

# COMP9313: Big Data Management



**Lecturer: Xin Cao**

**Course web site: <http://www.cse.unsw.edu.au/~cs9313/>**

# **Chapter 6.2: Mining Data Streams II**

## **Part 3: Filtering Data Streams**

# Filtering Data Streams

- ❖ Each element of data stream is a tuple
- ❖ Given a list of keys **S**
- ❖ **Determine which tuples of stream are in S**
- ❖ Obvious solution: Hash table
  - But suppose we **do not have enough memory** to store all of **S** in a hash table
    - ▶ E.g., we might be processing millions of filters on the same stream

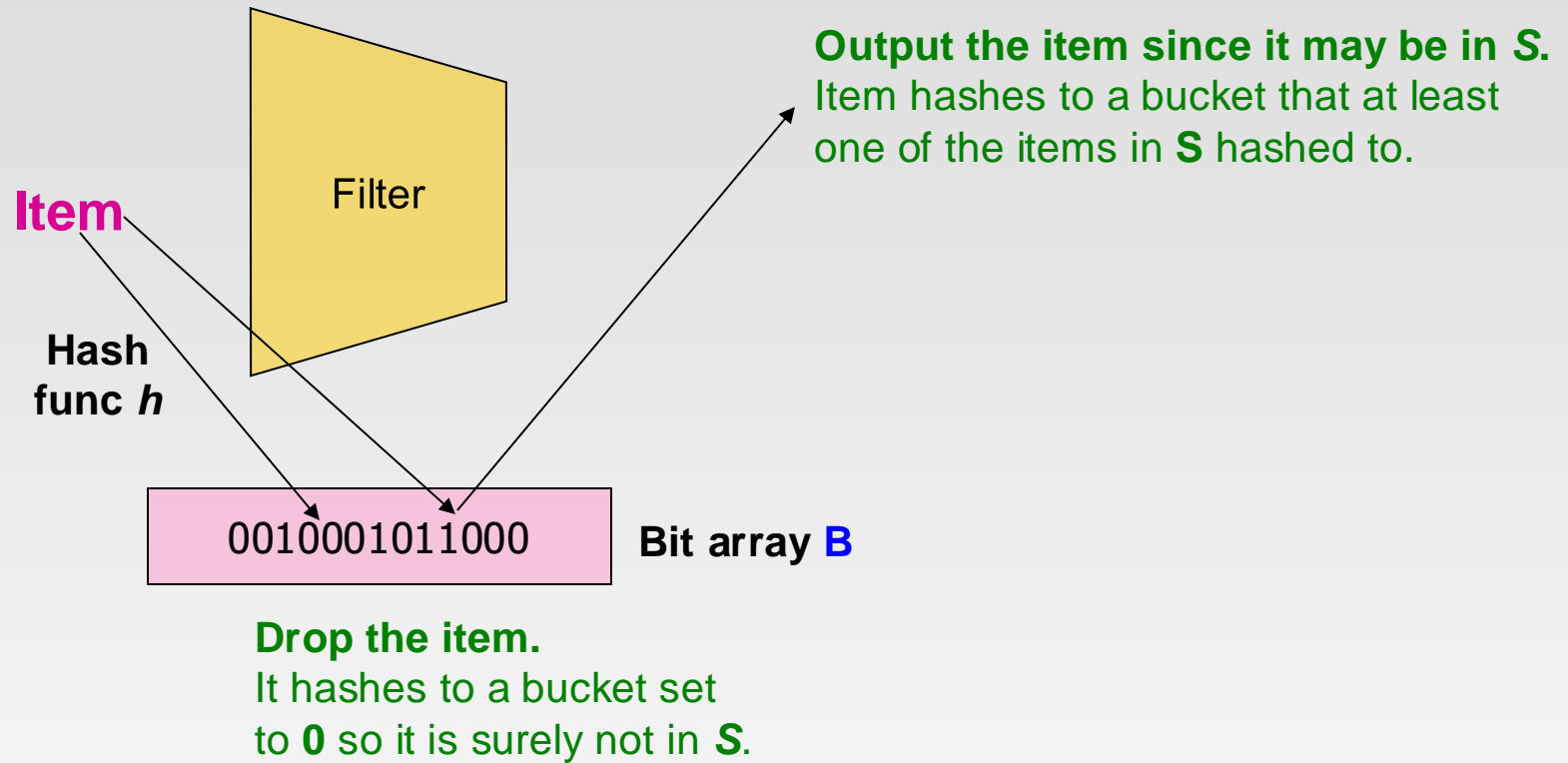
# Applications

- ❖ Example: Email spam filtering
  - We know 1 billion “good” email addresses
  - If an email comes from one of these, it is **NOT** spam
- ❖ Publish-subscribe systems
  - You are collecting lots of messages (news articles)
  - People express interest in certain sets of keywords
  - Determine whether each message matches user’s interest

# First Cut Solution (1)

- ❖ Given a set of keys  $S$  that we want to filter
- ❖ Create a bit array  $B$  of  $n$  bits, initially all 0s
- ❖ Choose a hash function  $h$  with range  $[0, n)$
- ❖ Hash each member of  $s \in S$  to one of  $n$  buckets, and set that bit to 1, i.e.,  $B[h(s)] = 1$
- ❖ Hash each element  $a$  of the stream and output only those that hash to bit that was set to 1
  - Output  $a$  if  $B[h(a)] == 1$

# First Cut Solution (2)



## ❖ Creates false positives but no false negatives

- If the item is in  $S$  we surely output it, if not we may still output it

# First Cut Solution (3)

- ❖  $|S| = 1$  billion email addresses  
 $|B| = 1\text{GB} = 8$  billion bits
- ❖ If the email address is in  $S$ , then it surely hashes to a bucket that has the bit set to 1, so it always gets through (*no false negatives*)
  - False negative: a result indicates that a condition failed, while it actually was successful
- ❖ Approximately  $1/8$  of the bits are set to 1, so about  $1/8$ th of the addresses not in  $S$  get through to the output (*false positives*)
  - False positive: a result that indicates a given condition has been fulfilled, when it actually has not been fulfilled
  - Actually, less than  $1/8$ th, because more than one address might hash to the same bit
  - Since the majority of emails are spam, eliminating  $7/8$ th of the spam is a significant benefit

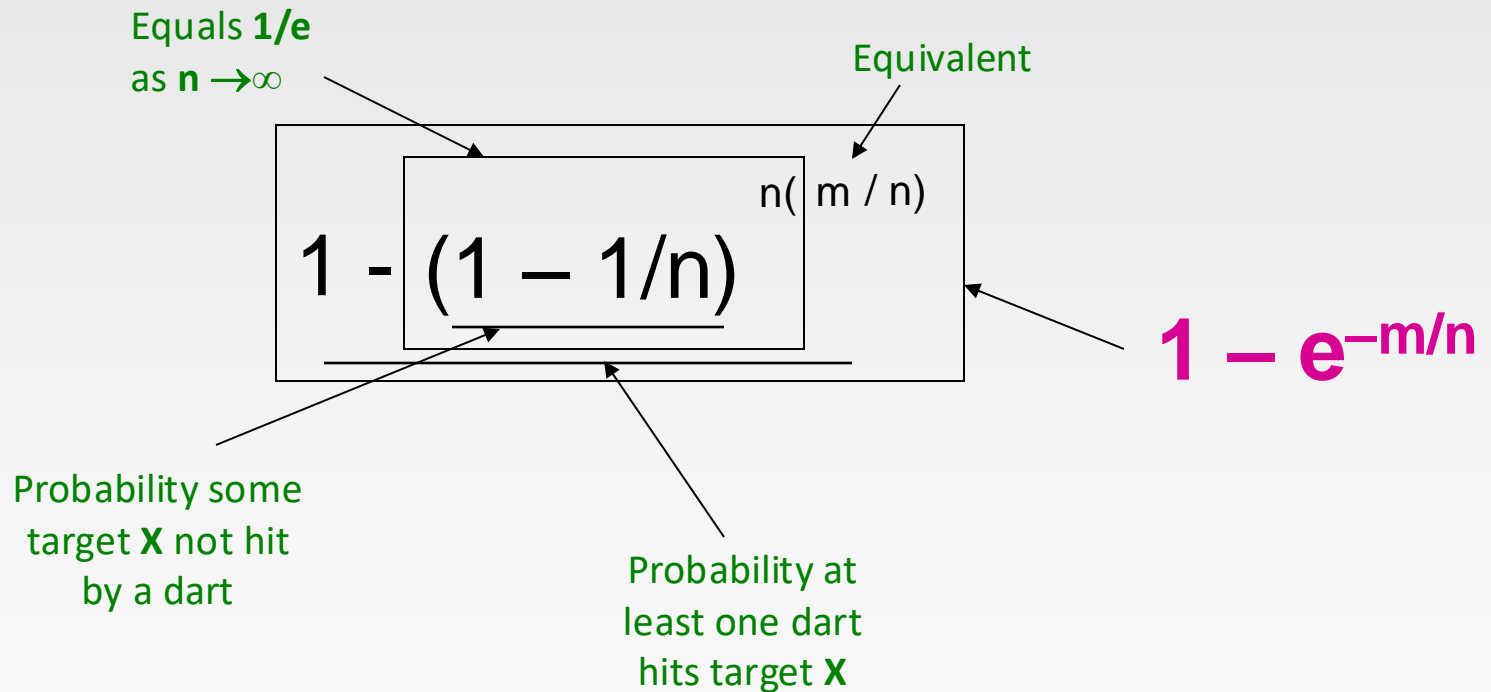


# Analysis: Throwing Darts (1)

- ❖ More accurate analysis for the number of **false positives**
- ❖ **Consider:** If we throw  $m$  darts into  $n$  equally likely targets, **what is the probability that a target gets at least one dart?**
- ❖ **In our case:**
  - **Targets** = bits/buckets
  - **Darts** = hash values of items

# Analysis: Throwing Darts (2)

- ❖ We have  $m$  darts,  $n$  targets
- ❖ What is the probability that a target gets at least one dart?



# Analysis: Throwing Darts (3)

❖ Fraction of 1s in the array B

$$= \text{probability of false positive} = 1 - e^{-m/n}$$

❖ **Example:**  $10^9$  darts,  $8 \cdot 10^9$  targets

➤ Fraction of 1s in B =  $1 - e^{-1/8} = 0.1175$

▶ Compare with our earlier estimate:  $1/8 = 0.125$

# Bloom Filter

- ❖ Consider:  $|\mathbf{S}| = m$ ,  $|\mathbf{B}| = n$
- ❖ Use  $k$  independent hash functions  $h_1, \dots, h_k$
- ❖ Initialization:
  - Set  $\mathbf{B}$  to all 0s
  - Hash each element  $s \in \mathbf{S}$  using each hash function  $h_i$ , set  $\mathbf{B}[h_i(s)] = 1$  (for each  $i = 1, \dots, k$ )
- ❖ Run-time:
  - When a stream element with key  $x$  arrives
    - ▶ If  $\mathbf{B}[h_i(x)] = 1$  for all  $i = 1, \dots, k$  then declare that  $x$  is in  $\mathbf{S}$ 
      - That is,  $x$  hashes to a bucket set to 1 for every hash function  $h_i(x)$
    - ▶ Otherwise discard the element  $x$

# Bloom Filter

Start with an  $n$  bit array, filled with 0s.

$B$

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Hash each item  $x_j$  in  $S$  for  $k$  times. If  $H_i(x_j) = a$ , set  $B[a] = 1$ .

$B$

0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

To check if  $y$  is in  $S$ , check  $B$  at  $H_i(y)$ . All  $k$  values must be 1.

$B$

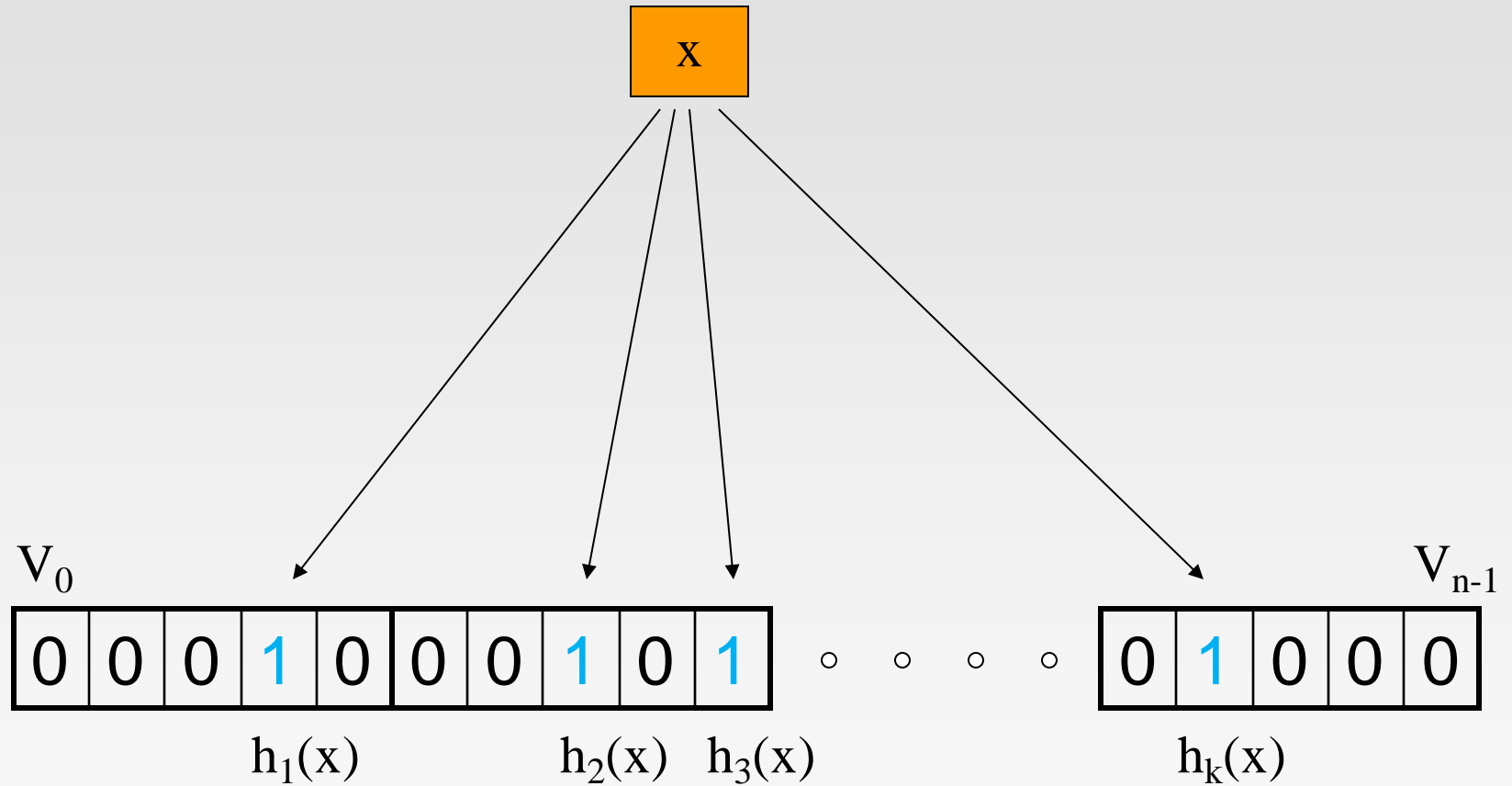
0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Possible to have a false positive; all  $k$  values are 1, but  $y$  is not in  $S$ .

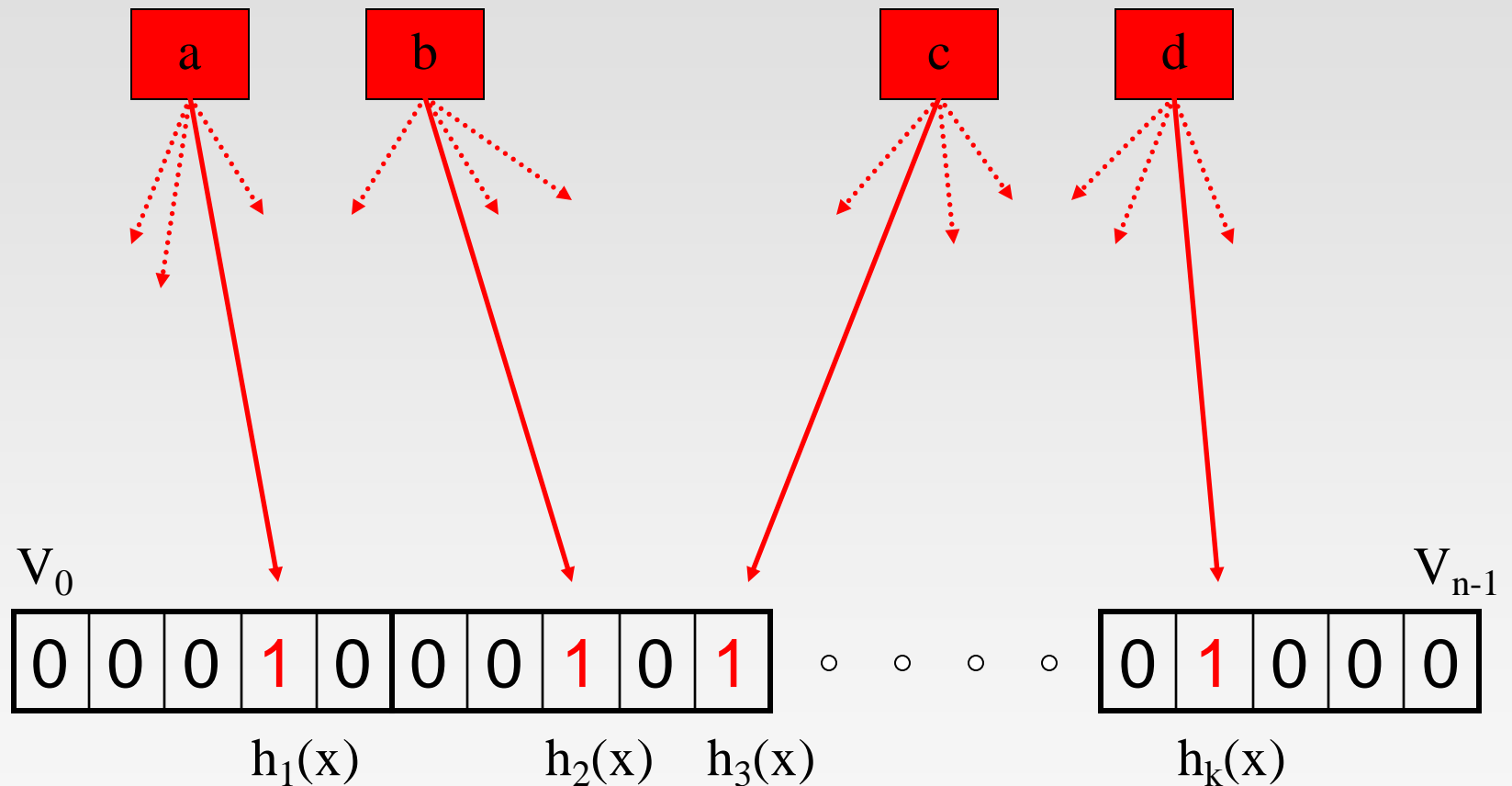
$B$

0	1	0	0	1	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

# Bloom Filter Hashing



# Bloom Errors



x didn't appear, yet its bits are already set

# Bloom Filter Example

- ❖ Consider a Bloom filter of size  $m=10$  and number of hash functions  $k=3$ . Let  $H(x)$  denote the result of the three hash functions.

- ❖ The 10-bit array is initialized as below

0	1	2	3	4	5	6	7	8	9
0	0	0	0	0	0	0	0	0	0

- ❖ Insert  $x_0$  with  $H(x_0) = \{1, 4, 9\}$

0	1	2	3	4	5	6	7	8	9
0	1	0	0	1	0	0	0	0	1

- ❖ Insert  $x_1$  with  $H(x_1) = \{4, 5, 8\}$

0	1	2	3	4	5	6	7	8	9
0	1	0	0	1	1	0	0	1	1

- ❖ Query  $y_0$  with  $H(y_0) = \{0, 4, 8\} \Rightarrow ???$

- ❖ Query  $y_1$  with  $H(y_1) = \{1, 5, 8\} \Rightarrow ???$  **False positive!**

- ❖ Another Example: <https://lilmlib.github.io/bloomfilter-tutorial/>



# Bloom Filter – Analysis

- ❖ What fraction of the bit vector  $B$  are 1s?
  - Throwing  $k \cdot m$  darts at  $n$  targets
  - So fraction of 1s is  $(1 - e^{-km/n})$
- ❖ But we have  $k$  independent hash functions and we only let the element  $x$  through if all  $k$  hash element  $x$  to a bucket of value 1
- ❖ So, false positive probability =  $(1 - e^{-km/n})^k$

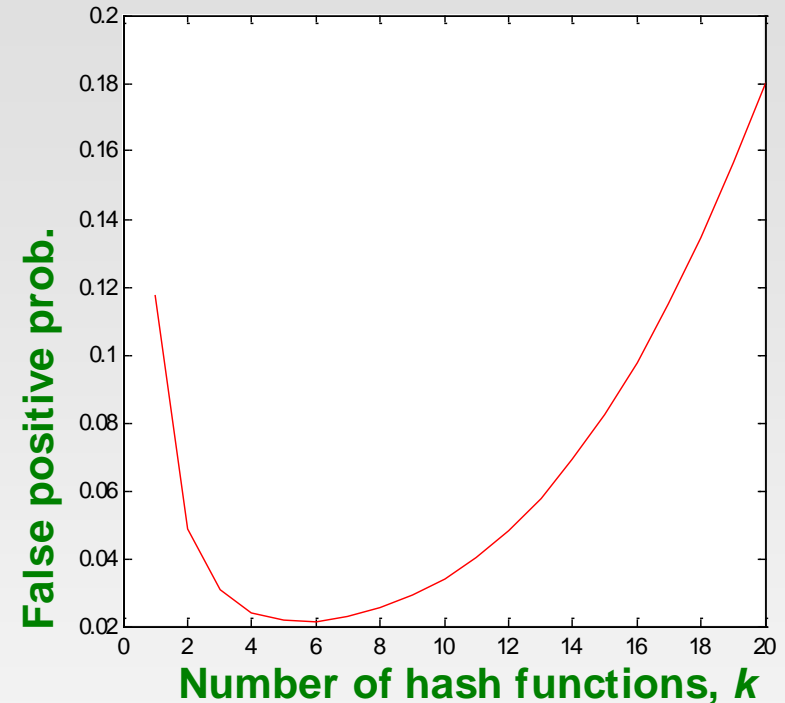
# Bloom Filter – Analysis (2)

❖  $m = 1$  billion,  $n = 8$  billion

➤  $k = 1$ :  $(1 - e^{-1/8}) = 0.1175$

➤  $k = 2$ :  $(1 - e^{-1/4})^2 = 0.0493$

❖ What happens as we keep increasing  $k$ ?



❖ “Optimal” value of  $k$ :  $n/m \ln(2)$

➤ In our case: Optimal  $k = 8 \ln(2) = 5.54 \approx 6$

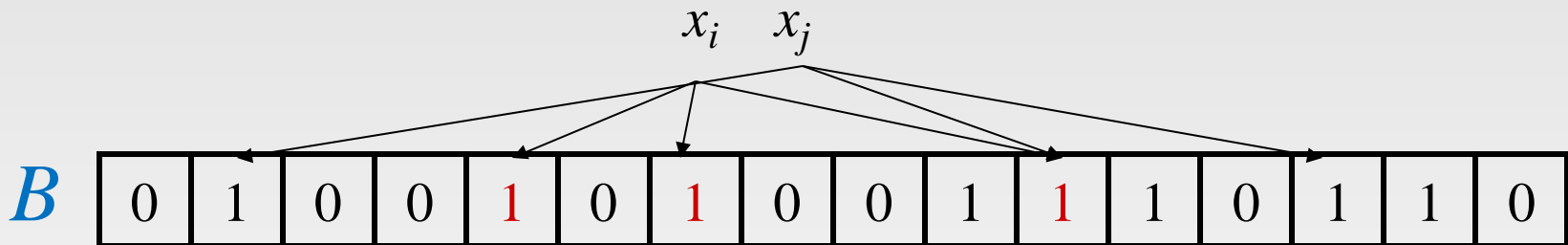
▶ Error at  $k = 6$ :  $(1 - e^{-6/8})^6 = 0.02158$

# Bloom Filter: Wrap-up

- ❖ Bloom filters guarantee no false negatives, and use limited memory
  - Great for pre-processing before more expensive checks
- ❖ Suitable for hardware implementation
  - Hash function computations can be parallelized
- ❖ Is it better to have 1 big **B** or  $k$  small **B**s?
  - It is the same:  $(1 - e^{-km/n})^k$  vs.  $(1 - e^{-m/(n/k)})^k$
  - But keeping 1 big **B** is simpler

# Handling Deletions

- ❖ Bloom filters can handle insertions, but not deletions.
- ❖ If deleting  $x_i$  means resetting 1s to 0s, then deleting  $x_j$  will “delete”  $x_i$ .



- ❖ Can Bloom filters handle deletions?
  - Use Counting Bloom Filters to track insertions/deletions

# Counting Bloom Filters

Start with an  $n$  bit array, filled with 0s.

$B$

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Hash each item  $x_j$  in  $S$  for  $k$  times. If  $H_i(x_j) = a$ , add 1 to  $B[a]$ .

$B$

0	3	0	0	1	0	2	0	0	3	2	1	0	2	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

To delete  $x_j$  decrement the corresponding counters.

$B$

0	2	0	0	0	0	2	0	0	3	2	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Can obtain a corresponding Bloom filter by reducing to 0/1.

$B$

0	1	0	0	0	0	1	0	0	1	1	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

## **Part 4: Finding Frequent Elements (Majority and Heavy Hitters)**

# The Majority Problem

- ❖ Given a stream of elements, find the majority if there is one
  - A majority element in the data stream (assume that we have received  $n$  elements already) is an element that appears more than  $n/2$  times
- ❖ A A B C D B A A B B A A A A A C C C D A B A A A
  - Answer: A
- ❖ It is trivial if we have enough memory
  - For each received element, keep a counter for it. Once receiving it again, increase the counter
  - Can use the binary search tree/hashmap to store the elements
  - $O(n \log n)/O(n)$  complexity and  $O(n)$  space
- ❖ What if we only have limited memory?

# Boyer-Moore Voting Algorithm

- ❖ This algorithm takes  $O(n)$  time and  $O(1)$  space
- ❖ Basic idea of the algorithm is if we cancel out each occurrence of an element  $e$  with all the other elements that are different from  $e$ , then  $e$  will exist till the end. Then, we can check if it is indeed the majority element.
- ❖ Thus, the algorithm contains two phases:
  - First pass: find the possible candidate (the element that has the largest frequency in the stream)
  - Second pass: compute its frequency and verify that it is  $> n/2$



# Boyer-Moore Voting Algorithm

## ❖ Phase 1:

- Loop through each element and maintains a count of majority element, and a majority index, maj\_index
- If the next element is same then increment the count, if the next element is not same then decrement the count.
- if the count reaches 0 then changes the maj\_index to the current element and set the count again to 1.

```
maj_index = 0
count = 1
for i in range(len(A)):
    if A[maj_index] == A[i]:
        count += 1
    else:
        count -= 1
        if count == 0:
            maj_index = i
            count = 1
return A[maj_index]
```

# Boyer-Moore Voting Algorithm

- ❖ Example: given a stream as  $A[] = [2, 2, 3, 5, 2, 2, 6]$
- $\text{maj\_index} = 0$ ,  $\text{count} = 1 \rightarrow$  candidate 2?
  - Same as  $a[\text{maj\_index}] \Rightarrow \text{count} = 2$
  - Different from  $a[\text{maj\_index}] \Rightarrow \text{count} = 1$
  - Different from  $a[\text{maj\_index}] \Rightarrow \text{count} = 0$
  - Since  $\text{count} = 0$ , change candidate for majority element to 5  $\Rightarrow \text{maj\_index} = 3$ ,  $\text{count} = 1$
  - Different from  $a[\text{maj\_index}] \Rightarrow \text{count} = 0$
  - Since  $\text{count} = 0$ , change candidate for majority element to 2  $\Rightarrow \text{maj\_index} = 4$
  - Same as  $a[\text{maj\_index}] \Rightarrow \text{count} = 2$
  - Different from  $a[\text{maj\_index}] \Rightarrow \text{count} = 1$
  - Finally, candidate for majority element is 2

# Boyer-Moore Voting Algorithm

- ❖ Phase 2: Just compute the count of the element in the stream for verification

```
count = 0
for i in range(len(A)):
    if A[i] == cand:
        count += 1
if count > len(A)/2:
    return True
else:
    return False
```

- ❖ We can see that this algorithm still requires two passes of the stream, which is actually not possible in most streaming applications.
- ❖ If only one pass and  $O(1)$  space allowed, not possible to get the majority element!

Input is an array: <https://leetcode.com/problems/majority-element/>

# Heavy Hitters

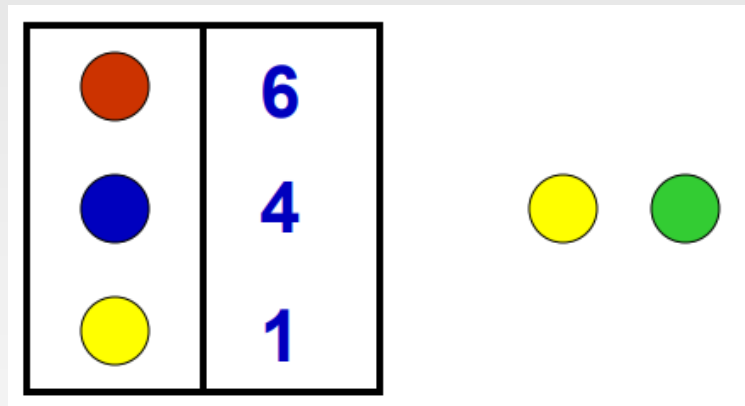
- ❖ A more general problem: find all elements with counts  $> n/k$  ( $k \geq 2$ )
  - There can be at most  $k-1$  such values; and there might be none
  - Trivial if we have enough storage
  
- ❖ Applications
  - Computing popular products. For example,  $A$  could be all of the page views of products on amazon.com yesterday. The heavy hitters are then the most frequently viewed products
  - Computing frequent search queries. For example,  $A$  could be all of the searches on Google yesterday. The heavy hitters are then searches made most often
  - Identifying heavy TCP flows. Here,  $A$  is a list of data packets passing through a network switch, each annotated with a source-destination pair of IP addresses. The heavy hitters are then the flows that are sending the most traffic. This is useful for, among other things, identifying denial-of-service attacks

# Approximate Heavy Hitters

- ❖ There is no exact algorithm that solves the Heavy Hitters problems in one pass while using a sublinear amount of auxiliary space
- ❖ Relaxation, the  $\epsilon$ -approximate heavy hitters problem:
  - If an element has count  $> n/k$ , it must be reported, together with its estimated count with (absolute) error  $< \epsilon n$
  - If an element has count  $< (1/k - \epsilon) n$ , it cannot be reported
  - For elements in between, don't care
- ❖ In fact, we will estimate all counts with at most  $\epsilon n$  error

# Misra-Gries Algorithm

- ❖ Keep  $k-1$  different candidates in hand (thus with space  $O(k)$ )
- ❖ For each element in stream:
  - If item is monitored, increase its counter
  - Else, if  $< k-1$  items monitored, add new element with count 1
  - Else, decrease all counts by 1, and delete element with count 0



- ❖ Each decrease can be charged against  $k$  arrivals of different items, so no item with frequency  $N/k$  is missed
- ❖ But false positive (elements with count smaller than  $n/k$ ) may appear in the result

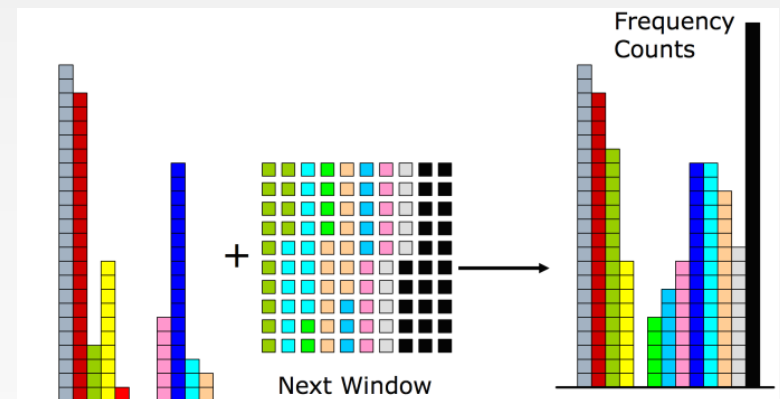
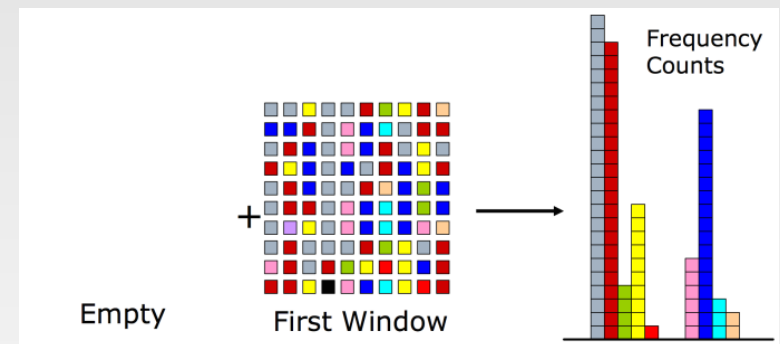
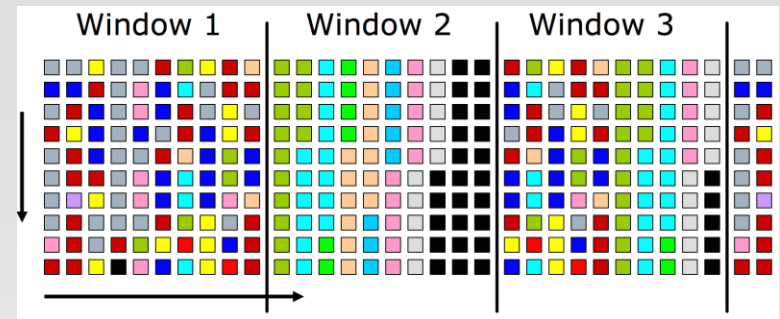
# Misra-Gries Algorithm

- ❖  $[1, 1, 2, 3, 4, 5, 1, 1, 1, 5, 3, 3, 1, 1, 2]$  with  $k=3$ , we want to find element that occurred more than  $15/3 = 5$  times.



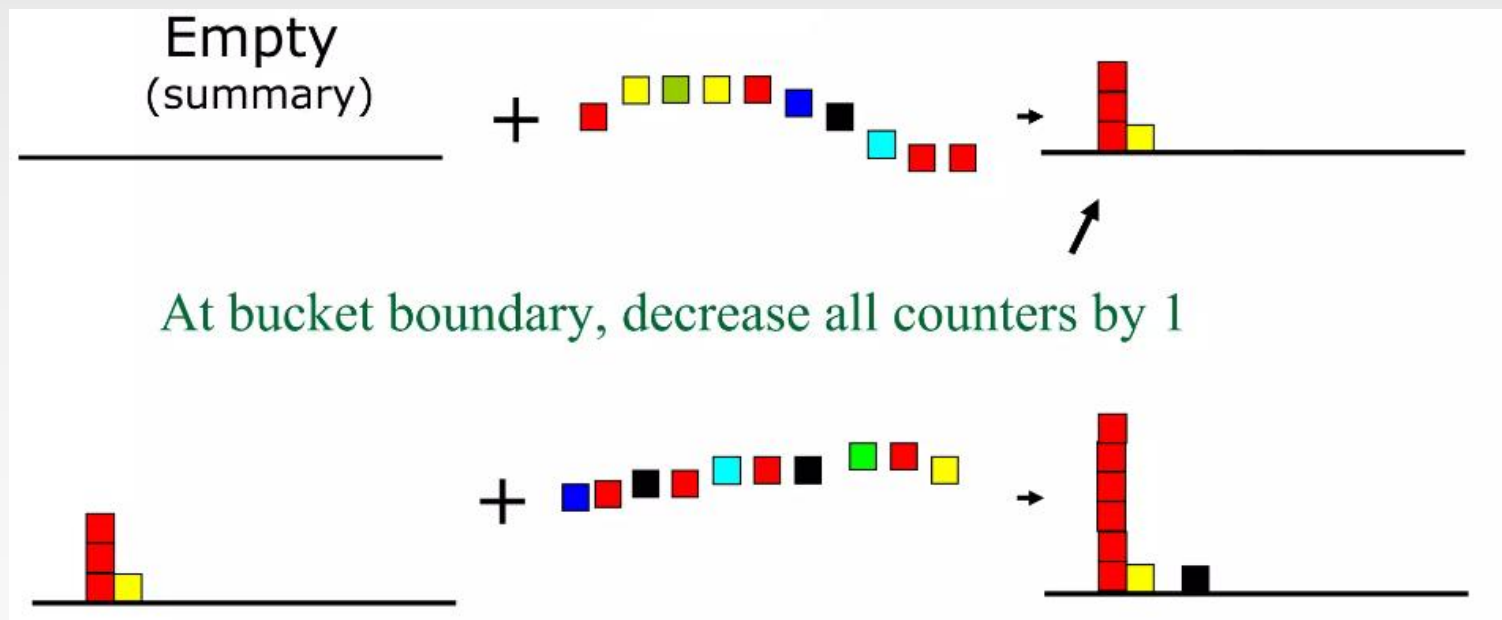
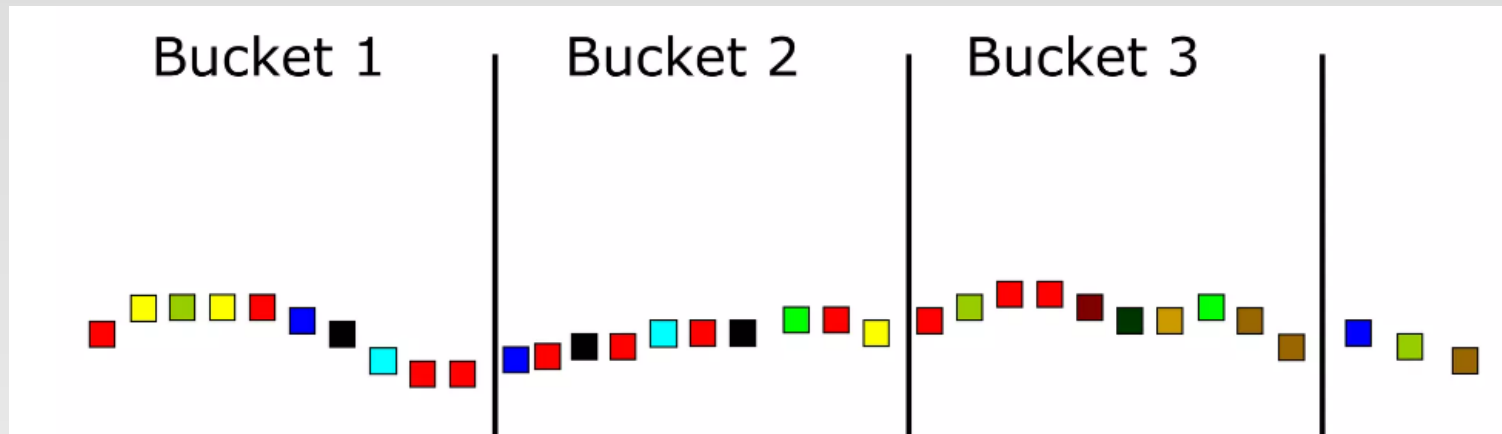
# Lossy Counting

- ❖ Step 1: Divide the incoming data stream into windows, and each window contains  $1/\epsilon$  elements
- ❖ Step 2: Increment the frequency count of each item according to the new window values. After each window, decrement all counters by 1. Drop elements with counter 0.
- ❖ Step 3: Repeat – Update counters and after each window, decrement all counters by 1



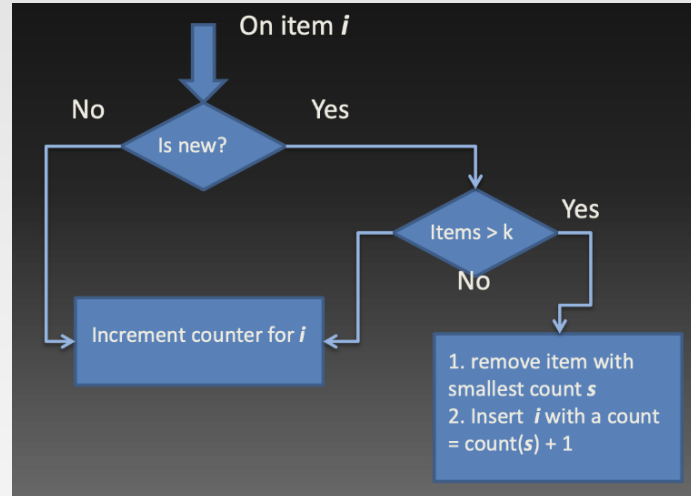


# Lossy Counting



# The Space-Saving Algorithm

- ❖ Keep  $k = 1/\epsilon$  item names and counts, initially zero
- ❖ On seeing new item:
  - If it has a counter, increment counter
  - If not, replace item with least count, increment count

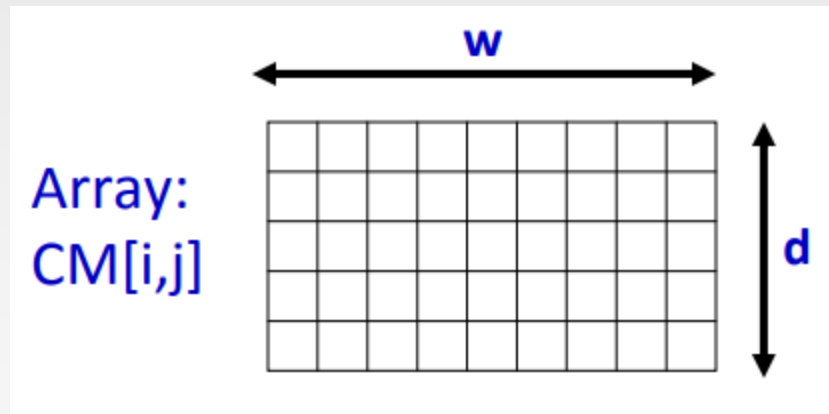


[http://romania.amazon.com/techon/presentations/DataStreamsAlgorithms\\_FlorinManolache.pdf](http://romania.amazon.com/techon/presentations/DataStreamsAlgorithms_FlorinManolache.pdf)

- ❖ Analysis:
  - Smallest counter value, min, is at most  $\epsilon n$
  - True count of an uncounted item is between 0 and min
  - Any item  $x$  whose true count  $> \epsilon n$  is stored
- ❖ So: Find all items with count  $> \epsilon n$ , error in counts  $\leq \epsilon n$

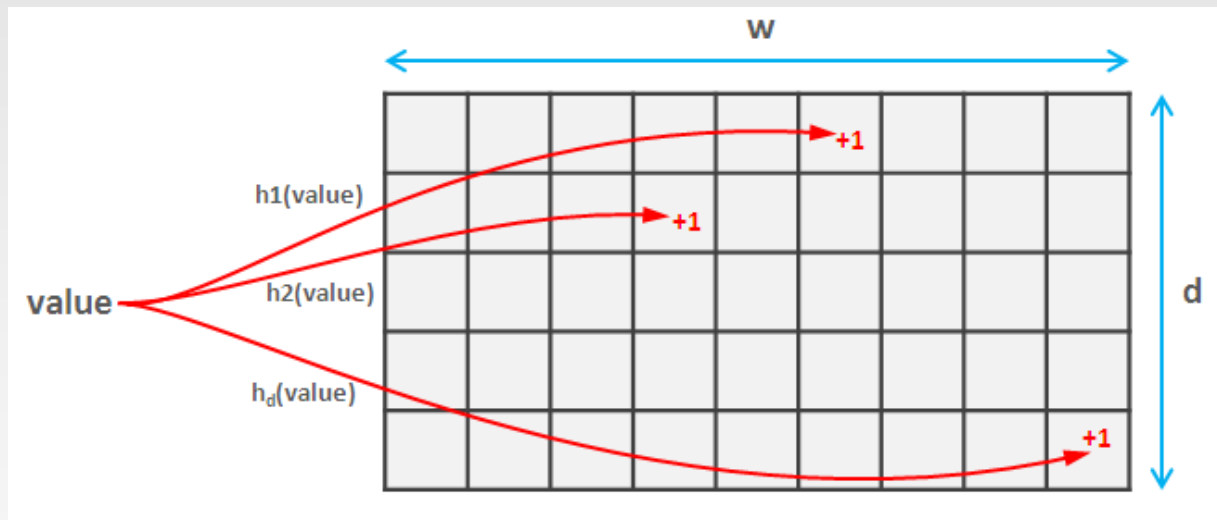
# Count-Min Sketch

- ❖ In general, model input stream as a vector  $x$  of dimension  $U$ 
  - $x[i]$  is frequency of element  $i$
- ❖ The count-min sketch has two parameters, the number of buckets  $w$  and the number of hash functions  $d$
- ❖ Creates a small summary as an array of  $w \times d$  in size
- ❖ Use  $d$  hash function to map vector entries to  $[1..w]$



# Count-Min Sketch

- ❖ The count-min-sketch supports two operations:  $\text{Inc}(x)$  and  $\text{Count}(x)$
- ❖ The operation  $\text{Count}(x)$  is supposed to return the frequency count of  $x$ , meaning the number of times that  $\text{Inc}(x)$  has been invoked in the past
- ❖ The code for  $\text{Inc}(x)$  is simply:
  - for  $i = 1, 2, \dots, d$ : increment  $\text{CMS}[i][h_i(x)]$



- ❖ The code for  $\text{Count}(x)$  is simply:
  - return  $\min_{i=1}^d \text{CMS}[i][h_i(x)]$

# References

- ❖ Chapter 4, Mining of Massive Datasets.
- ❖ [Finding Frequent Items in Data Streams](#)

**End of Chapter 6.2**