

大数据近似算法

- ❖ 研究背景与计算模型
- ❖ 随机采样算法
- ❖ 基于计数的近似算法
- ❖ 基于哈希的近似算法
- ❖ 研究成果简介



基于计数的近似算法

- ❖ 多数问题(Majority)
- ❖ Misra-Gries(MG)算法
- ❖ MG合并算法



多数问题(Majority)

- ❖ 输入: N 个元素
- ❖ 输出: **Majority**, 即出现次数超过 $N/2$ 的元素
- ❖ 算法1: 扫描所有元素, 给每个出现过的元素分配一个计数器记录其出现次数, 如果某个元素出现次数超过 $N/2$, 即为**Majority**。

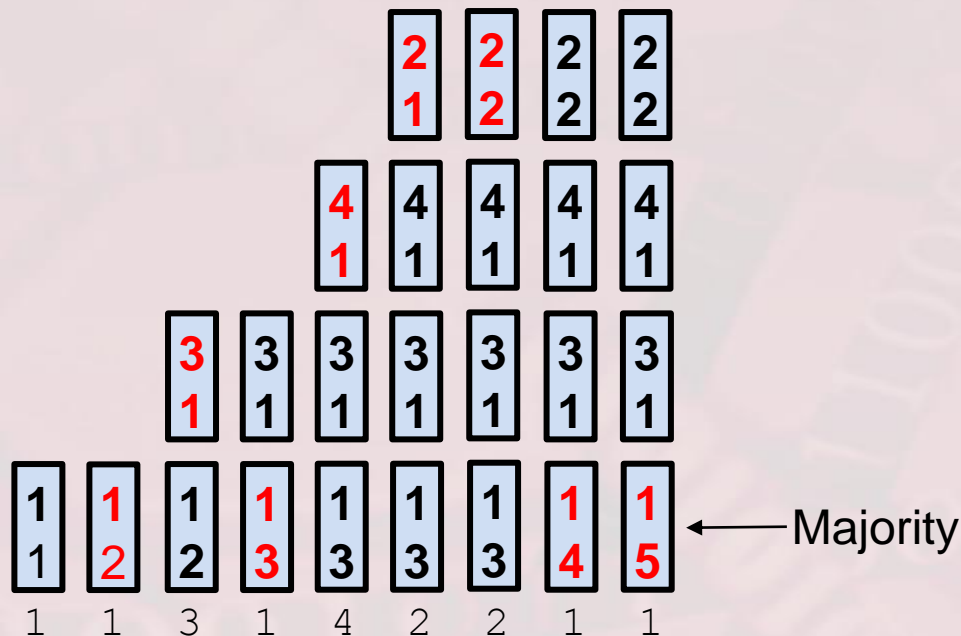
Majority =1

1 1 3 1 4 2 2 1 1



多数问题(Majority)

❖ 算法运行过程:



多数问题(Majority)

- ❖ 输入: N 个元素
- ❖ 输出: **Majority**, 即出现次数超过 $N/2$ 的元素
- ❖ 算法1: 扫描所有元素, 给每个出现过的元素分配一个计数器记录其出现次数, 如果某个元素出现次数超过 $N/2$, 即为**Majority**。
- ❖ 时间复杂度: $O(N)$ (扫描一遍)
- ❖ 空间复杂度: $O(N)$



多数问题(Majority)

- ❖ 输入: N 个元素
- ❖ 输出: **Majority**, 即出现次数超过 $N/2$ 的元素
- ❖ 算法2: 对每个出现过的元素, 扫描一遍 N 个元素, 并记录其出现次数。如果某个元素出现次数超过 $N/2$, 即为**Majority**。
- ❖ 时间复杂度: $O(N^2)$ (如果 N 个元素都不相同, 需要扫描 N 遍)
- ❖ 空间复杂度: $O(1)$



多数问题(Majority)

- ❖ 输入: N 个元素
- ❖ 输出: **Majority**, 即出现次数超过 $N/2$ 的元素
- ❖ 算法3: 扫描所有元素, 储存一个计数器与对应元素。当扫描到的元素与存储元素相同时, 给计数器加1; 当扫描到的元素与存储元素不同时, 给计数器减一; 当计数器归零时, 重新开始。
- ❖ 性质: 如果存在**Majority**, 则扫描完毕时, 存储元素一定是**Majority**。

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



多数问题(Majority)

- ❖ 输入：N个元素
- ❖ 输出：Majority, 即出现次数超过 $N/2$ 的元素
- ❖ 算法3：扫描所有元素，储存一个计数器与对应元素。当扫描到的元素与存储元素相同时，给计数器加1；当扫描到的元素与存储元素不同时，给计数器减一；当计数器归零时，重新开始。
- ❖ 扫描第二遍，确定候选元素是否为Majority
- ❖ 时间复杂度： $O(N)$, 空间复杂度： $O(1)$



Majority算法分析

❖ 定理：如果一个元素出现次数超过 $N/2$ ，那么它一定会被存储。

- N = 数据流中的元素个数
- 每次减一操作可以看成将两个不同元素抵消（当前计数器中的元素以及当前扫描到的元素）
- 总共有 N 个元素可以抵消
- 最多能执行 $N/2$ 次减一操作
- 如果一个元素出现次数超过 $N/2$ ，那么它最多被减 $N/2$ 次，因此一定会被存储。

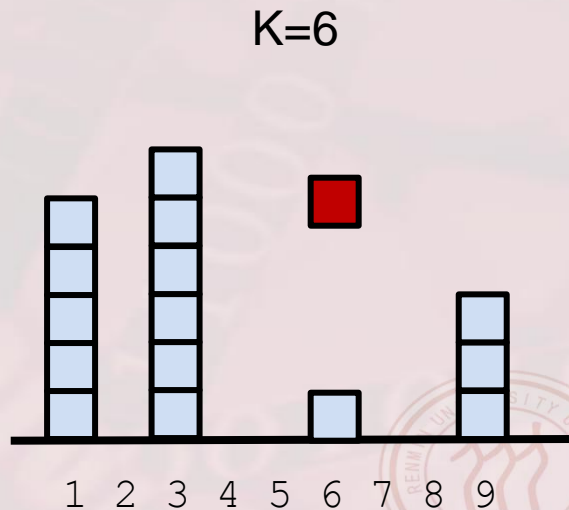


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过N/k的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

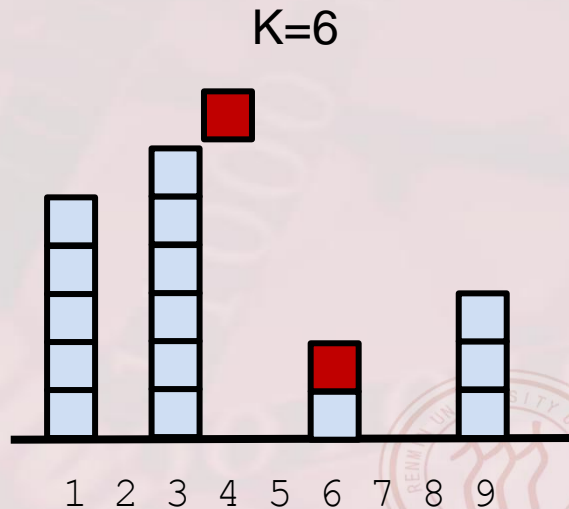


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满足k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

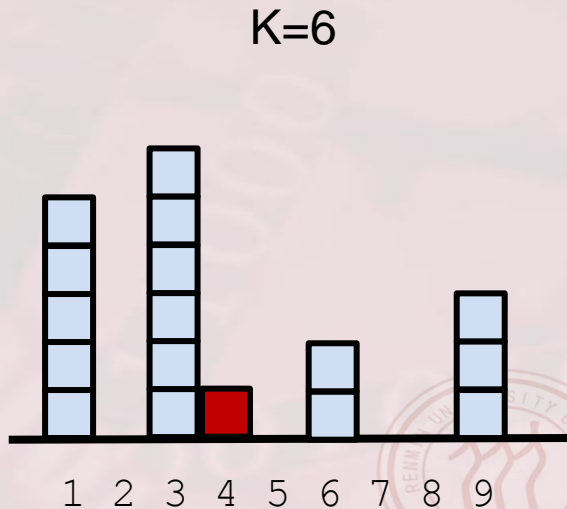


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满足k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

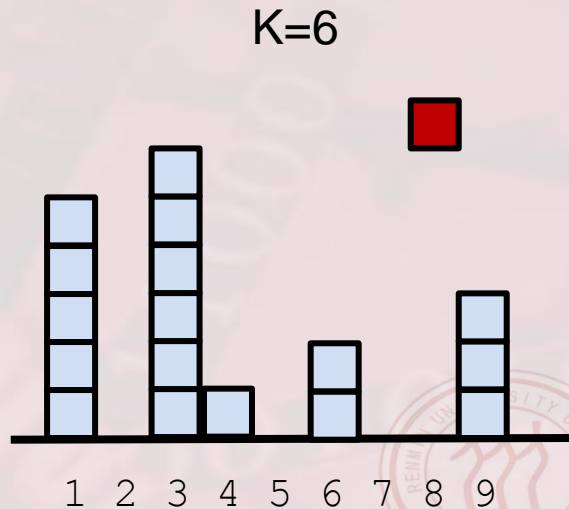


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

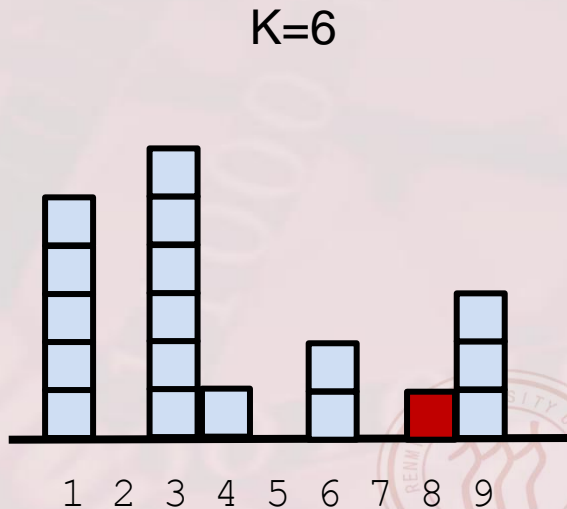


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

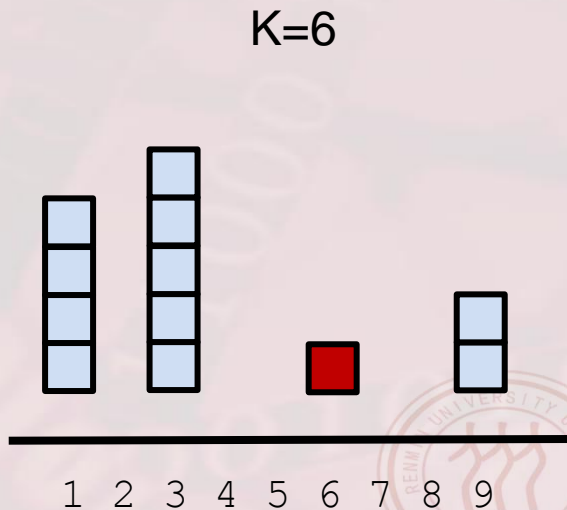


Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。



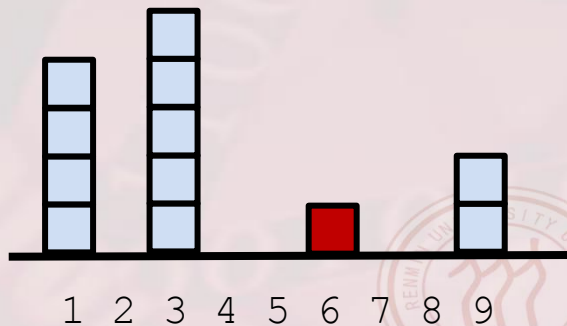
Misra-Gries算法

- ❖ 输入：N个元素，整数k
- ❖ 输出：出现次数超过 N/k 的元素

算法：保存k-1个元素与其对应的计数器。对于一个新元素

- 如果该元素已被记录，则将该元素对应的计数器加一
- 否则，如果未满k-1个元素，则记录新元素，计数器设为一
- 否则，所有计数器减一，移除计数器为0的元素。

K=6



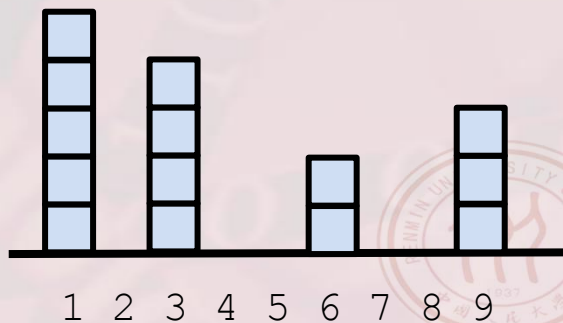
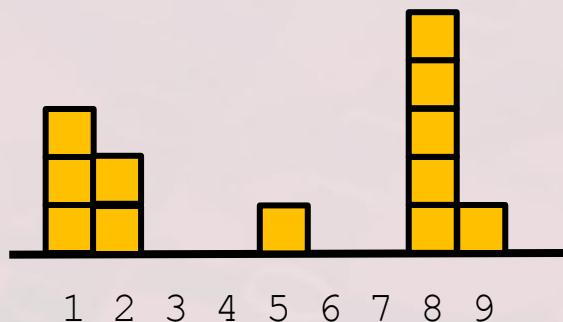
Misra-Gries算法分析

- ❖ 定理：如果一个元素未被记录，它的出现次数不会超过 N/k
 - N = 数据流中的元素个数
 - 误差 \leq 所有元素减一出现的次数
 - 每次减一，我们减去了 k 个元素：1 个新元素和 $k-1$ 个老元素
 - 最多能减掉 N 个元素
 - 最多能产生 N/k 次所有元素减一操作
 - 1% 误差 \rightarrow 99 个计数器



Misra-Gries合并算法

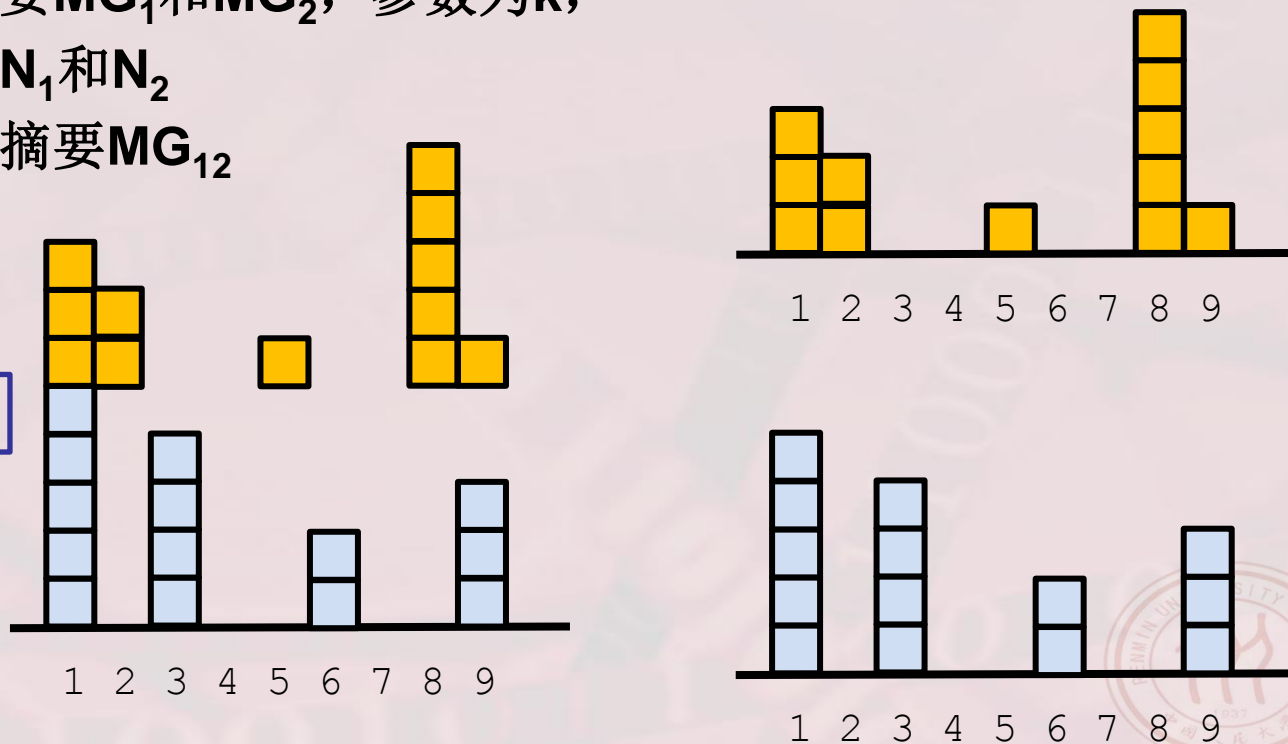
- ❖ 输入：两个摘要 MG_1 和 MG_2 ，参数为 k ，
- ❖ 元数据大小为 N_1 和 N_2
- ❖ 输出：合集的摘要 MG_{12}



Misra-Gries合并算法

- ❖ 输入：两个摘要 MG_1 和 MG_2 ，参数为 k ，
- ❖ 元数据大小为 N_1 和 N_2
- ❖ 输出：合集的摘要 MG_{12}

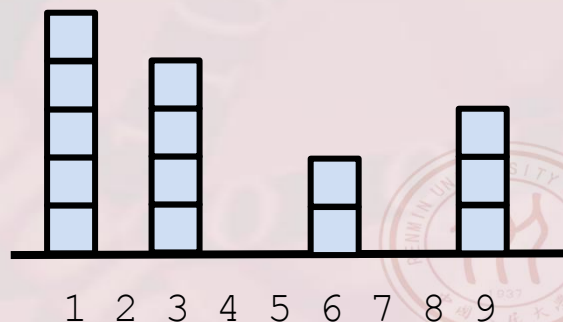
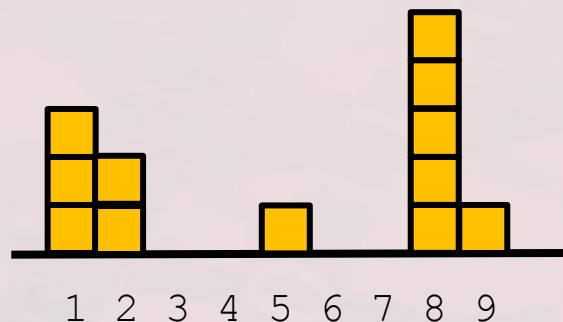
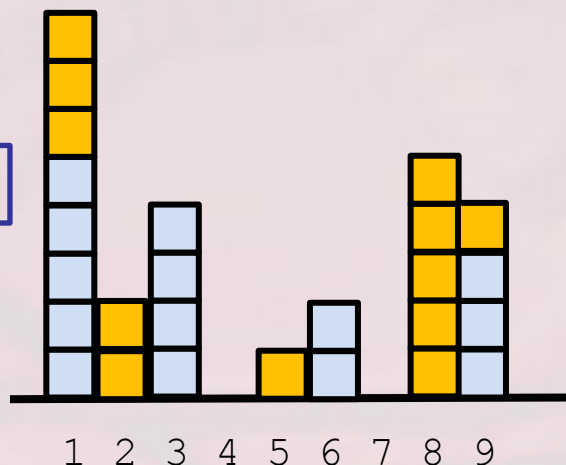
1. 对应计数器相加



Misra-Gries合并算法

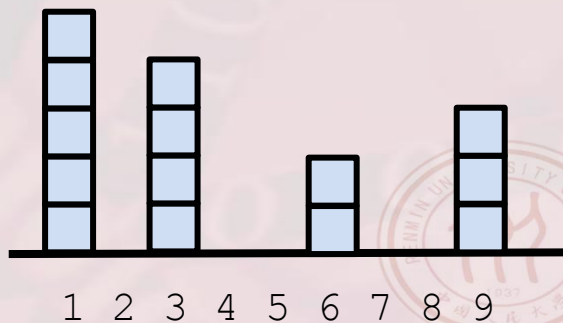
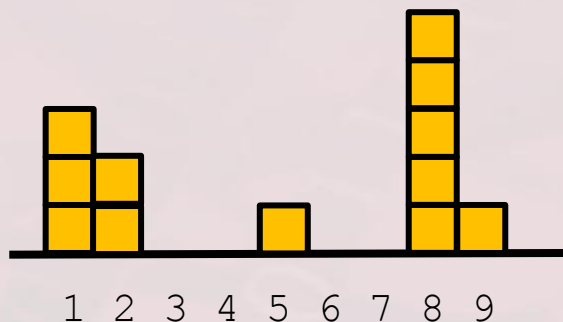
- ❖ 输入：两个摘要 MG_1 和 MG_2 ，参数为 k ，
- ❖ 元数据大小为 N_1 和 N_2
- ❖ 输出：合集的摘要 MG_{12}

1. 对应计数器相加

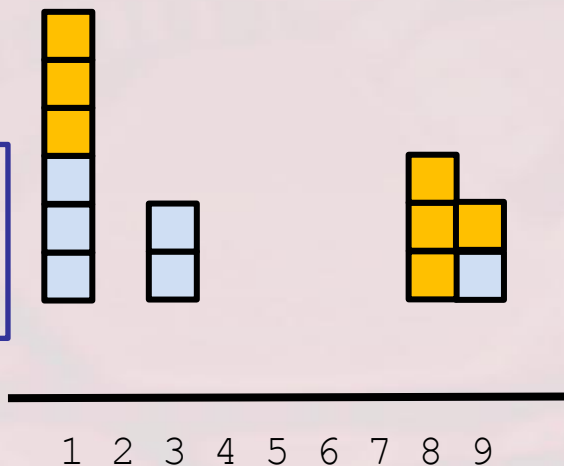


Misra-Gries合并算法

- ❖ 输入：两个摘要 MG_1 和 MG_2 ，参数为 k ，
- ❖ 元数据大小为 N_1 和 N_2
- ❖ 输出：合集的摘要 MG_{12}



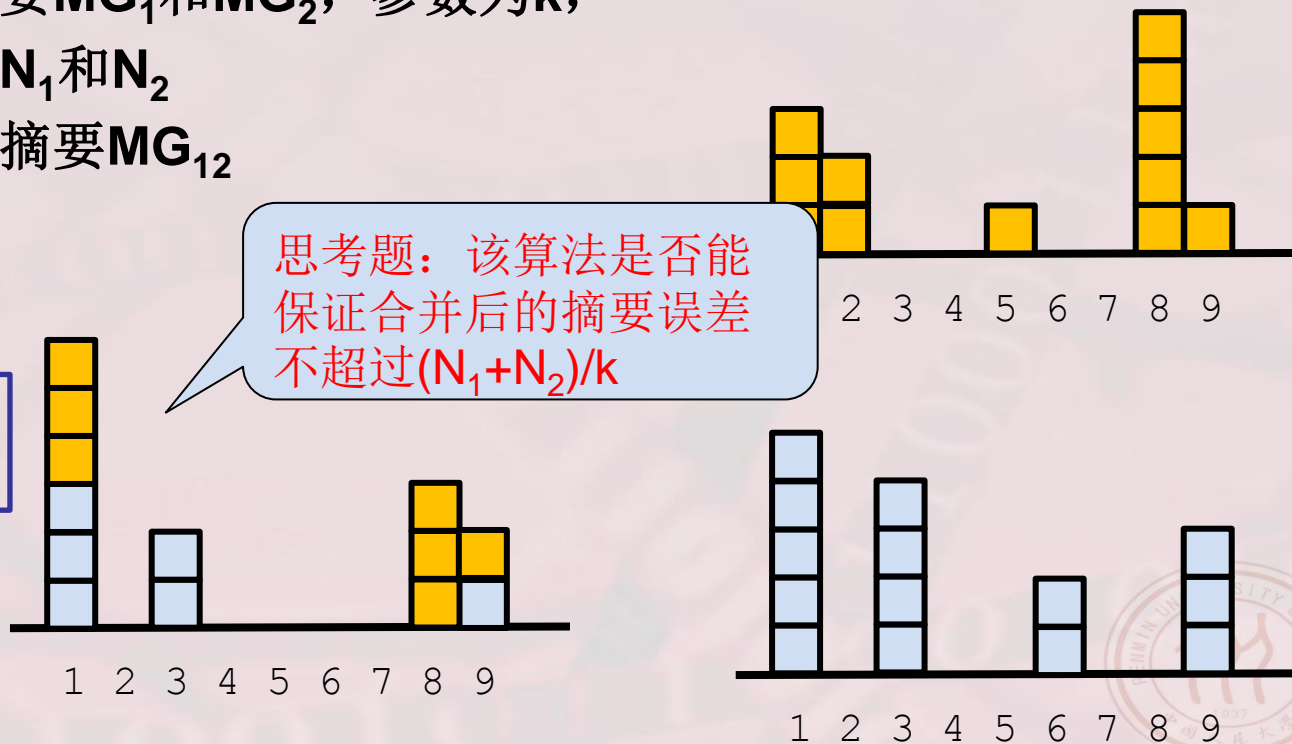
2. 找出第 k 大计数器 C_k ，
从所有计数器中减去
 C_k



Misra-Gries合并算法

- ❖ 输入：两个摘要 MG_1 和 MG_2 ，参数为 k ，
- ❖ 元数据大小为 N_1 和 N_2
- ❖ 输出：合集的摘要 MG_{12}

3. 移除负数或者零值计数器



Misra-Gries算法分析

❖ 定理：如果一个元素未被记录，它的出现次数不会超过 $(N-M)/k$

- M 为 MG 算法中计数器之和

- N = 数据流中的元素个数

- 误差 \leq 所有元素减一出现的次数

- 每次减一，我们减去了 k 个元素：1 个新元素和 $k-1$ 个老元素

- 最多能减掉 $N-M$ 个元素

- 最多能产生 $(N-M)/k$ 次所有元素减一操作

- 1% 误差 \rightarrow 99 个计数器



Misra-Gries合并算法分析

- ❖ 定理：如果一个元素未被记录，它的出现次数不会超过 $(N-M)/k$
 - 合并算法从计数器之和减去了 kC_k

