid \sum_{i=1}^n i a_i \equiv \sum_{i=1}^n i^2 a_i \equiv 0 \pmod{l}\}$$

Ошибки замещения

- Односторонние ошибки. Например, если в двоичном канале возможны только замещения вида $0 \rightarrow 1$ или только $1 \rightarrow 0$. Канал связи в этом случае называется *несимметричным*.
- Двусторонние (симметричные) ошибки. Если возможно замещение символов $b_1 \rightarrow b_2$, то возможно и замещение $b_2 \rightarrow b_1$.
- Канал связи, в котором ошибки только симметричные, называется *симметричным*.

Ошибки замещения

Коды Варшамова—Тененгольца могут исправлять единичные односторонние ошибки замещения:

- Если \mathbf{a}' получается из \mathbf{a} замещением i -го символа с 0 на 1, то $S' \equiv i \pmod{n+1}$
- Если \mathbf{a}' получается из \mathbf{a} замещением i -го символа с 1 на 0, то $S' \equiv -i \pmod{n+1}$

(Исправлять оба типа ошибок одновременно коды Варшамова—Тененгольца не могут.)

Обозначение кодов

Далее будем изучать коды, исправляющие ошибки замещения, значит, метрика по умолчанию — метрика Хемминга.

Обозначение кода с заданными параметрами

Если C — q -ичный код с длиной слов n , числом слов M и кодовым расстоянием d , то пишут:

« C является $(n, M, d)_q$ -кодом»

Если код двоичный, то символ q не указывают.

Граница сферической упаковки

Теорема. (Граница Хемминга (R.W. Hamming), граница сферической упаковки)

Для любого $(n, M, d)_q$ -кода имеем

$$M \leq \frac{q^n}{|S_{\lfloor (d-1)/2 \rfloor}(\mathbf{0})|}$$

В двоичном случае

$$M \leq \frac{2^n}{\sum_{k=0}^{\lfloor (d-1)/2 \rfloor} \binom{n}{k}}$$

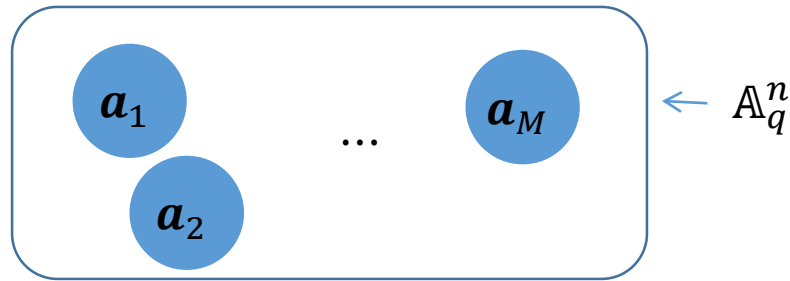
Коды, достигающие эту границу, называются *совершенными* или *плотно упакованными*.

Граница сферической упаковки

Доказательство теоремы:

Пусть $C = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_M\}$ — $(n, M, d)_q$ -код.

Так как $d(C) = d$, то шары радиуса $\lfloor (d - 1)/2 \rfloor$ с центрами в кодовых словах не пересекаются:



$$\text{Отсюда } q^n \geq \sum_{j=1}^M |S_{\lfloor (d-1)/2 \rfloor}(\mathbf{a}_j)| = M \cdot |S_{\lfloor (d-1)/2 \rfloor}(\mathbf{0})|$$

«Анти-Хемминг»

Теорема. (В некотором смысле, обратная границе Хемминга)

Пусть числа $q, n, M, d \in \mathbb{N}$ таковы, что

$$M \leq \frac{q^n}{|S_d(\mathbf{0})|}.$$

Тогда существует $(n, M, d)_q$ -код.

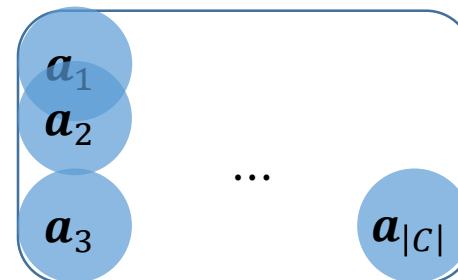
«Анти-Хемминг»

Доказательство теоремы:

Пусть $C = \{\mathbf{a}_1, \dots, \mathbf{a}_{|C|}\}$ — код максимальной мощности с кодовым расстоянием d и длиной слов n .

Тогда шары радиуса d с центрами в кодовых словах покрывают целиком множество \mathbb{A}_q^n (иначе код C можно было пополнить любым из слов, не лежащих ни в одном из этих шаров).

Отсюда $\sum_{j=1}^{|C|} |S_d(\mathbf{a}_j)| \geq q^n$,
следовательно $|C| \geq M$.



Граница Синглтона

Теорема. (R.C. Singleton)

Для любого $(n, M, d)_q$ -кода имеем

$$|C| \leq q^{n-d+1}$$

Коды, на которых достигается граница Синглтона, называются *MDS-кодами* (maximum distance separable codes).

Граница Синглтона

Доказательство:

Пусть $C = \{\mathbf{a}_1, \dots, \mathbf{a}_M\}$ — $(n, M, d)_q$ -код.

Рассмотрим слова $\{\mathbf{a}'_i\}_{i=1}^M$, где \mathbf{a}'_i получено из \mathbf{a}_i отбрасыванием $(d - 1)$ последних координат.

Так как $d(\mathbf{a}_i, \mathbf{a}_j) \geq d$ для любых i, j , то все слова \mathbf{a}'_i различны. Их количество не превосходит числа всех q -ичных слов длины $(n - d + 1)$.

Поэтому и $M \leq q^{n-d+1}$.

На лекции мы рассмотрели:

- Исправление ошибок, кодовое расстояние
- Коды Варшамова—Тененгольца
- Простые границы мощностей кодов