

# 海量数据流处理技术



童咏昕

计算机学院 软件开发环境国家重点实验室

[yxtong@buaa.edu.cn](mailto:yxtong@buaa.edu.cn)

# 课程提纲

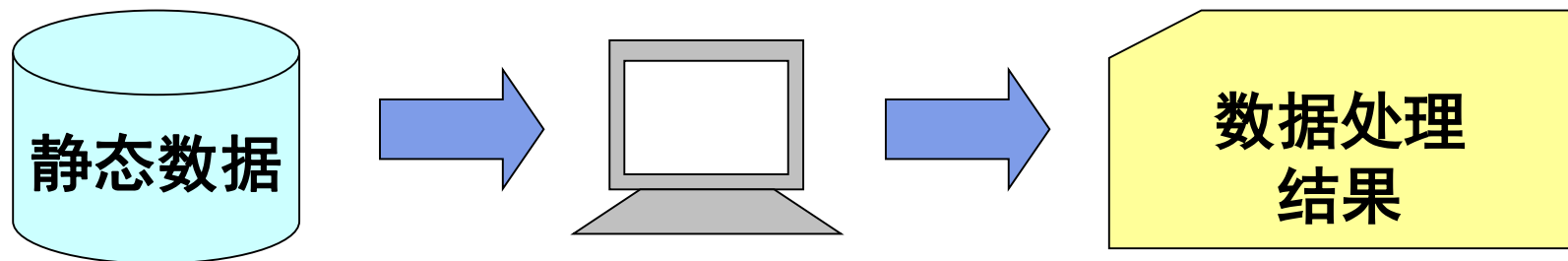
- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# 课程提纲

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

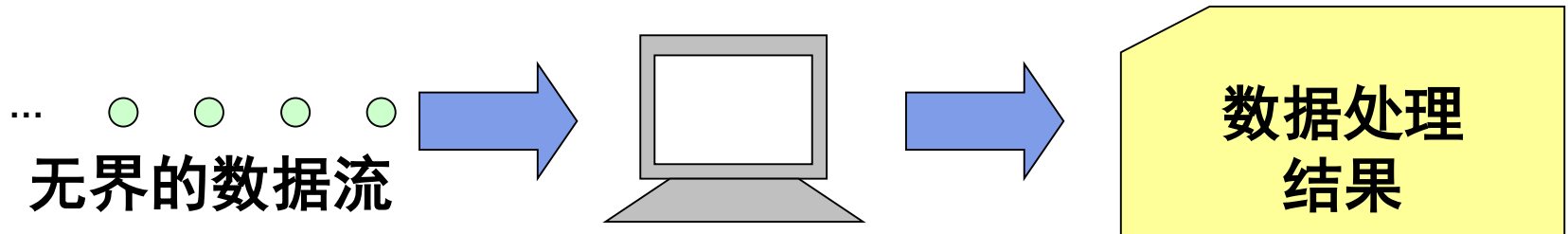
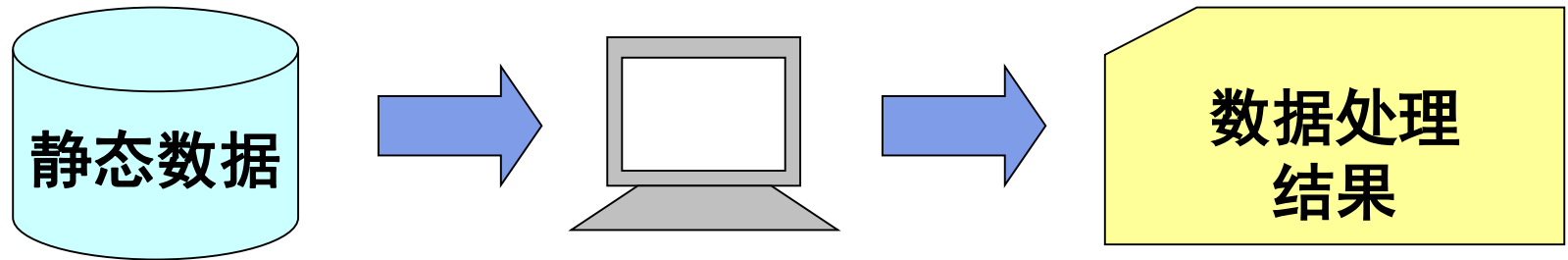
# 静态数据 vs. 流状数据

---



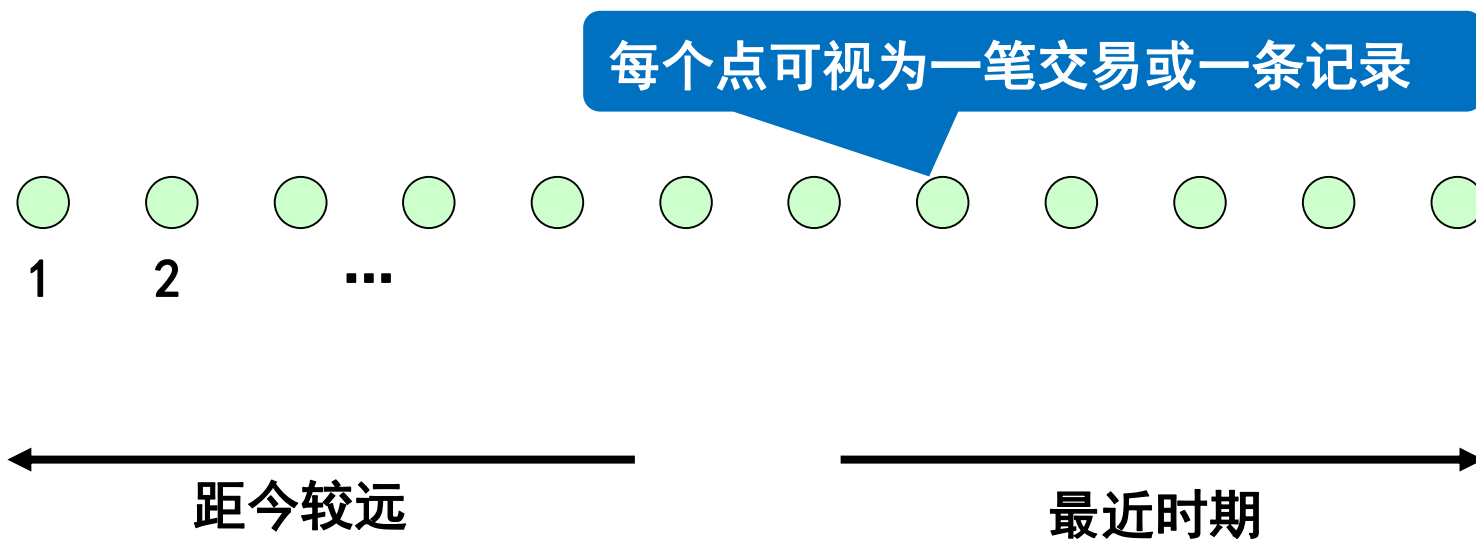
# 静态数据 vs. 流状数据

---



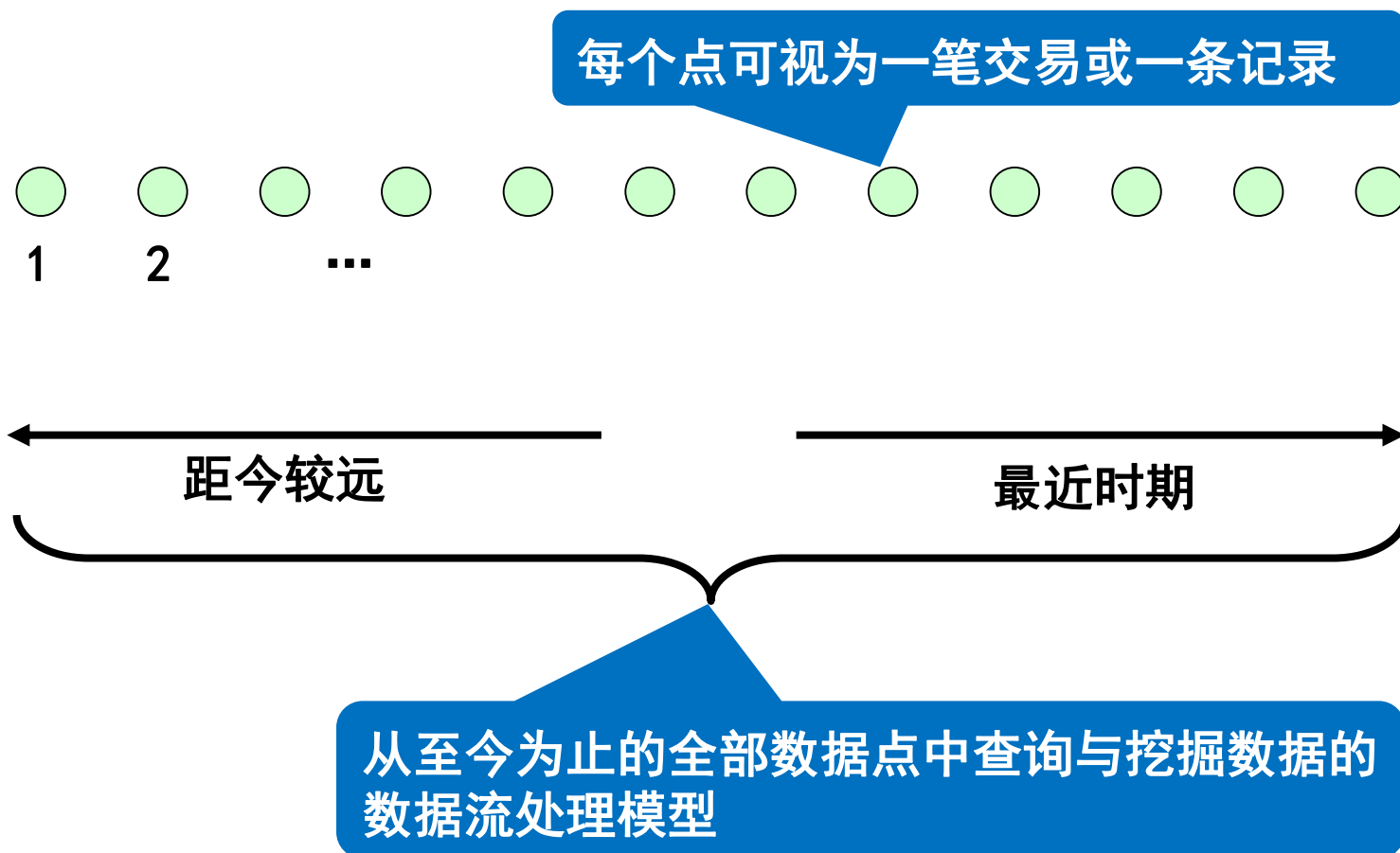
# 数据流概念

---



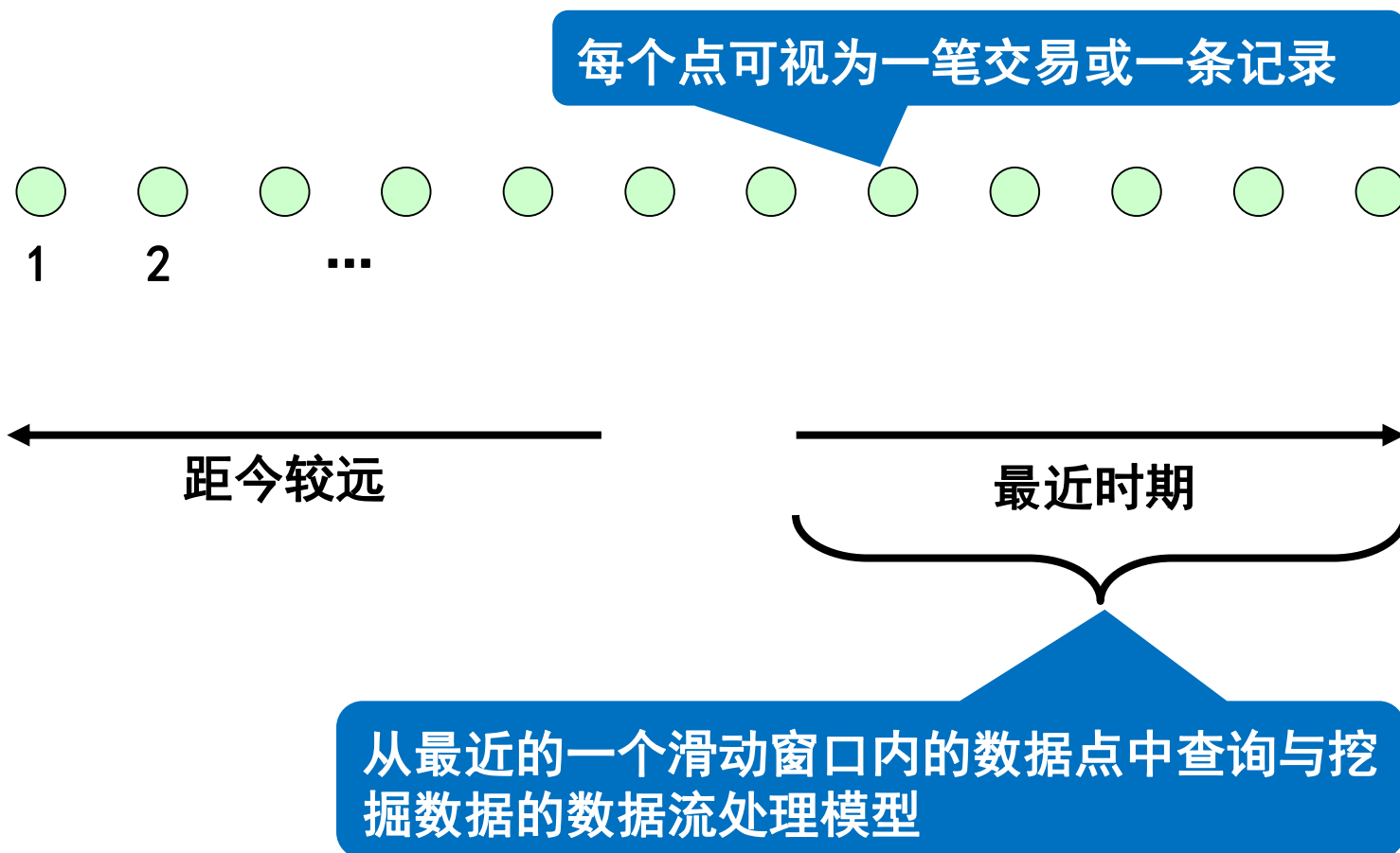
# 数据流模型

- 完整数据流模型 (Entire Data Stream Model)



# 数据流模型

- 滑动窗口数据流模型 (Sliding-Window Model)





# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# 数据流处理的基础性问题

---

**差异项种类查询问题：** 给定一个长度为 $N$ 的数据流，该问题旨在完整数据流模型下统计差异项（即不同项）种类的个数！

# 数据流处理的基础性问题

---

**差异项种类查询问题：** 给定一个长度为 $N$ 的数据流，该问题旨在完整数据流模型下统计差异项（即不同项）种类的个数！

**差异项频率查询问题：** 给定一个长度为 $N$ 的数据流，该问题完整数据流模型下统计每个差异项（即不同项）的出现频率！

# 数据流处理的基础性问题

**差异项种类查询问题：** 给定一个长度为 $N$ 的数据流，该问题旨在完整数据流模型下统计差异项（即不同项）种类的个数！

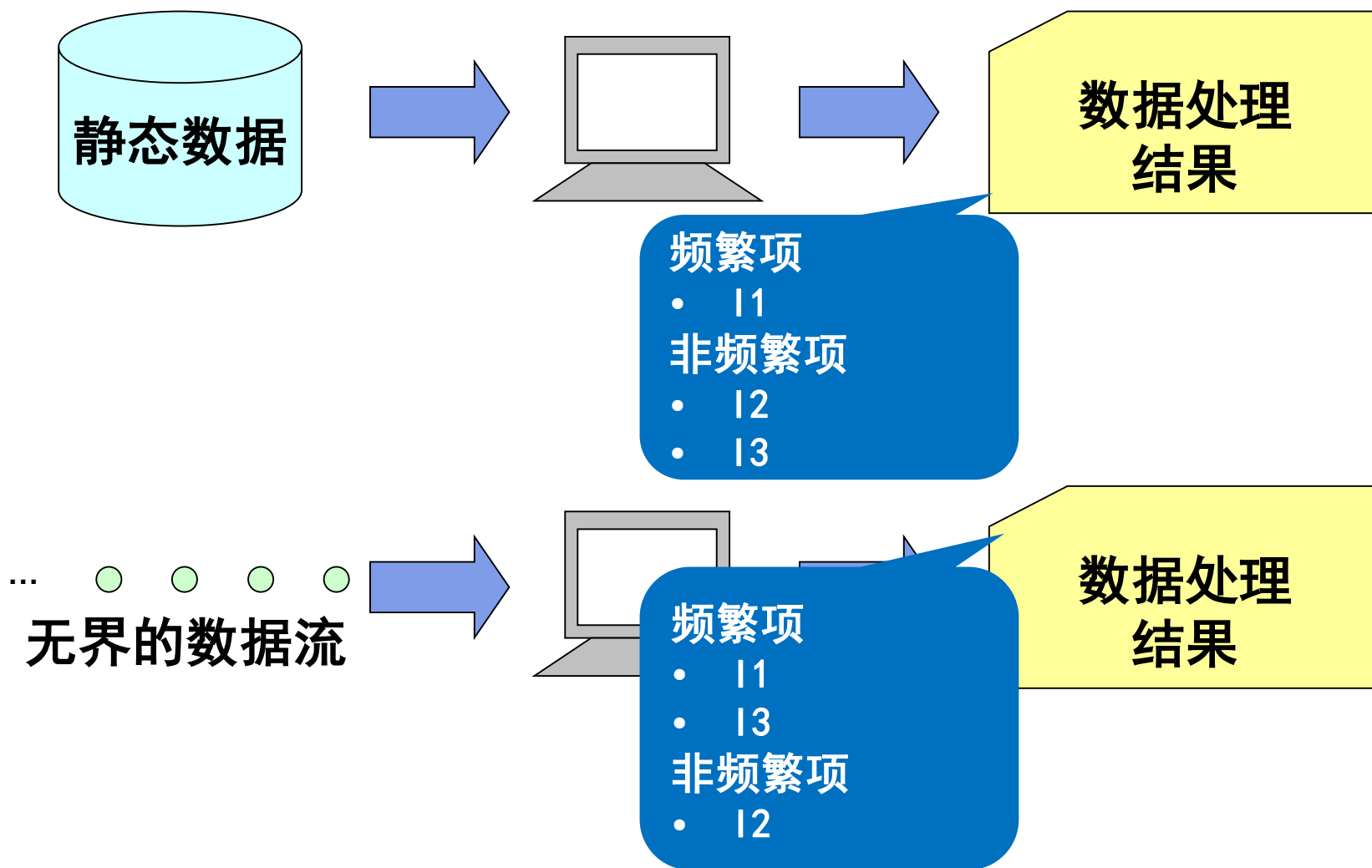
**频繁项查询问题：** 给定一个长度为 $N$ 的数据流和频繁阈值 $s$ ，该问题旨在完整数据流模型下发现全部频繁项，即每个项出现频率大于 $sN$ ！

**差异项频率查询问题：** 给定一个长度为 $N$ 的数据流，该问题完整数据流模型下统计每个差异项（即不同项）的出现频率！

# 数据流处理算法特性

	传统静态数据处理	动态数据流处理
数据类型	有限规模的静态数据	无限规模、动态且高速抵达的流状态数据
存储方式	硬盘/存储有限	存储空间不足
响应效率	非实时查询	实时查询
查询结果	精准结果	近似(或精准)结果

# 静态数据查询 vs. 数据流查询



# 假阳性 (False Positive)

- E. g.
  - 期望输出
    - 频繁项
      - $I_1$
    - 非频繁项
      - $I_2$
      - $I_3$
  - 算法输出
    - 频繁项
      - $I_1$
      - $I_3$
    - 非频繁项
      - $I_2$

假阳性

- 某项被分类为频繁项
- 但该项事实上是非频繁项

哪个项是假阳性?

$I_3$

还有么?      没了

假阳性数量 = 1

如果称算法不存在假阳性

所有非频繁项在算法输出中都被分类为非频繁项

# 假阴性 (False Negative)

- E. g.
  - 期望输出
    - 频繁项
      - $I_1$
    - 非频繁项
      - $I_2$
      - $I_3$
  - 算法输出
    - 频繁项
      - $I_1$
      - $I_3$
    - 非频繁项
      - $I_2$

假阴性

- 某项被分类为非频繁项
- 但该项事实上是频繁项

哪个项是假阴性?

$I_3$

还有么?      没了

假阴性数量 = 1

假阳性数量 = 0

如果称算法不存在假阴性

所有频繁项在算法输出中都被分类为频繁项



# 数据流查询的近似误差

所有项真实统计信息

- I<sub>1</sub>: 10
- I<sub>2</sub>: 8
- I<sub>3</sub>: 12

N: 流中项总数

$$N = 20$$

$$\varepsilon = 0.2$$

$$\varepsilon N = 4$$

静态数据

数据处理  
结果

项	真实出现次数	估计出现次数
I <sub>1</sub>	10	10
I <sub>2</sub>	8	4
I <sub>3</sub>	12	10

Diff. D

D ≤ εN ?

0

Yes

4

Yes

2

Yes

数据处理  
结果

...  
无界的数据流

所有项估计统计信息

- I<sub>1</sub>: 10
- I<sub>2</sub>: 4
- I<sub>3</sub>: 10

# 容错概要 ( $\epsilon$ -deficient synopsis)

- 数据流长度 $N$ : 为流的当前长度(或项出现的总次数)
- 近似误差 $\epsilon$ : 为输入参数(0和1之间的实数)
- 一个算法如果满足以下条件, 则符合 $\epsilon$ -deficient synopsis
  - 条件 1: 不存在假阴性.  

在算法的输出中, 所有真正频繁的项都被输出为频繁项
  - 条件 2: 估计频率和真实频率的差值不超过 $\epsilon N$ .
  - 条件 3: 所有真实频率少于 $(s-\epsilon)N$ 的项都会被分类为不频繁项.

**近似频繁项查询问题:** 给定一个长度为 $N$ 的数据流、频繁阈值 $s$ 和近似误差 $\epsilon$ , 该问题旨在完整数据流模型下发现全部的近似频繁项, 即每个项出现频率大于 $(s-\epsilon)N$ !

# 课程提纲

---

- 数据流简介
- 基础性问题
- **确定性算法**
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# 热身：过半数Majority算法 ( $>1/2$ )

---

- 问题：如何找出数据流S中出现次数超过 $1/2$ 的元素
- $|S|$ 很大, 无法记录每个元素的出现次数

# 热身：过半数Majority算法 ( $>1/2$ )

---

- 输入：
  - 数据流  $S = x_1, x_2, \dots, x_n, x_i \in U$
- 输出：流  $S$  中出现频率超过  $s=1/2$  的元素
- 算法：MAJORITY
  - 如果存在多数元素，则MAJORITY一定返回该元素

# Majority算法

---

- MAJORITY: Initialize(x)
  - $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$
- MAJORITY: Update(x)
  - If  $\text{value} = 0$  then
    - $\text{key} \leftarrow x, \text{value} = 1$
  - Else if  $\text{key} = x$  then
    - $\text{value} = \text{value} + 1$
  - Else
    - $\text{value} = \text{value} - 1$
- MAJORITY: Return()
  - If  $\text{value} > 0$  then return key
  - Else return  $\emptyset$

# Majority算法

- MAJORITY: Initialize( $x$ )

- $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$

- MAJORITY: Update( $x$ )

- If  $\text{value} = 0$  then

- $\text{key} \leftarrow x, \text{value} \leftarrow 1$

- Else if  $\text{key} = x$

- $\text{value} = \text{value} + 1$

- Else

- $\text{value} = \text{value} - 1$

- MAJORITY: Return()

- If  $\text{value} > 0$  then return  $\text{key}$

- Else return  $\emptyset$

初始化key为空, 计数器为0



# Majority算法

- MAJORITY: Initialize(x)
  - $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$
- MAJORITY: Update(x)
  - If  $\text{value} = 0$  then
    - $\text{key} \leftarrow x, \text{value} = 1$
  - Else if  $\text{key} = x$  then
    - $\text{value} = \text{value} + 1$
  - Else
    - $\text{value} = \text{value} - 1$
- MAJORITY: Return()
  - If  $\text{value} > 0$  then return  $\text{key}$
  - Else return  $\emptyset$

如果计数器为0, 记录当前元素  
并设计数器为1

# Majority算法

- MAJORITY: Initialize(x)
  - $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$
- MAJORITY: Update(x)
  - If  $\text{value} = 0$  then
    - $\text{key} \leftarrow x, \text{value} = 1$
  - Else if  $\text{key} = x$  then
    - $\text{value} = \text{value} + 1$
  - Else
    - $\text{value} = \text{value} - 1$
- MAJORITY: Return
  - If  $\text{value} > 0$  then return  $\text{key}$
  - Else return  $\emptyset$

如果计数器不为0, 且x与key记录元素相同, 增加计数器

# Majority算法

- MAJORITY: Initialize(x)
  - $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$
- MAJORITY: Update(x)
  - If  $\text{value} = 0$  then
    - $\text{key} \leftarrow x, \text{value} \leftarrow 1$
  - Else if  $\text{key} = x$  then
    - $\text{value} = \text{value} + 1$
  - Else
    - $\text{value} = \text{value} - 1$
- MAJORITY: Return()
  - If  $\text{value} > 0$  then return key
  - Else return  $\emptyset$

如果计数器不为0, 且x与key记录元素不同, 减小计数器

# Majority算法

- MAJORITY: Initialize(x)

- $\text{key} \leftarrow \emptyset, \text{value} \leftarrow 0$

- MAJORITY: Update(x)

- If  $\text{value} = 0$  then

- $\text{key} \leftarrow x, \text{value} = 1$

- Else if  $\text{key} = x$  then

- $\text{value} = \text{value} + 1$

- Else

- $\text{value} = \text{value} - 1$

- MAJORITY: Return()

- If  $\text{value} > 0$  then return key

- Else return  $\emptyset$

如果计数器不为0, 则key中记录元素为出现次数超过1/2元素

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key =  $\emptyset$ , value = 0

两种元素

- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1



# Majority算法实例

---

- a, b, a, b, a, b, a, a, b
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 2
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1



# Majority算法实例

---

- a, b, a, b, a, b, a, b, a, a, b
- key = a, value = 1
- 存在多数元素a
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

- a, b, a, c, a, e, a, d, a, f, d
- key =  $\emptyset$ , value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

多种元素且无任意元素频率过半

MAJORITY不适用

# Majority算法实例

---

- **a**, b, a, c, a, e, a, d, a, f, d
- key = **a**, value = **1**
- If value = 0 then
  - **key**  $\leftarrow$  **x**, **value** = **1**
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1



# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = a, value = 0
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, f, d
- key = d, value = 1
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# Majority算法实例

---

- a, b, a, c, a, e, a, d, a, a, d
- key = d, value = 1
- d并非多数元素
- 需要检查. 但对于流数据, 无法检查
- If value = 0 then
  - key  $\leftarrow$  x, value = 1
- Else if key = x then
  - value = value + 1
- Else
  - value = value - 1

# 过半数Majority算法扩展 ( $>1/k$ )

---

- 问题：如何找出数据流S中出现频率超过 $s=1/k$ 的元素
  - $|S|$ 很大, 无法记录每个元素的出现次数

# 过半数Majority算法扩展( $>1/k$ )

- 问题：如何找出数据流 $S$ 中出现频率超过 $s=1/k$ 的元素
  - $|S|$ 很大, 无法记录每个元素的出现次数
- 输入：
  - 数据流 $S=x_1, x_2, \dots, x_n, x_i \in U$
- 输出： $S$ 中所有出现频率超过 $s=1/k$ 的元素
- 算法：
  - MG算法
  - Space-Saving算法

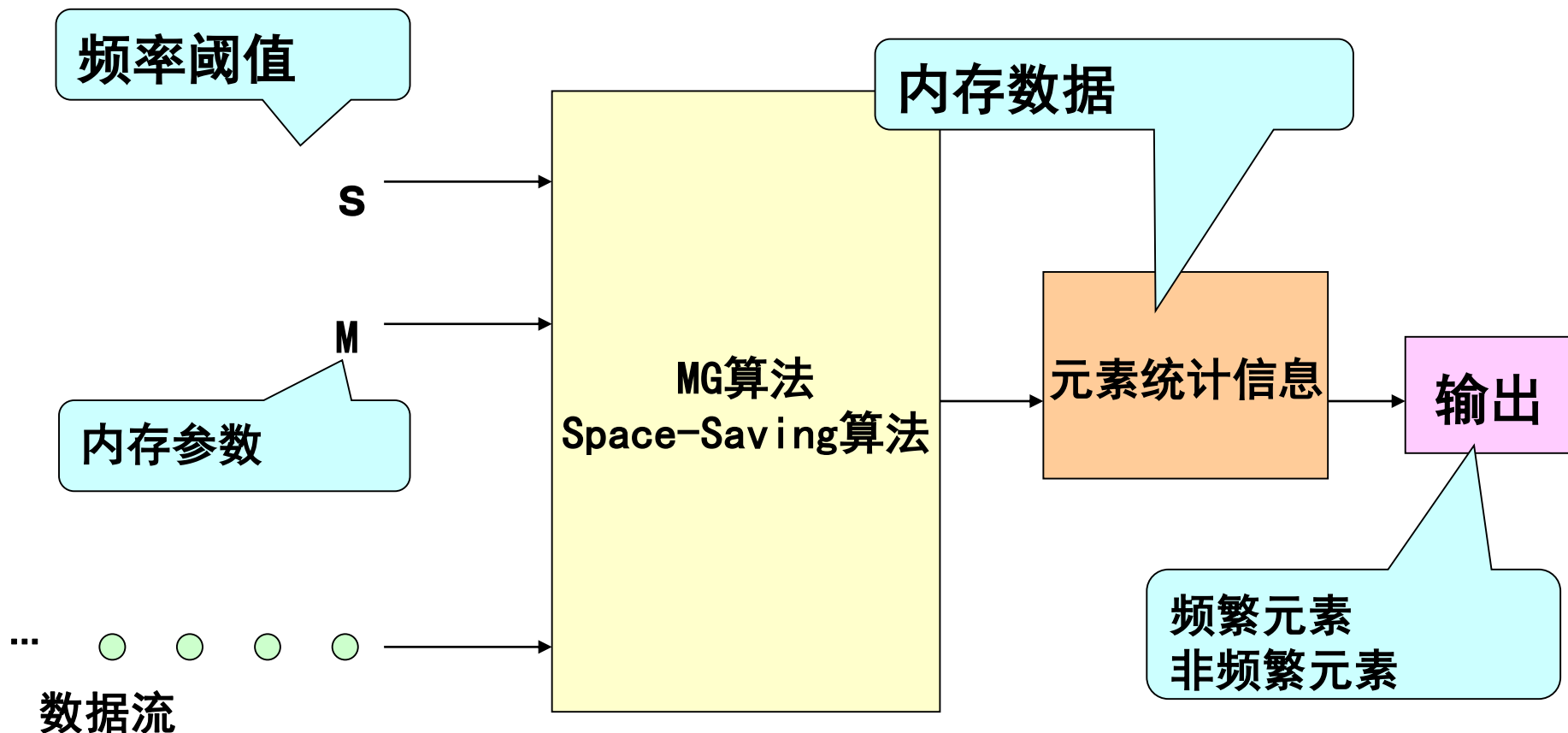


# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# MG算法与Space-Saving算法框架<sup>58</sup>



# Misra-Gr i e s (MG) 算法

---

- MG:Update (x, w)
  - If  $x \in T$ 
    - $w_x \leftarrow w_x + w$
  - Else
    - $T \leftarrow T \cup \{x\}$
    - $w_x \leftarrow w$
    - If  $|T| > k$  then
      - $m = \min(\{w_x : x \in T\})$
      - For all  $j \in T$  do
        - $w_j \leftarrow w_j - m$
        - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# Misra-Gries (MG) 算法

- MG: Update ( $x, w$ )

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$

- Else

- $T \leftarrow T \cup \{x\}$

- $w_x \leftarrow w$

- If  $|T| > k$  then

- $m = \min(\{w_x : x \in T\})$

- For all  $j \in T$  do

- $w_j \leftarrow w_j - m$

- If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

对于已经在T中的元素, 直接更新其权重

# Misra-Gries (MG) 算法

- MG: Update ( $x, w$ )

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $T \leftarrow T \cup \{x\}$

- $w_x \leftarrow w$

- If  $|T| > k$  then

- $m = \min(\{w_x : x \in T\})$

- For all  $j \in T$  do

- $w_j \leftarrow w_j - m$

- If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

对于不在T中的元素, 将其加入T并初始化其权重

# Misra-Gries (MG) 算法

- MG: Update ( $x, w$ )

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $T \leftarrow T \cup \{x\}$

- $w_x \leftarrow w$

- If  $|T| > k$  then

- $m = \min(\{w_x : x \in T\})$

- For all  $j \in T$  do

- $w_j \leftarrow w_j - m$

- If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

如果  $|T| > k$ , 将  $T$  中每个元素的权重减去  $T$  中元素的最小权重, 并删除权重为 0 的元素

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(-, 0), (-, 0), (-, 0)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- $a, b, a, c, d, e, a, d, f, a, d$
- $T = \{ (a, 1), (-, 0), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$



# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 1), (b, 1), (-, 0)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 2), (b, 1), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (b, 1), (c, 1)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (b, 1), (c, 1), (d, 1)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{ (a, 1), (b, 0), (c, 0), (d, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $T \leftarrow T \cup \{x\}$

- $w_x \leftarrow w$

- If  $|T| > k$  then

- $m = \min(\{w_x : x \in T\})$

- For all  $j \in T$  do

- $w_j \leftarrow w_j - m$

- If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 1), (-, 0), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 1), (e, 1), (-, 0)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 2), (e, 1), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$



# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (e, 1), (d, 1)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (e, 1), (d, 1), (f, 1)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{ (a, 1), (e, 0), (d, 0), (f, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $T \leftarrow T \cup \{x\}$

- $w_x \leftarrow w$

- If  $|T| > k$  then

- $m = \min(\{w_x : x \in T\})$

- For all  $j \in T$  do

- $w_j \leftarrow w_j - m$

- If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 1), (-, 0), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 2), (-, 0), (-, 0) \}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (d, 1), (-, 0)\}$

内存参数  $|T|=3$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $T \leftarrow T \cup \{x\}$
  - $w_x \leftarrow w$
  - If  $|T| > k$  then
    - $m = \min(\{w_x : x \in T\})$
    - For all  $j \in T$  do
      - $w_j \leftarrow w_j - m$
      - If  $w_j = 0$  then  $T \leftarrow T \setminus \{j\}$

# MG算法实例

---

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (d, 1), (-, 0)\}$
- T中每个元素的估计误差不超过 $W/|T|=11/3$ 
  - W为流中所有元素权重加和
- 使用T找出流中出现频率超过1/2的元素
- 流中出现频率超过1/2的元素为a
  - $a: 2 + 11/3 > 11/2$

# Space-Saving算法

---

- Space-Saving: Update ( $x, w$ )
  - If  $x \in T$ 
    - $w_x \leftarrow w_x + w$
  - Else
    - $y \leftarrow \operatorname{argmin}_{y \in T} w_y$
    - $w_x \leftarrow w_y + w$
    - $T \leftarrow T \cup \{x\} \setminus \{y\}$



# Space-Saving算法

- Space-Saving: Update ( $x, w$ )

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_j$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

对于已经在T中的元素, 直接  
更新其权重

# Space-Saving算法

- Space-Saving: Update ( $x, w$ )

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_y$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

对于不在T中的元素, 将其权重加上T中元素的最小权重, 并取代T中权重最小的元素

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(-, 0), (-, 0), (-, 0), (-, 0)\}$

内存参数  $|T|=4$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $y \leftarrow \operatorname{argmin}_{y \in T} w_y$
  - $w_x \leftarrow w_y + w$
  - $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 1), (-, 0), (-, 0), (-, 0) \}$

内存参数  $|T|=4$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $y \leftarrow \operatorname{argmin}_{y \in T} w_y$  ( $y = -$ )
  - $w_x \leftarrow w_y + w$  ( $w_y = 0$ )
  - $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{(a, 1), (b, 1), (-, 0), (-, 0)\}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_j \quad (y = -)$

- $w_x \leftarrow w_y + w \quad (w_y = 0)$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{ (a, 2), (b, 1), (-, 0), (-, 0) \}$

内存参数  $|T|=4$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $y \leftarrow \operatorname{argmin}_{y \in T} w_y$
  - $w_x \leftarrow w_y + w$
  - $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (b, 1), (c, 1), (-, 0)\}$

内存参数  $|T|=4$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $y \leftarrow \operatorname{argmin}_{y \in T} w_y$  ( $y = -$ )
  - $w_x \leftarrow w_y + w$  ( $w_y = 0$ )
  - $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{(a, 2), (b, 1), (c, 1), (d, 1)\}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_j$  ( $y = -$ )

- $w_x \leftarrow w_y + w$  ( $w_y = 0$ )

- $T \leftarrow T \cup \{x\} \setminus \{y\}$



# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 2), (e, 2), (c, 1), (d, 1)\}$

内存参数  $|T|=4$

- If  $x \in T$ 
  - $w_x \leftarrow w_x + w$
- Else
  - $y \leftarrow \operatorname{argmin}_{y \in T} w_j \quad (w_y = 1)$
  - $w_x \leftarrow w_y + w \quad (w_y = 1)$
  - $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{ (a, 3), (e, 2), (c, 1), (d, 1) \}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_y$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{(a, 3), (e, 2), (c, 1), (d, 2)\}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_y$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{(a, 3), (e, 2), \text{~~(e, 1)~~, (d, 2)}\}$

- $T = \{(a, 3), (e, 2), \text{(f, 2)}, (d, 2)\}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_j \quad (y = c)$

- $w_x \leftarrow w_y + w \quad (w_y = 1)$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{ (a, 4), (e, 2), (f, 2), (d, 2) \}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_y$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

- a, b, a, c, d, e, a, d, f, a, d

- $T = \{(a, 4), (e, 2), (f, 2), (d, 3)\}$

内存参数  $|T|=4$

- If  $x \in T$

- $w_x \leftarrow w_x + w$

- Else

- $y \leftarrow \operatorname{argmin}_{y \in T} w_y$

- $w_x \leftarrow w_y + w$

- $T \leftarrow T \cup \{x\} \setminus \{y\}$

# Space-Saving算法实例

---

- a, b, a, c, d, e, a, d, f, a, d
- $T = \{(a, 4), (e, 2), (f, 2), (d, 3)\}$
- T中每个元素出现次数估计误差不超过 $W/|T|=11/4$ 
  - W为流中所有元素权重加和
- 使用T找出流中出现频率超过1/3的元素
- 流中出现频率超过1/3的元素为a, e, f, d
  - a:  $4 + 11/4 > 11/3$
  - e:  $2 + 11/4 > 11/3$ ,
  - f:  $2 + 11/3 > 11/3$
  - d:  $3 + 11/4 > 11/3$

# MG算法与Space-Saving算法关系

---

- $T_{MG} = \{(a, 2), (d, 1), (-, 0)\}$
- $T_{Space-Saving} = \{(a, 4), (e, 2), (f, 2), (d, 3)\}$
- $T_{MG} = \{(a, 2), (d, 1), (-, 0)\}$
- $T_{Space-Saving} = \{(a, 2), (e, 0), (f, 0), (d, 1)\}$
- $T_{MG} = \{(a, 2), (d, 1)\}$
- $T_{Space-Saving} = \{(a, 2), (d, 1)\}$
- MG算法  $|T|=k$  与 Space-Saving算法  $|T|=k+1$  等价



# MG算法与Space-Saving算法分析

---

- 内存参数为 $M$ , 即 $|T|=M$ 
  - 对每个元素频率的估计误差最大为 $1/M$
- 使用 $T$ 寻找流中出现频率超过 $s$ 的元素
  - 当满足 $E \leq 1/M \leq \epsilon$  时
  - 流中所有出现频率超过 $s$ 的元素都不会漏掉
    - 即不存在假阴性 (No false negative)

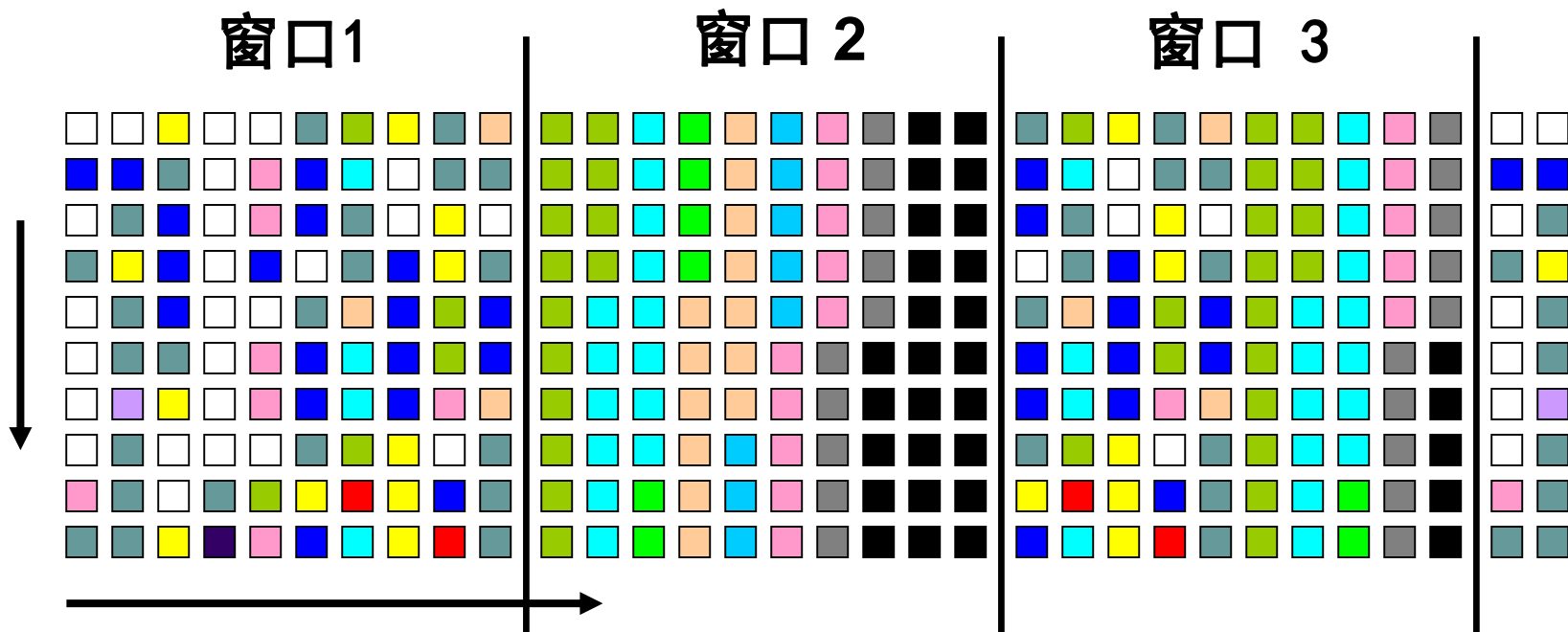
# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# Lossy Counting 算法核心思想

- 第一步：基于窗口的数据流切分

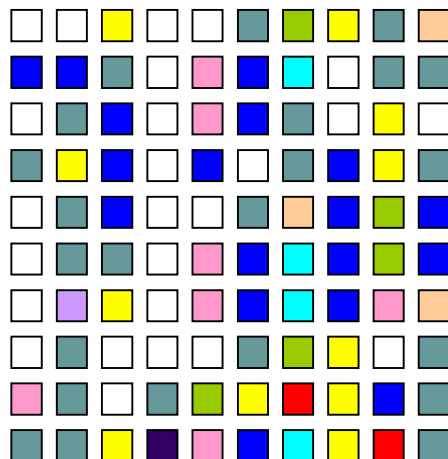


# Lossy Counting 算法核心思想

- 第二步：每到窗口边界，每项计数减一

空集

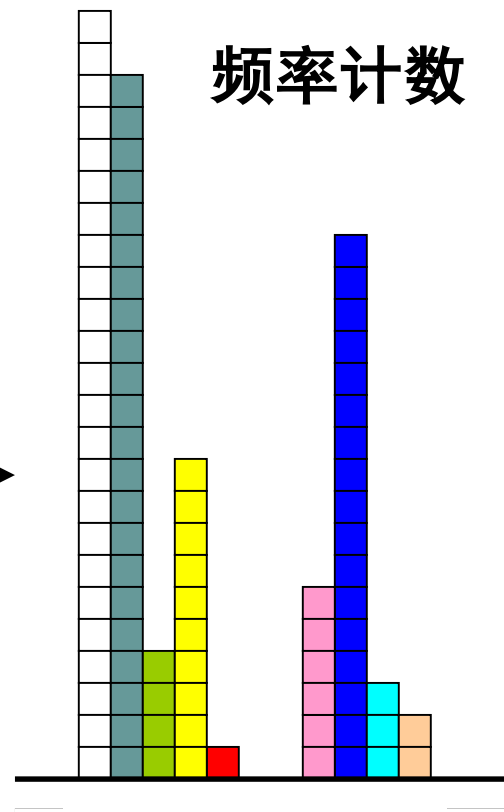
+



第一个窗口

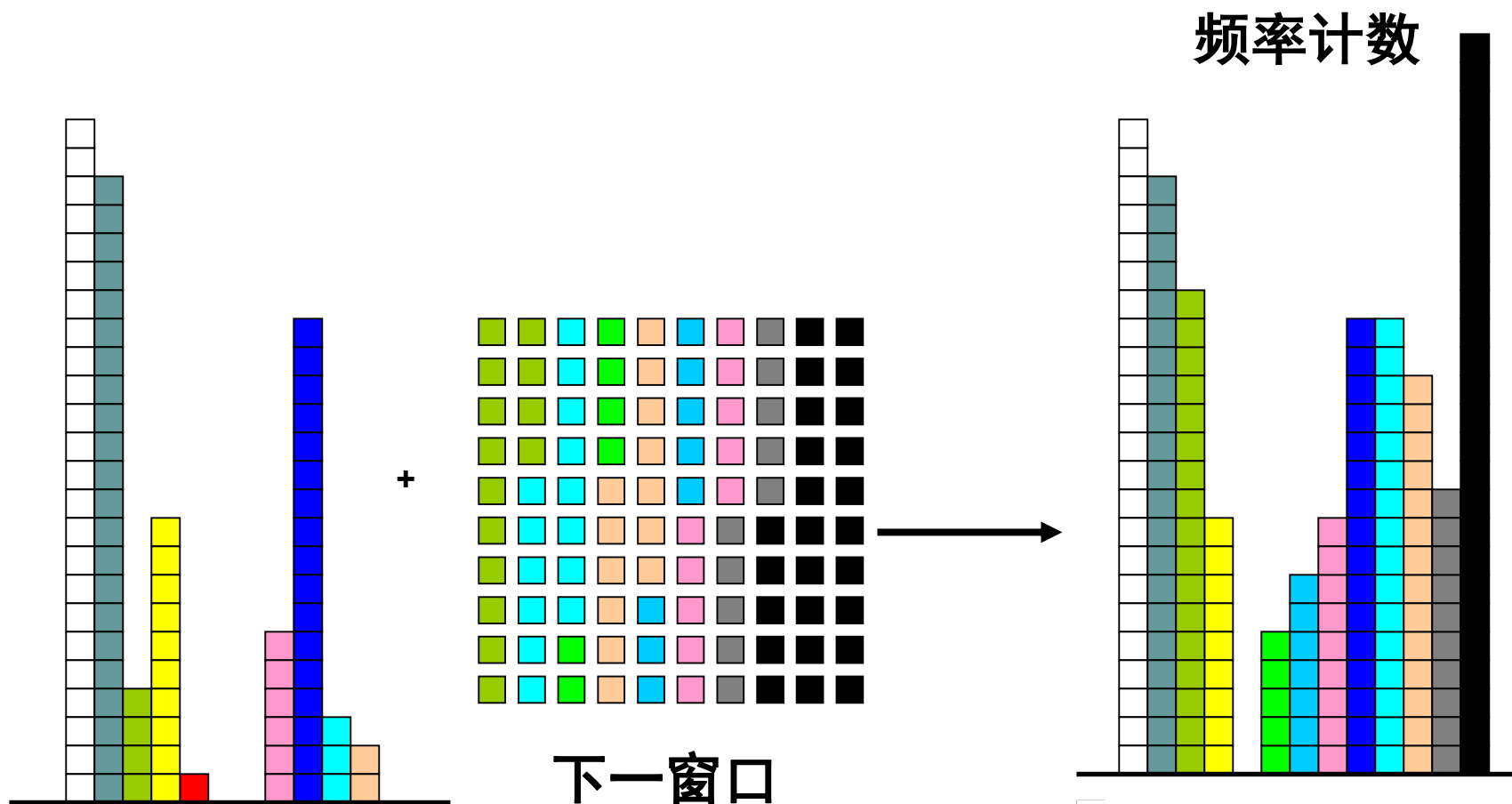


频率计数



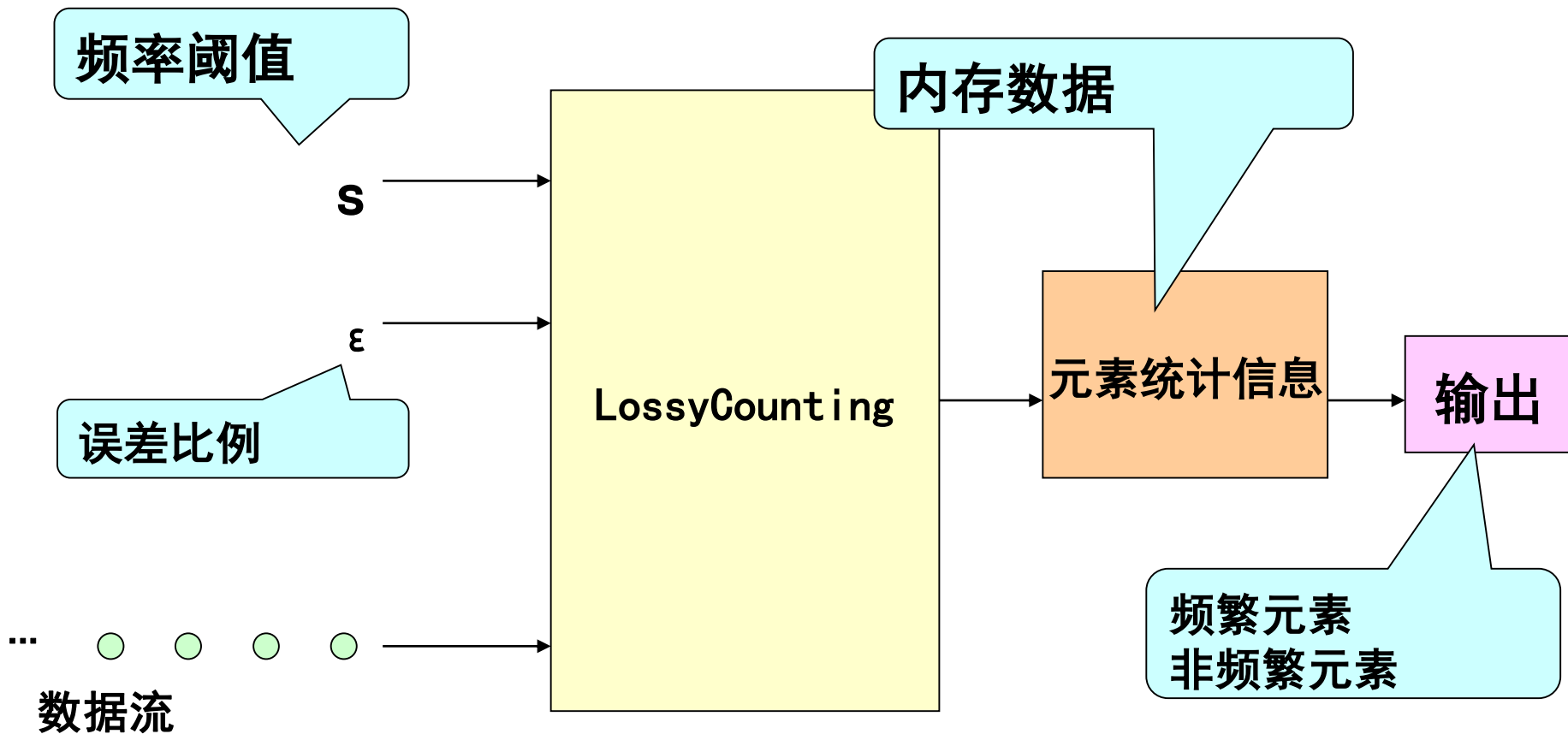
# Lossy Counting 算法核心思想

- 第一步：每到窗口边界，每项计数减一



# Lossy Counting算法框架

102



# Lossy Counting算法

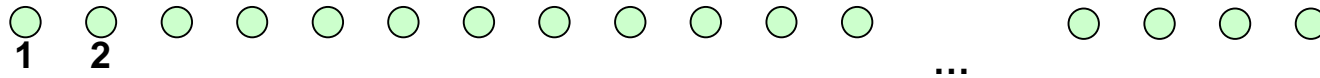
**问题定义：**统计流数据中出现频率大于给定值 $s$ 的元素。

- **算法框架**

- 根据参数设置流数据窗口大小。
- 对于顺序到来的一个窗口内的数据更新内存中数据结构。
- 每执行完一个窗口内的数据，更新内存中的数据结构（以减少空间），更新参数。

# Lossy Counting算法

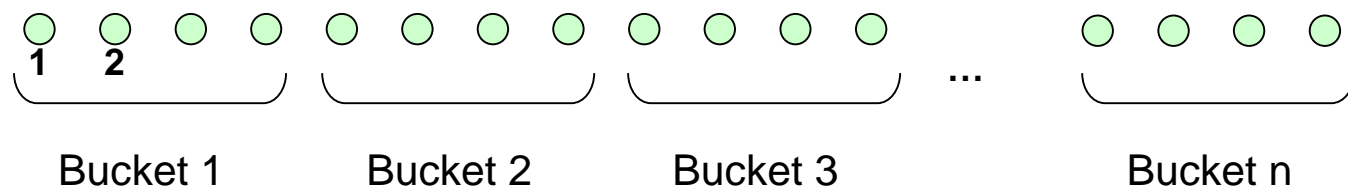
问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。





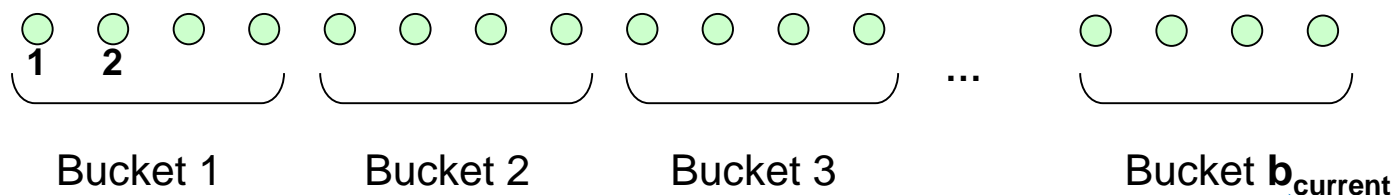
# Lossy Counting算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。



# Lossy Counting算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。



Width  $w = \left\lceil \frac{1}{\epsilon} \right\rceil$

$b_{\text{current}} = \left\lceil \frac{N}{w} \right\rceil$

# Lossy Counting算法

**问题定义：**统计流数据中出现频率大于给定值 $s$ 的元素。

- **算法介绍**

- 内存中保存部分元素、其出现次数以及该元素的最大误差 $(e, f, \Delta)$ ，对于数据流顺序到来的元素，若数据出现的元素已存在，将次数加一；若不存在，则将 $(e, 1, \Delta)$ 放入内存。
  - 窗口大小不变，动态更新内存存储信息。
  - 数据结构：字典
  - 返回数据结构中频度大于 $sN - \epsilon N$ 的元素。
  - 确定性算法
  - 至多使用 $\frac{1}{\epsilon} \log(\epsilon N)$ 个key-value对

- **相关参数**

- 窗口大小参数 $w = \lceil 1/\epsilon \rceil$

# Lossy Counting算法

问题定义：统计流数据中出现频率大于给定值  $s$  的元素。

1. **D: Empty set**
  - Will contain  $(e, f, \Delta)$
2. **When data  $e$  arrives,**
  - If  $e$  exists in  $D$ ,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in  $D$ ,
3.  **$s$ , Remove some entries in  $D$  whenever  $N \equiv 0 \pmod w$**   
(i.e., whenever it reaches the bucket boundary)  
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  
 $f + \Delta \leq b_{\text{current}}$
4. **[Output]** Get a list of items where  
 $f + \epsilon N \geq sN$

元素

元素自被添加之后出现的次数

$f$  的最大可能误差

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小 $w=\lceil 1/\epsilon \rceil=5$ ,窗口个数 $b=1$

a

S

元素	频度	误差

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=1$

a b

s

元素	频度	误差
a	1	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. s, Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=1$

a b a

S

元素	频度	误差
a	2	0
b	1	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小 $w=\lceil 1/\epsilon \rceil=5$ ,窗口个数 $b=1$

a b a c

S

元素	频度	误差
a	2	0
b	1	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$



# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=1$

a b a c d

S

元素	频度	误差
a	2	0
b	1	0
c	1	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=1$

a b a c d

S

元素	频度	误差
a	2	0
b	1	0
c	1	0
d	1	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. s, Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=1$

a b a c d

S

元素	频度	误差
a	2	0
b	1	0
c	1	0
d	1	0

$2+0>1$
$1+0\leq 1$ , drop
$1+0\leq 1$ , drop
$1+0\leq 1$ , drop

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e

S

元素	频度	误差
a	2	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a

S

元素	频度	误差
a	2	0
e	1	1

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a d

S

元素	频度	误差
a	3	0
e	1	1

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a d f

S

元素	频度	误差
a	3	0
e	1	1
d	1	1

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a d f a

S

元素	频度	误差
a	3	0
e	1	1
d	1	1
f	1	1

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$



# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a d f a

S

元素	频度	误差
a	4	0
e	1	1
d	1	1
f	1	1

$3+0>2$
$1+1\leq 2$ , drop
$1+1>2$ , drop
$1+1>2$ , drop

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data e arrives,
  - If e exists in D,
    - Increment f in  $(e, f, \Delta)$
  - If e does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=2$

a b a c d | e a d f a | d

S

元素	频度	误差
a	4	0

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3.  $s$ , Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=3$

a b a c d | e a d f a | d

S

元素	频度	误差
a	4	0
d	1	2

1. D: Empty set
  - Will contain  $(e, f, \Delta)$
2. When data  $e$  arrives,
  - If  $e$  exists in D,
    - Increment  $f$  in  $(e, f, \Delta)$
  - If  $e$  does not exist in D,
3. Remove some entries in D whenever  $N \equiv 0 \pmod w$   
The rule of deletion is:  
 $(e, f, \Delta)$  is deleted if  $f + \Delta \leq b_{\text{current}}$
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Lossy Counting算法

$\epsilon=0.2$ , 窗口大小  $w=\lceil 1/\epsilon \rceil=5$ , 窗口个数  $b=3$

a b a c d | e a d f a | d

s

元素	频度	误差
a	4	0
d	1	2

[Output] Get a list of items where  
 $f + \epsilon N \geq sN$

$s=0.3$ , 取使得  $f \geq N(s - \epsilon)$  的元素, 为 **a**

# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# Sticky Sampling算法

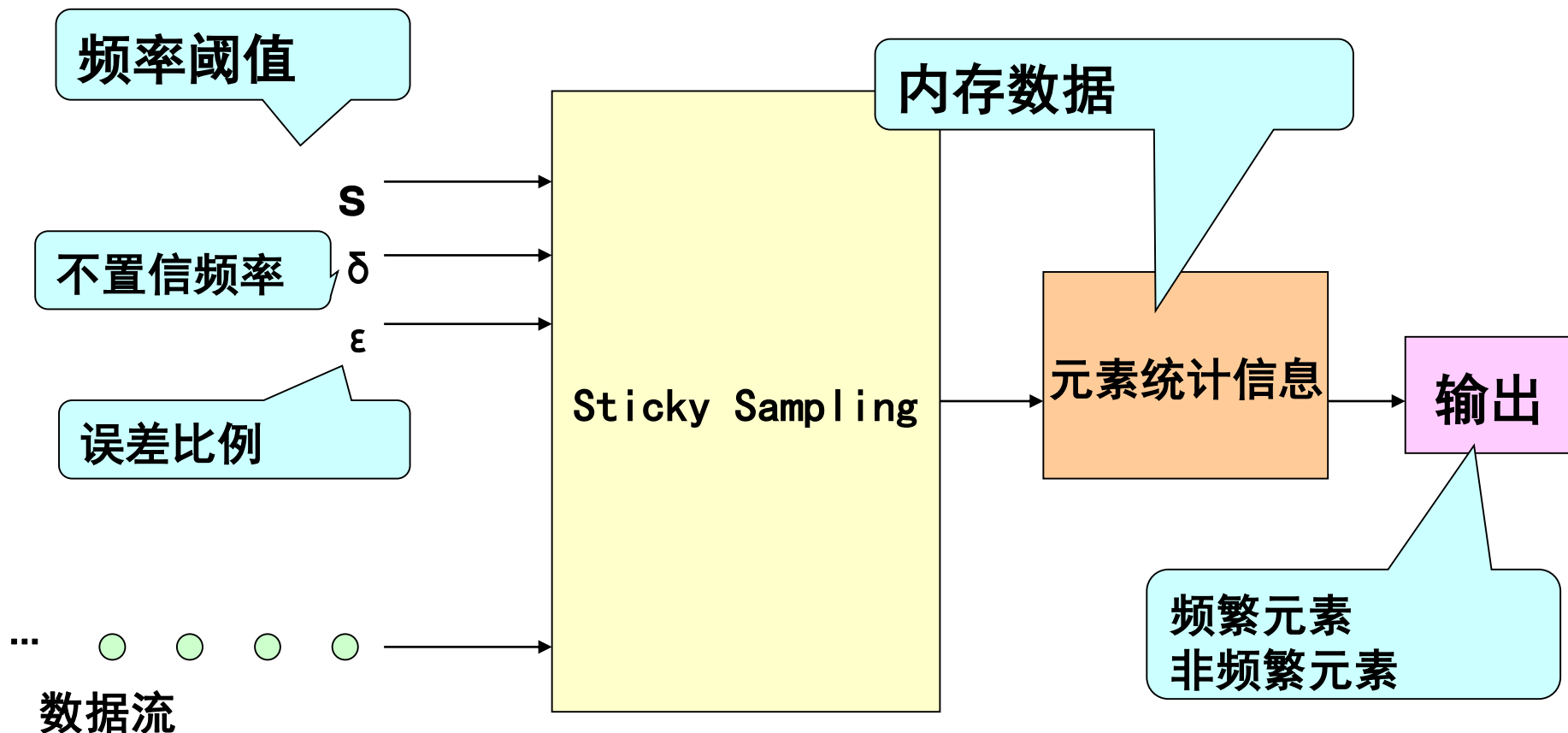
问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。

- 算法框架

- 基本数据流模型
- 根据参数设置流数据窗口大小。
- 对于顺序到来的一个窗口内的数据更新内存中数据结构。
- 每执行完一个窗口内的数据，更新内存中的数据结构（以减少空间），更新参数。

# Sticky Sampling算法框架

127



# Sticky Sampling算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。

## ● 算法介绍

- 内存中保存部分元素及其出现次数  $(e, f)$ ，对于数据流顺序到来的元素，数据出现的元素已存在，将次数加一；若不存在，做一次抽样决定是否将其放入内存。
  - 每处理完一个窗口的数据，抽样频率减半，窗口大小加倍，更新数据结构。
  - 返回数据结构中频度大于  $sN - \epsilon N$  的元素。
  - 至多使用  $\frac{1}{\epsilon} \log(s^{-1} \delta^{-1})$  个 key-value 对



# Sticky Sampling算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。

- 相关参数

- 抽样率频率参数 $r$ ，值随数据流的增加而变化。
- 置信频率参数 $\delta$
- 样窗口参数  $t = \frac{1}{\epsilon} \log(s^{-1} \delta^{-1})$ ，窗口大小依次为  $2t, 2t, 4t, 8t, \dots$

<b>2t</b>	<b>1</b>
<b>2t</b>	<b>2</b>
<b>4t</b>	<b>4</b>
<b>...</b>	<b>...</b>

# Sticky Sampling算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。

- 相关参数

- 抽样率频率参数 $r$ ，值随数据流的变化而变化。

- 置信频率参数 $\delta$

- 样窗口参数  $t = \frac{1}{\varepsilon} \log(s^{-1} \delta^{-1})$ ，窗口大小依次为  $2t, 2t, 4t, 8t, \dots$

e.g.

$s = 0.02$

$\varepsilon = 0.01$

$\delta = 0.1$

$t = 622$

1~1244

$2t$

1

1245~2488

$2t$

2

2489~4976

$4t$


4

...

...

# Sticky Sampling算法

问题定义：统计流数据中出现频率大于给定值 $s$ 的元素。

1.  $S$ : empty list  
→ will contain  $(e, f)$   

2. When data  $e$  arrives,
  - if  $e$  exists in  $S$ , increment  $f$  in  $(e, f)$
  - if  $e$  does not exist in  $S$ , add entry  $(e, 1)$  with prob.  $1/r$  (where  $r$ : sampling rate)
3. When  $r$  changes,
  - For each entry  $(e, f)$ ,
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement  $f$  in  $(e, f)$
4. [Output] Get a list of items where  
 $f + \epsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第1个窗口  $r=1$ , 窗口大小2, 抽样频率  $1/r=1$

a

0.3 < 1/r, insert



s

元素	频度

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob. 1/r
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第1个窗口  $r=1$ , 窗口大小2, 抽样频率  $1/r=1$

a b

0.6 < 1/r, insert



s

元素	频度
a	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob. 1/r
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第1个窗口  $r=1$ , 窗口大小2, 抽样频率  $1/r=1$

a b

s

元素	频度
a	1
b	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第1个窗口  $r=1$ , 窗口大小2, 抽样频率  $1/r=1$

a b

s

元素	频度	
a	1	$0 < 1$
b	1	$0 < 1$



1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\epsilon=0.5$

第2个窗口  $r=2$ , 窗口大小2, 抽样频率  $1/r=1/2$

a b | a

不需要sample

s

元素	频度
a	1
b	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$



# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第2个窗口  $r=2$ , 窗口大小2, 抽样频率  $1/r=1/2$

a b | a c

0.4 < 1/r, insert



s

元素	频度
a	2
b	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob. 1/r
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\epsilon=0.5$

第2个窗口  $r=2$ , 窗口大小2, 抽样频率  $1/r=1/2$

a b | a c

s

元素	频度
a	2
b	1
c	1

1 < 2

1 ≥ 1, drop

0 < 1



1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第3个窗口  $r=4$ , 窗口大小4, 抽样频率  $1/r=1/4$

a b | a c | d

0.1 < 1/r, insert



s

元素	频度
a	1
c	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob. 1/r
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\epsilon=0.5$

第3个窗口  $r=4$ , 窗口大小4, 抽样频率  $1/r=1/4$

a b | a c | d e

0.4 > 1/r, drop



s

元素	频度
a	1
c	1
d	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\epsilon=0.5$

第3个窗口  $r=4$ , 窗口大小4, 抽样频率  $1/r=1/4$

a b | a c | d e a

不需要sample

s

元素	频度
a	1
c	1
d	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\epsilon=0.5$

第3个窗口  $r=4$ , 窗口大小4, 抽样频率  $1/r=1/4$

a b | a c | d e a d

不需要sample

s

元素	频度
a	2
c	1
d	1

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \epsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第3个窗口  $r=4$ , 窗口大小4, 抽样频率  $1/r=1/4$

a b | a c | d e a d

s

元素	频度	
a	2	$0 < 2$
c	1	$1 \geq 1$ , drop
d	2	$0 < 2$



1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第4个窗口  $r=8$ , 窗口大小8, 抽样频率  $1/r=1/8$

a b | a c | d e a d | f

$0.2 > 1/r$ , drop

s

元素	频度
a	2
c	1
d	2

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$



# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第4个窗口  $r=8$ , 窗口大小8, 抽样频率  $1/r=1/8$

a b | a c | d e a d | f a

不需要sample

s

元素	频度
a	2
c	1
d	2

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第4个窗口  $r=8$ , 窗口大小8, 抽样频率  $1/r=1/8$

a b | a c | d e a d | f a d

不需要sample

s

元素	频度
a	3
c	1
d	2

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第4个窗口  $r=8$ , 窗口大小8, 抽样频率  $1/r=1/8$

a b | a c | d e a d | f a d

s

元素	频度
a	3
c	1
d	3

1. S: empty list  
→ will contain (e, f)
2. When data e arrives,
  - if e exists in S, increment f in (e, f)
  - if e does not exist in S, add entry (e, 1) with prob.  $1/r$
3. When r changes,
  - For each entry (e, f),
    - Repeatedly toss a coin with  $P(\text{head}) = 1/r$  until the outcome of the coin toss is head
    - If the outcome of the toss is tail,
      - Decrement f in (e, f)
4. [Output] Get a list of items where  $f + \varepsilon N \geq sN$

# Sticky Sampling算法

抽样窗口参数  $t=1$ ,  $\varepsilon=0.5$

第4个窗口  $r=8$ , 窗口大小8, 抽样频率  $1/r=1/8$

a b | a c | d e a d | f a d

[Output] Get a list of items where  
 $f + \varepsilon N \geq sN$

$s=0.6$ , 取使得  $f \geq (s - \varepsilon)N$  的元素, 结果为 a,d

s

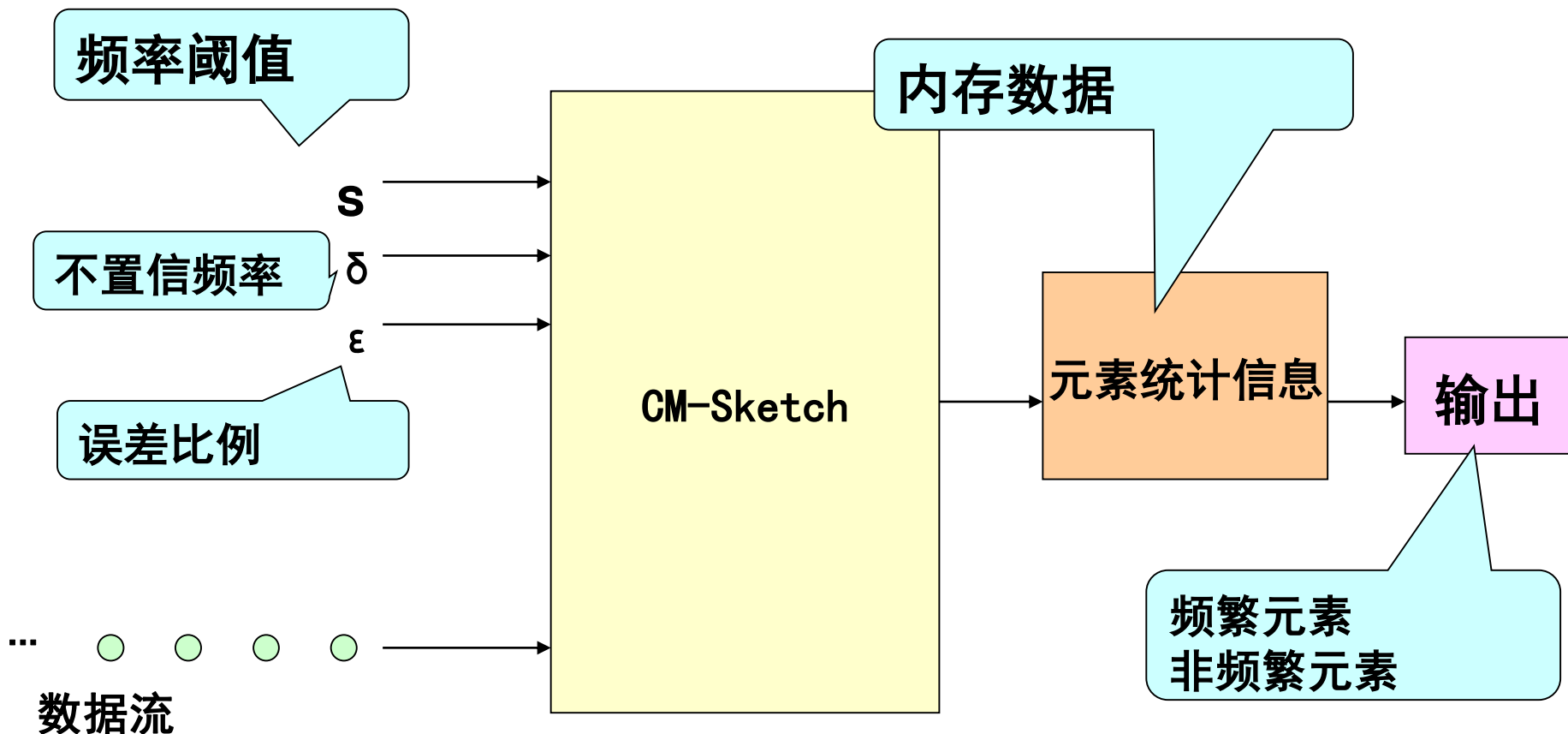
元素	频度
a	3
c	1
d	3

# 课程提纲

---

- 数据流简介
- 基础性问题
- 确定性算法
  - Majority算法
  - Misra-Gries (MG) 算法
  - Space-Saving算法
  - Lossy Counting算法
- 随机性算法
  - Sticky Sampling算法
  - Count-Min (CM) Sketch算法

# CM Sketch算法框架



# CM Sketch

问题定义：估计流数据中元素出现的总次数。

- 算法设计

- 使用Hash函数将元素映射到长度更小的数组中。
- 使用多个独立的Hash函数减小Hash碰撞导致的误差。
- 多个映射结果中最小的一个作为估计结果 (Count Min).
- 要求哈希函数  $h: \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, w\}$  碰撞频率不超过  $1/m$ .

$$\Pr\{x_i, x_j \in \{1, 2, \dots, n\} \wedge x_i \neq x_j \wedge P(x_i) = P(x_j)\} < 1/m$$

- 数据结构：二维数组

$m$ 为元素种类

- 元素的估计值  $Est(e) = \min_{1 \leq i \leq d} \{CM[i, h_i(e)]\}$

- 相关参数

- 数组深度  $w = \lceil e/\epsilon \rceil$  为数组列数
- 哈希函数个数  $d = \lceil \ln(1/\delta) \rceil$  为数组行数

# CM Sketch

**问题定义：统计流数据中出现频率大于给定值s的元素。**

- **CM伪代码**

**CM\_Sketch() :**

**$h_1, h_2, \dots, h_d$  //Hash Function**

**$CM \leftarrow \text{array}[d, w]$**

**For every element  $(e, f)$  :**

**for  $i$  from 1 to  $d$ :**

**$CM[i, h_i(e)] \leftarrow CM[i, h_i(e)] + f$**



# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

CM

0	0	0	0
0	0	0	0
0	0	0	0

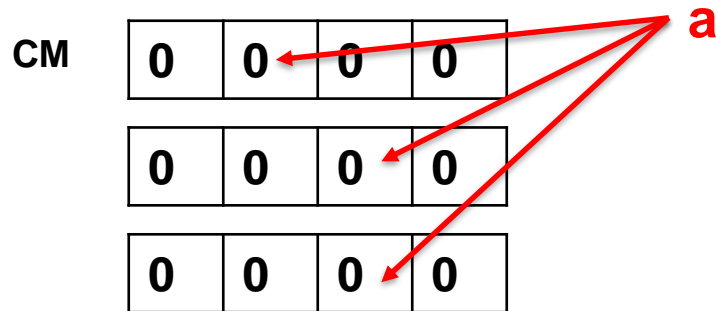
**a**

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

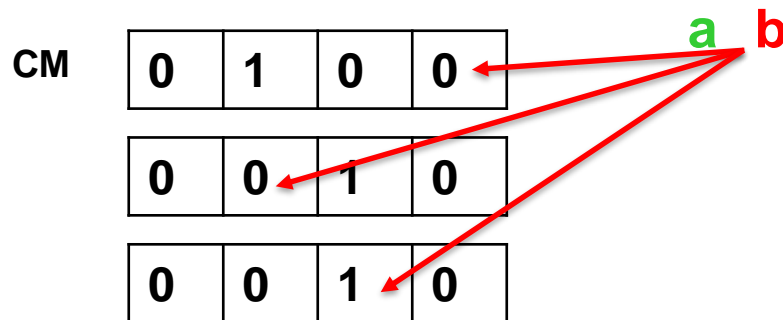


$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

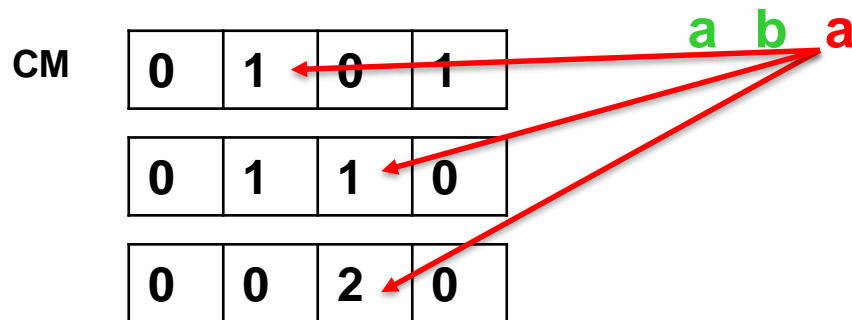


$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

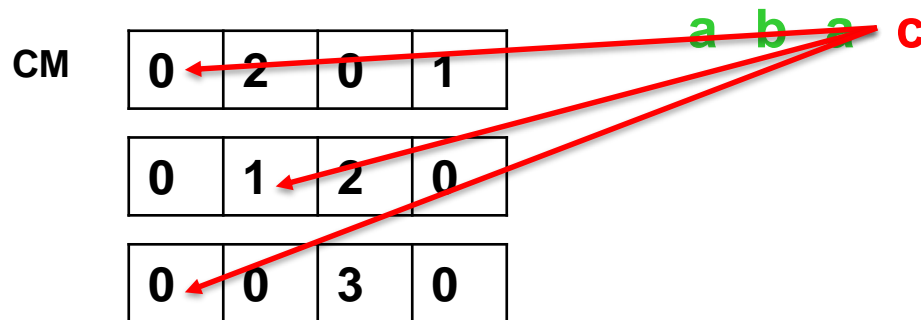


$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

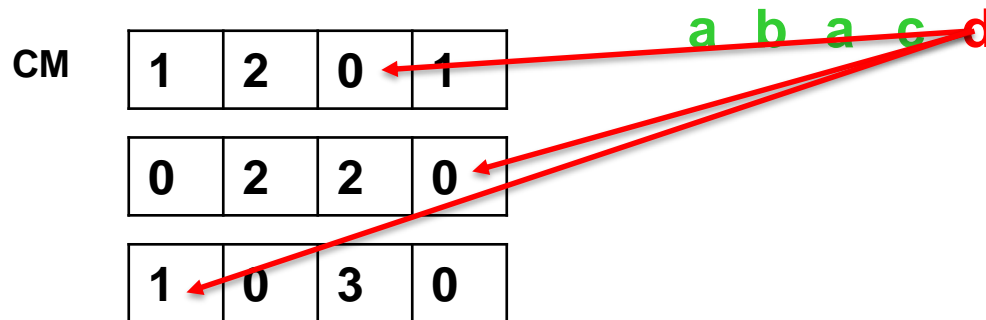


$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

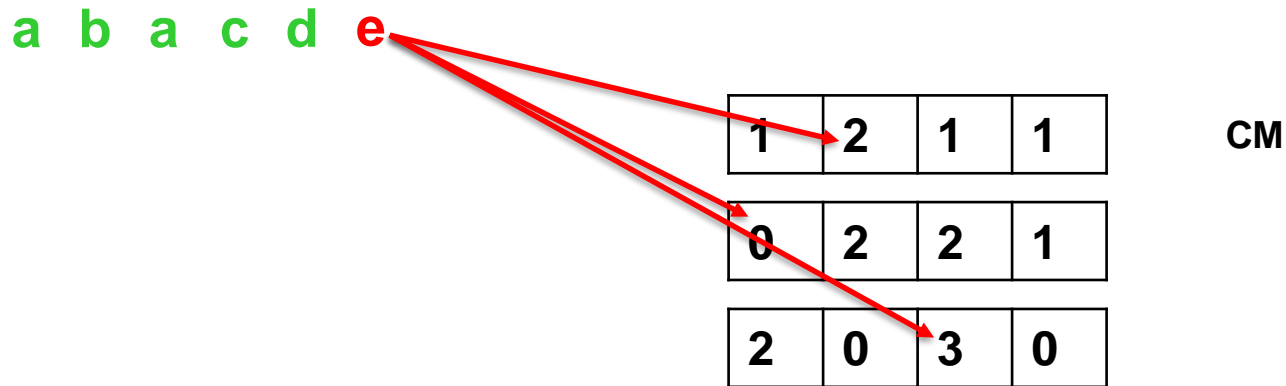


$$CM[i, hi(e)] \leftarrow CM[i, hi(e)] + f$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$



$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$\epsilon=1, \delta=0.1$

$w=\lceil e/\epsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$

a b a c d e a

1	3	1	1
1	2	2	1
2	0	4	0

CM

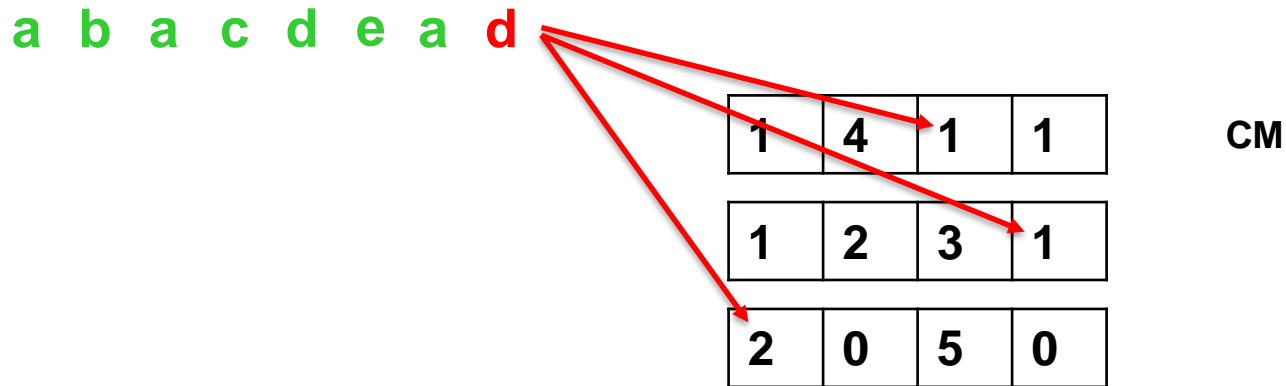
$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$



# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

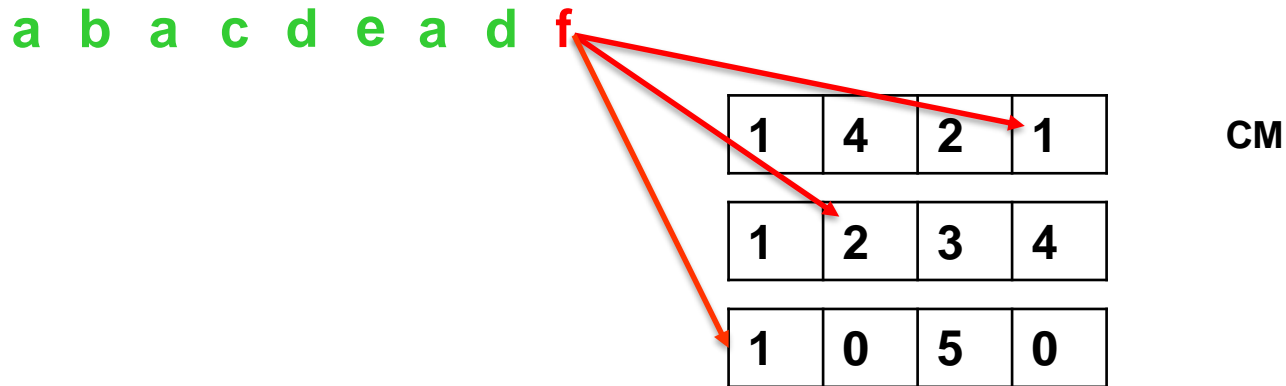


$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$



$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

a b a c d e a d f a

1	4	2	2
1	3	3	4
2	0	5	0

CM

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

a b a c d e a d f a d

1	5	2	2
1	3	4	4
2	0	6	0

CM

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

a b a c d e a d f a d

1	5	3	2
1	3	4	5
3	0	6	0

CM

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

$$Est(e) = \min_{1 \leq i \leq d} \{CM[i, h_i(e)]\}$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

a b a c d e a d f a d

1	5	3	2
1	3	4	5
3	0	6	0

CM

Est(a) = 4

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

$$Est(e) = \min_{1 \leq i \leq d} \{CM[i, h_i(e)]\}$$

# CM Sketch

$$\varepsilon=1, \delta=0.1$$

$$w=\lceil e/\varepsilon \rceil=4, d=\lceil \ln(1/\delta) \rceil=3$$

a b a c d e a d f a d

1	5	3	2
1	3	4	5
3	0	6	0

CM

Est(d) = 3

$$CM[i, h_i(e)] = CM[i, h_i(e)] + 1$$

$$Est(e) = \min_{1 \leq i \leq d} \{CM[i, h_i(e)]\}$$

# 数据流算法算法比较

	算法	$\epsilon$ -Deficient Synopsis	内存消耗
确定性 算法	Majority [1]	100%	1
	Misra-Gries [5]	100% ( $E \leq \epsilon$ )	M
	Space-Saving [4]	100% ( $E \leq \epsilon$ )	M
	Lossy Counting [3]	100%	$\lceil 1/\epsilon \log(\epsilon N) \rceil$
随机性 算法	Sticky Sampling [3]	$1-\delta$	$\lceil 2/\epsilon \ln(s^{-1}\delta^{-1}) \rceil$
	CM-Sketch [2]	$1-\delta$	$\lceil e/\epsilon \rceil \times \lceil \ln(1/\delta) \rceil$



# 参考文献

---

1. G. Cormode and M. Hadjieleftheriou. Finding frequent items in data Streams. In VLDB 2008.
2. G. Cormode and S. Muthukrishnan. An improved data stream summary: The count-min sketch and its applications. Journal of Algorithms, 2005.
3. G. Manku and R. Motwani. Approximate frequency counts over data streams. In VLDB 2002.
4. A. Metwally, D. Agrawal, and A. E. Abbadi. An integrated efficient solution for computing frequent and top-k elements in data streams. ACM Transactions on Database Systems, 2006.
5. J. Misra and D. Gries. Finding repeated elements. Science of Computer Programming, 1982.

---

谢谢