

بسم الله الرحمن الرحيم

جلسه ششم

خلاصه سازی برای مه داده



كران پايين



کران پایین برای تعداد اعداد متفاوت F0

الگوریتم‌های ارائه شده (مثلا FM++):

تقریبی

تصادفی

الگوریتم‌های ارائه شده (مثلا FM++):

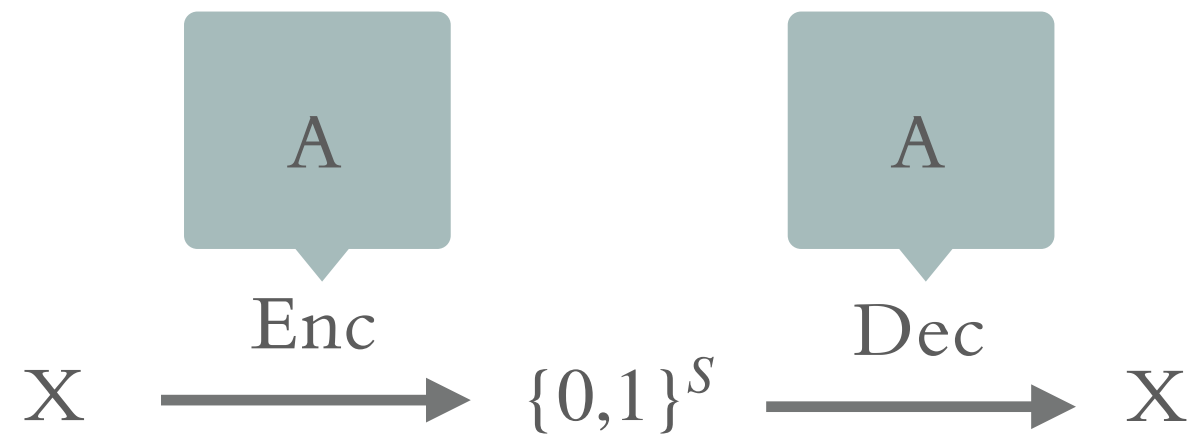
تقریبی

تصادفی

آیا هر دو لازمند؟

تکنیک برای تولید کران پایین

هدف: نمی‌تواند حافظه A کم باشد



پس: $S \geq \log |X|$

کران پایین برای حافظه الگوریتم دقیق قطعی

Theorem 3.1.1. *Suppose \mathcal{A} is a deterministic streaming algorithm that computes the number of distinct elements exactly. Then \mathcal{A} uses at least n bits of memory.*

کران پایین برای حافظه الگوریتم دقیق قطعی

Theorem 3.1.1. *Suppose \mathcal{A} is a deterministic streaming algorithm that computes the number of distinct elements exactly. Then \mathcal{A} uses at least n bits of memory.*

$\text{Enc}(x): x \in \{0,1\}^n$

۱ – یک رشته با اعدادی که اندیششان در ورودی هست (مرتب)

۲ – رشته را به الگوریتم A بدهیم.

۳ – خروجی (M) : حافظه A

کران پایین برای حافظه الگوریتم دقیق قطعی

Theorem 3.1.1. *Suppose \mathcal{A} is a deterministic streaming algorithm that computes the number of distinct elements exactly. Then \mathcal{A} uses at least n bits of memory.*

$\text{Enc}(x): x \in \{0,1\}^n$

۱ – یک رشته با اعدادی که اندیششان در ورودی هست (مرتب)

۲ – رشته را به الگوریتم A بدهیم.

۳ – خروجی (M) : حافظه A

$\text{Dec}(M)$

```
 $s \leftarrow \mathcal{A}.\text{query}()$  // support size of  $x$ , i.e.  $|\{i : x_i \neq 0\}|$   
 $x \leftarrow (0, 0, \dots, 0)$   
for  $i = 1 \dots n$ :  
     $\mathcal{A}.\text{update}(i)$  // append  $i$  to the stream  
     $r \leftarrow \mathcal{A}.\text{query}()$  // will either be  $s$  or  $s+1$   
    if  $r = s$ : // Encoder must have included  $i$ , so it wasn't a new distinct element  
         $x_i \leftarrow 1$   
     $s \leftarrow r$   
return  $x$ 
```

الگوریتم قطعی تقریبی – کران پایین

Theorem 3.1.6. *Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

کدهای اصلاح خطا

اندازه بلوک

کد اصلاح خطا: $C \subset [q]^l$

الفبا

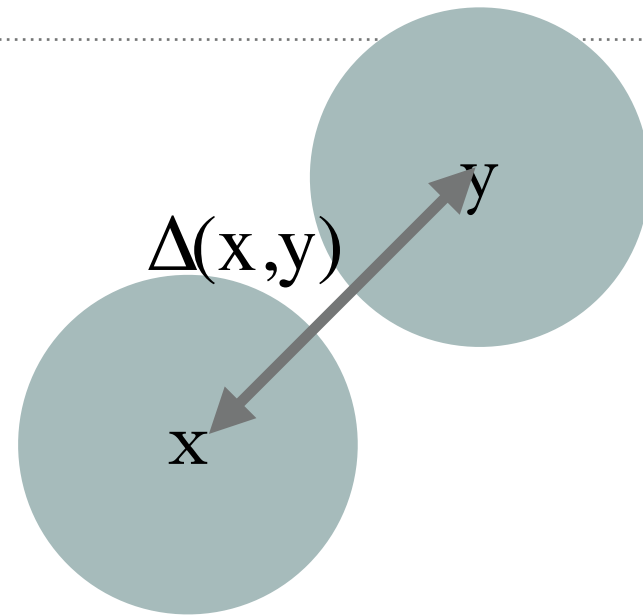
Hamming distance: $\Delta(x, y) := |\{i : x_i \neq y_i\}|$

relative Hamming distance : $\delta(x, y) = \Delta(x, y)/l$.

distance of the code : $\min_{c, c' \in C} \Delta(c, c')$

relative distance : $\min_{c, c' \in C} \delta(c, c')$

کدهای اصلاح خطا



اندازه بلوک

کد اصلاح خطا: $C \subset [q]^l$

الفبا

Hamming distance: $\Delta(x, y) := |\{i : x_i \neq y_i\}|$

relative Hamming distance : $\delta(x, y) = \Delta(x, y)/l$.

distance of the code : $\min_{c, c' \in C} \Delta(c, c')$

relative distance : $\min_{c, c' \in C} \delta(c, c')$

چرا مهم است؟

Theorem 3.1.4. *For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.*

Theorem 3.1.4. *For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.*

C: شامل n تا ci هر کدام با توزیع یکنواخت از $[q]^\ell$

Theorem 3.1.4. *For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.*

\mathcal{C} : شامل n تا c_i هر کدام با توزیع یکنواخت از $[q]^\ell$

$Y_{i,j}$: تعداد شباهت c_i و c_j

$$1/q = E[Y_{i,j}]$$

Theorem 3.1.4. *For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.*

\mathcal{C} : شامل n تا c_i هر کدام با توزیع یکنواخت از $[q]^\ell$

$Y_{i,j}$: تعداد شباهت c_i و c_j

$$1/q = E[Y_{i,j}]$$

$$P[Y_{i,j} > 6\ell/q] < 2\exp(-25/3 \cdot \ell/q)$$

Theorem 3.1.4. *For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.*

\mathcal{C} : شامل n تا c_i هر کدام با توزیع یکنواخت از $[q]^\ell$

$Y_{i,j}$: تعداد شباهت c_i و c_j

$$1/q = E[Y_{i,j}]$$

$$P[Y_{i,j} > 6\ell/q] < 2\exp(-25/3 \cdot \ell/q) < \frac{1}{n^2}$$

با انتخاب 1 بزرگ

Theorem 3.1.4. For any integers $q, n > 1$, there exists a code \mathcal{C} with $|\mathcal{C}| = n$ and block length $\ell = O(q \log n)$ with relative distance $1 - 6/q$.

\mathcal{C} : شامل n تا c_i هر کدام با توزیع یکنواخت از $[q]^\ell$

$Y_{i,j}$: تعداد شباهت c_i و c_j

$$1/q = E[Y_{i,j}]$$

$$P[Y_{i,j} > 6\ell/q] < 2\exp(-25/3 \cdot \ell/q) < \frac{1}{n^2}$$

با انتخاب 1 بزرگ

احتمال بد بودن $1 > 1$

Corollary 3.1.5. *For any integer $n > 0$ and any integers $\ell, q > 1$ such that $n = q\ell$, there exists a subset $\mathcal{B}_{q,\ell}$ of $\{0, 1\}^n$ satisfying the following properties:*

1. *Every $c \in \mathcal{B}_{q,\ell}$ has support size ℓ , i.e. $|\{i : c_i \neq 0\}| = \ell$.*

2. *For $c \neq c' \in \mathcal{B}_{q,\ell}$, $|\{i : c_i = c'_i\}| \leq 6\ell/q$.*

3. $|\mathcal{B}_{q,\ell}| = \exp(\Omega(\ell/q))$.

و ci=1

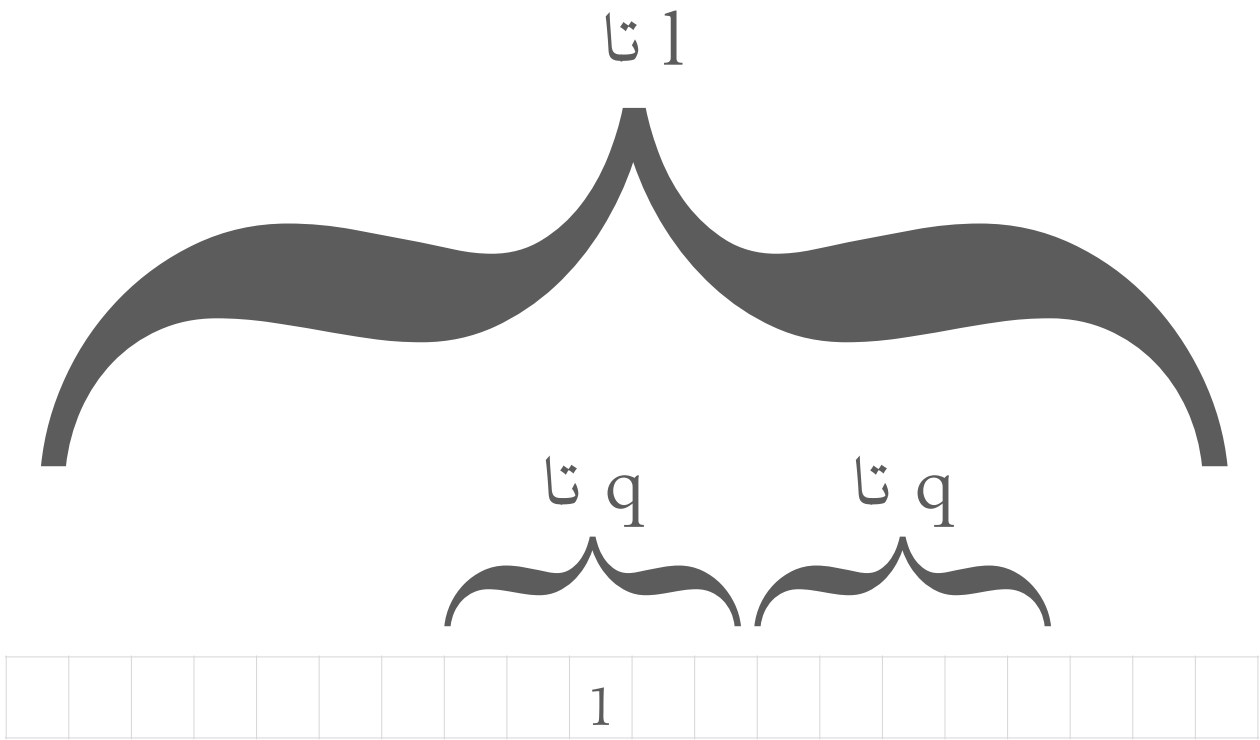
$$P[Y_{i,j} > 6\ell/q] < 2\exp(-25/3 \cdot \ell/q)$$

Corollary 3.1.5. For any integer $n > 0$ and any integers $\ell, q > 1$ such that $n = q\ell$, there exists a subset $\mathcal{B}_{q,\ell}$ of $\{0, 1\}^n$ satisfying the following properties:

- 1. Every $c \in \mathcal{B}_{q,\ell}$ has support size ℓ , i.e. $|\{i : c_i \neq 0\}| = \ell$.
- 2. For $c \neq c' \in \mathcal{B}_{q,\ell}$, $|\{i : c_i = c'_i\}| \leq 6\ell/q$.
- 3. $|\mathcal{B}_{q,\ell}| = \exp(\Omega(\ell/q))$.

و ci=1

$$P[Y_{i,j} > 6\ell/q] < 2\exp(-25/3 \cdot \ell/q)$$



Theorem 3.1.6. *Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

Theorem 3.1.6. *Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

$$\begin{array}{ccc} \text{Enc} : B_{q,1} \rightarrow \{0, 1\}^S & \longrightarrow & S \geq \log |B_{q,1}| = \Omega(n) \\ q = 100 \text{ and } l = n/q & & \end{array}$$

Theorem 3.1.6. *Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

$$\begin{array}{ccc} \text{Enc} : B_{q,1} \rightarrow \{0, 1\}^S & \longrightarrow & S \geq \log |B_{q,1}| = \Omega(n) \\ q = 100 \text{ and } l = n/q & & \end{array}$$

$$\text{Enc} : x \rightarrow \{i \mid x[i]=1\} \rightarrow A \rightarrow \text{حافظه } A$$

Theorem 3.1.6. *Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

$$\begin{array}{ccc} \text{Enc} : B_{q,1} \rightarrow \{0, 1\}^S & \longrightarrow & S \geq \log |B_{q,1}| = \Omega(n) \\ q = 100 \text{ and } 1 = n/q & & \end{array}$$

$$\text{Enc} : x \rightarrow \{i \mid x[i]=1\} \rightarrow A \rightarrow \text{حافظه } A$$

Dec:

```

for  $c \in \mathcal{B}_{q,t}$ :
     $\mathcal{A}.\text{init}(M)$  // initialize  $\mathcal{A}$ 's memory to  $M$ 
    for  $i = 1, 2, \dots, n$ :
        if  $c_i = 1$ :
             $\mathcal{A}.\text{update}(i)$ 
    if  $\mathcal{A}.\text{query}() \leq 1.9\ell$ 
        return  $c$ 

```

Theorem 3.1.6. Suppose \mathcal{A} is a deterministic streaming algorithm that always outputs a value \tilde{t} when queried such that $t \leq \tilde{t} \leq 1.9t$, where t is the number of distinct elements. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.

$$\begin{array}{ccc} \text{Enc} : B_{q,1} \rightarrow \{0, 1\}^S & \longrightarrow & S \geq \log |B_{q,1}| = \Omega(n) \\ q = 100 \text{ and } 1 = n/q & & \end{array}$$

$$\text{Enc} : x \rightarrow \{i \mid x[i]=1\} \rightarrow A \rightarrow \text{حافظه } A$$

Dec:

```

for  $c \in \mathcal{B}_{q,t}$ :
     $\mathcal{A}.\text{init}(M)$  // initialize  $\mathcal{A}$ 's memory to  $M$ 
    for  $i = 1, 2, \dots, n$ :
        if  $c_i = 1$ :
             $\mathcal{A}.\text{update}(i)$ 
    if  $\mathcal{A}.\text{query}() \leq 1.9\ell$ 
        return  $c$ 

```

اگر $c=x$:

وگرنه:

$$|c \cup x| = |c| + |x| - |c \cap x| > 2l - 6l/q = 1.94 \times l$$

الگوریتم تصادفی دقیق – کران پایین

Theorem 3.1.7. *Suppose \mathcal{A} is a randomized streaming algorithm that outputs the exact number of distinct elements with success probability at least $2/3$ for the last query in any fixed sequence of stream updates and queries. Then \mathcal{A} uses at least cn bits of memory for some constant $c > 0$.*

کم کردن احتمال خطا

احتمال خطا: 10^{-6}
حافظه: $O(S)$

$A \rightarrow A'$

احتمال درستی: $2/3$
حافظه: S

۱

$\in B_{q,l}$

همان قبلی

۲



منبع
تصادفی

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$\mathbb{E} Y_X \leq n/10^6$$

Y_X : تعداد بیت‌های متفاوت x و x'

منبع
تصادفی

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$\mathbb{E} Y_x \leq n/10^6$$

Y_x : تعداد بیت‌های متفاوت x و x'

حالت بد: $Y_x > 2l/q = 2n/q^2$

$$Y_x > 2n/10^6 : Z_x$$

منبع
تصادفی

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$E Y_X \leq n/10^6$$

Y_X : تعداد بیت‌های متفاوت x و x'

حالت بد: $Y_X > 2l/q = 2n/q^2$

$$Y_x > 2n/10^6 : Z_X$$

$$P(Z_X = 1) < 1/2$$

پس ←

منبع
تصادفی

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$E Y_X \leq n/10^6$$

Y_X : تعداد بیت‌های متفاوت x و x'

حالت بد: $Y_X > 2l/q = 2n/q^2$

$$Y_x > 2n/10^6 : Z_X$$

$$P(Z_X = 1) < 1/2$$

پس ←

$$E \sum_{x \in B_{q,l}} Z_x < \frac{1}{2} |B_{q,l}|$$

پس ←

منبع تصادفی

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$E Y_X \leq n/10^6$$

Y_X : تعداد بیت‌های متفاوت x و x'

حالت بد: $Y_X > 2l/q = 2n/q^2$

$$Y_x > 2n/10^6 : Z_X$$

$$P(Z_X = 1) < 1/2$$

پس ←

$$E \sum_{x \in B_{q,l}} Z_x < \frac{1}{2} |B_{q,l}|$$

پس ←

منبع تصادفی

حداقل نیمی از $B_{q,l}$ درست پاسخ می‌دهند

پس ←

الگوریتم تصادفی A : $A(r, \text{ورودی})$

$\in B_{q,l}$

همان قبلی

۲



$$E Y_X \leq n/10^6$$

Y_X : تعداد بیت‌های متفاوت x و x'

حالت بد: $Y_X > 2l/q = 2n/q^2$

$$Y_x > 2n/10^6 : Z_X$$

$$P(Z_X = 1) < 1/2$$

پس ←

$$E \sum_{x \in B_{q,l}} Z_x < \frac{1}{2} |B_{q,l}|$$

پس ←

منبع تصادفی

حداقل نیمی از $B_{q,l}$ درست پاسخ می‌دهند

پس ←

الگوریتم تصادفی A : $A(r, \text{ورودی})$

قطعی سازی ←