

Frequent elements

- **Heavy Hitters Problem.** Find all elements i that occurs more than m/k times for some fixed k .
- **Example.** Return all elements that occur more than $21/3$ times = 7.

4 4 1 2 4 4 3 1 1 2 5 9 7 4 1 3 4 1 4 4 1

Frequent elements

- **Heavy Hitters Problem.** Find all elements i that occurs more than m/k times for some fixed k .
- **Example.** Return all elements that occur more than $21/3$ times = 7.

4 4 1 2 4 4 3 1 1 2 5 9 7 4 1 3 4 1 4 4 1

Frequent elements

- **Heavy Hitters Problem.** Find all elements i that occurs more than m/k times for some fixed k .
- **Example.** Return all elements that occur more than $21/3$ times = 7.

4 4 1 2 4 4 3 1 1 2 5 9 7 4 1 3 4 1 4 4 1

- **Bad news.** Need $\Omega(n)$ space for one-pass algorithm.
- **Good news.**
 - Can estimate the frequency.
 - Can do better if we allow one-sided error:
 - Output all elements that occur more than m/k times.
 - Might also output other elements.

Frequent elements

- **Example.** $k = 3$.

4 4 1 2 4 4 3 1 1 2 5 9 7 4 1 3 4 1 4 4 1

counter 1: 4, 2

counter 2: 1, 1

```

Keep k-1 counters in an associative array A.
while (stream is not empty) do
  if j ∈ keys(A) then
    A[j] ← A[j] + 1
  else if |keys(A)| < k - 1 then
    A[j] ← 1
  else
    Decrement all counters by 1.
    Remove all elements with counter 0.
Output all elements in keys(A)

```

- **Space.** $O(k \cdot (\log n + \log m))$.

Misra-Gries Analysis

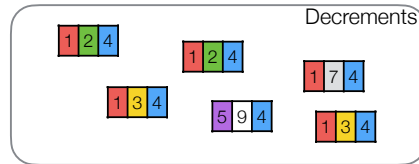
- **Lemma.** Any item with frequency more than m/k is in A by the end of the algorithm.
- **Lemma.** Let \hat{f}_i be the estimate of the frequency of element i . Then

$$f_i - \frac{m}{k} \leq \hat{f}_i \leq f_i .$$



counter 1: 4, 2

counter 2: 1, 1



Reservoir Sampling

- **Algorithm.**

```

put the first  $k$  elements into a "reservoir"  $R = \{r_1, r_2, \dots, r_k\}$ .
for  $i > k$  until the stream is empty do
    with probability  $k/i$  replace a random entry of  $R$  with  $a_i$ 
Return  $R$ .
    
```

- **Claim.** For all $t \geq i$, $P[a_i \in R_t] = k/t$, where R_t denotes the reservoir after time t .
- **Proof.** Consider element a_i .
 - $P[a_i \text{ chosen at time } i] = k/i$.
 - $P[a_i \text{ replaced at time } j] = (k/j) \cdot (1/k) = 1/j$.
 - $P[a_i \text{ not replaced at time } j] = 1 - 1/j = (j-1)/j$.
 - Thus

$$P[a_i \in R_t] = \frac{k}{i} \cdot \frac{i}{i+1} \cdot \frac{i+1}{i+2} \cdots \frac{t-1}{t} = \frac{k}{t} .$$

Reservoir Sampling