Lecture 6: The Data Stream Model

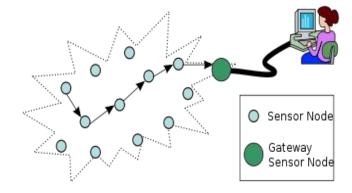
Elaine Shi

(Slides by David Woodruff)

Data Streams

 A stream is a sequence of data, that is too large to be stored in available memory

- Examples
 - Internet search logs
 - Network Traffic
 - Sensor networks



Scientific data streams (astronomical, genomics, physical simulations)...

Streaming Model 4 3 7 3 1 1 2

- Stream of elements a_1 , ..., a_i , ... each from an alphabet Σ and taking b bits to represent
- Single or small number of passes over the data

Streaming Model 4 3 7 3 1 1 2 ...

- Stream of elements $a_1,\,...,\,a_i,\,...$ each from an alphabet Σ and taking b bits to represent
- Single or small number of passes over the data
- Almost all algorithms are randomized and approximate
 - Usually necessary to achieve efficiency
 - Randomness is in the algorithm, not the input
- Goals: minimize space complexity (in bits), processing time

- Let $a_{[1:t]} = \langle a_1, ..., a_t \rangle$ be the first t elements of the stream
- Suppose $a_1, ..., a_t$ are integers in $\{-2^b+1, -2^b+2, ..., -1, 0, 1, 2, ..., 2^b-1\}$
 - Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32
- How many bits do we need to maintain $f(a_{[1:t]}) = \sum_{i=1,\dots,t} a_i$?

- Let $a_{[1:t]} = \langle a_1, ..., a_t \rangle$ be the first t elements of the stream
- Suppose $a_1, ..., a_t$ are integers in $\{-2^b+1, -2^b+2, ..., -1, 0, 1, 2, ..., 2^b-1\}$ (bit Sign Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900. 4. 32
- How many bits do we need to maintain $f(a_{[1:t]}) = \sum_{i=1,\dots,t} a_i$?
 - Outputs on example: 3, 4, 21, 25, 16, 48, 149, 152, -570, -567, 333, 337, 379, ...
- Outputs on examp
 O(b + log t)

- Let $a_{[1:t]} = \langle a_1, ..., a_t \rangle$ be the first t elements of the stream
- Suppose $a_1, ..., a_t$ are integers in $\{-2^b + 1, -2^b + 2, ..., -1, 0, 1, 2, ..., 2^b 1\}$
 - Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32
- How many bits do we need to maintain $f(a_{[1:t]}) = \sum_{i=1,...,t} a_i$?
 - Outputs on example: 3, 4, 21, 25, 16, 48, 149, 152, -570, -567, 333, 337, 379, ...
 - O(b + log t)
- How many bits do we need to maintain $f(a_{[1:t]}) = \max_{i=1,\dots,t} a_i$?



- Let $a_{[1:t]} = \langle a_1, ..., a_t \rangle$ be the first t elements of the stream
- Suppose $a_1, ..., a_t$ are integers in $\{-2^b + 1, -2^b + 2, ..., -1, 0, 1, 2, ..., 2^b 1\}$
 - Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32
- How many bits do we need to maintain $f(a_{[1:t]}) = \sum_{i=1,\dots,t} a_i$?
 - Outputs on example: 3, 4, 21, 25, 16, 48, 149, 152, -570, -567, 333, 337, 379, ...
 - O(b + log t)
- How many bits do we need to maintain $f(a_{[1:t]}) = \max_{i=1,...,t} a_i$?
 - Outputs on example: 3, 3, 17, 17, 17, 32, 101, 101, 101, 101, 900, 900, 900, ...
 - O(b) bits

- The median of all the numbers we've stored so far
 - Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32
 - Median: 3, 1, 3, 3, 3, 3, 4, 3, ...
 - This seems harder...

- The median of all the numbers we've stored so far
 - Example stream (3,1,17, 4, -9, 32, 101,3,-722, 3, 900, 4, 32
 - Median: 3, 1, 3, 3, 3, 3, 4, 3, ...
 - This seems harder...
- The number of distinct elements we've seen so far?
 - Outputs on example: 1, 2, 3, 4, 5, 6, 7, 7, 8, 8, 9, 9, 9, ...

- The median of all the numbers we've stored so far
 - Example stream: 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32
 - Median: 3, 1, 3, 3, 3, 3, 4, 3, ...
 - This seems harder...
- The number of distinct elements we've seen so far?
 - Outputs on example: 1, 2, 3, 4, 5, 6, 7, 7, 8, 8, 9, 9, 9, ...
- The elements that have appeared at least an ϵ -fraction of the time? These are the ϵ -heavy hitters
 - Cover today

Many Applications

 Internet router may want to figure out which IP connections are heavy hitters, e.g., the ones that use more than .01% of your bandwidth

• Or maybe the router wants to know the median (or 90-th percentile) of the file sizes being transferred

Hashing is a key technique

Finding €-Heavy Hitters

- S_t is the multiset of items at time t, so $S_0 = \emptyset$, $S_1 = \{a_1\}, \dots, S_i = \{a_1, \dots, a_i\}$, $count_t(e) = |\{i \in \{1, 2, \dots, t\} \text{ such that } a_i = e\}|$
- $e \in \Sigma$ is an ϵ -heavy hitter at time t)if $count_t(e) > \epsilon \cdot t$
- Given $\epsilon > 0$, can we output the ϵ -heavy hitters?
 - Let's output a set of size $\frac{1}{\epsilon}$ containing all the ϵ -heavy hitters
- Note: can output "false positives" but not allowed to output "false negatives", i.e., not allowed to miss any heavy hitter, but could output non-heavy hitters

Finding ∈-Heavy Hitters

- Example (E,D) (E,D) (E,E) (E
- At time 5, the element D is the only 1/3-heavy hitter
- At time 11, both B and D are 1/3-heavy hitters
- At time 15, there is no 1/3-heavy hitter
- At time 16, only E is a 1/3-heavy hitter

Can't afford to keep counts of all items, so how to maintain a short summary to output the ϵ -heavy hitters?

$\ell = \frac{1}{2}$ Finding a Majority Element

```
• First find a .5-heavy hitter, that is, a majority element:
 memory \ empty and counter \ 0
 when element at arrives
    if (counter == 0)
       memory \leftarrow (a_t) and counter \leftarrow (1)
    else
            memory
         counter + +
       else
         counter
         (discard at
• At end of the stream, return the element in memory
```

intially = count = 0 Memory = 3 Count = 1

Memory = 3, Count = 1

3 1 2 1 1

Memory = 3, Count = 1

Memory = 3, Count = 0

31211

Memory = 3, Count = 0

Memory = 3, Count = 1

Memory = 2, Count = 1

31211

Memory = 3, Count = 1 Memory = 3, Count = 0 Memory = 2, Count = 1 Memory = 2, Count = 0

31211

Memory = 3, Count = 0

Memory = 3, Count = 1

Memory = 2, Count = 1

Memory = 2, Count = 0

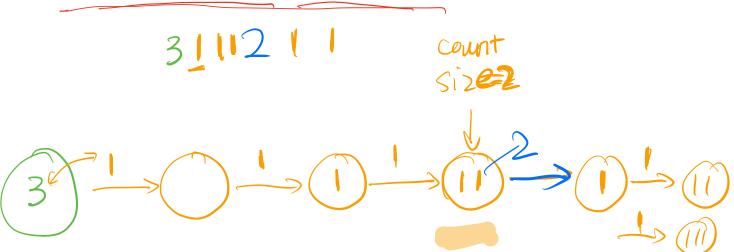
Memory # 1, Count =/

Output = (

Analysis of Finding a Majority Element

• If there is no majority element, we output a false positive, which is OK

• If there is a majority element, we will output it. Why?



Analysis of Finding a Majority Element

- If there is no majority element, we output a false positive, which is OK
- If there is a majority element, we will output it. Why?
- When we discard an element a_t, we throw away a different element
- When we throw away a copy of a majority element, we throw away another element
 - Either majority element is in memory, or majority element arrives in stream but some other item is in memory
- Majority element is more than half the total number of elements, so can't throw away all of them

Extending to ε -Heavy Hitters (ε)



Extending to ε-Heavy Hitters

Set
$$k = \left[\frac{1}{\epsilon}\right] - 1$$

Array T[1, ..., k], where each location can hold one element from Σ Array C[1, ..., k], where each location can hold a non-negative integer C[i] \leftarrow 0 and T[i] \leftarrow \bot for all i

If there is $j \in \{1, 2, ..., k\}$ such that $a_t = T[j]$, then C[j] + +Else if some counter C[j] = 0 then $T[j] \leftarrow a_t$ and $C[j] \leftarrow 1$ Else decrement all counters by 1 (and discard element a_t)

$$est_t(e) = C[j]$$
 if $e == T[j]$ for some j, and $est_t(e) = 0$ otherwise

Analyzing Counts

- Lemma: $0 \le \text{count}_t(e) \text{est}_t(e) \le \frac{t}{k+1} \le \epsilon \cdot t$
- Proof: $count_t(e) \ge est_t(e)$ since we never increase a counter for e unless we see e

If we don't increase est (e) by 1 for an update to e, we decrement k counters and discard the current update to e

- Either e is in memory and we decrement its count, or e is a stream update and we discard it

So we drop k+1 distinct stream updates, but there are t total updates, so we won't increase $est_t(e)$ by 1, when we should, at most $\frac{t}{k+1} \le \epsilon \cdot t$ times

Heavy Hitters Guarantee

• At any time t, all ϵ -heavy hitters e are in the array T. Why?

Heavy Hitters Guarantee

- At any time t, all ϵ -heavy hitters e are in the array T. Why?
- For an ϵ -heavy hitter e, we have $\operatorname{count}_{\mathsf{t}}(\mathsf{e}) \triangleright \epsilon \cdot \mathsf{t}$
- But $est_t(e) \ge count_t(e) \epsilon \cdot t$
- So $est_t(e) > 0$ so e is in array T
- Space is $O(k) \log(|\Sigma|) + \log t) = O(1/\epsilon) (\log(|\Sigma|) + \log t)$ bits

Heavy Hitters with Deletions

- Suppose we can delete elements e that have already appeared
- Example: (add, A) (add, B), (add, A), (del, B), (del, A), (add, C)
- Multisets at different times $S_0=\emptyset S_1=\{A\}, S_2=\{A,B\}, S_3=\{A,A,B\}, S_4=\{A,A\}, S_5=\{A\}, S_6=\{A,C\}, \dots$
- "active" set S_t has size $|S_t| \neq \sum_{e \in \Sigma} count_t(e)$ and can grow and shrink

- Query "What is $count_t(e)$?", should output $est_t(e)$ with: $\Pr[|est_t(e) count_t(e)| \le c |S_t|] \ge 1 \delta$
- Want space close to our previous $O(1/\epsilon)$ (log($|\Sigma|$) + log t) bits

- Query "What is $count_t(e)$?", should output $est_t(e)$ with: $Pr[|est_t(e) count_t(e)| \le \epsilon |S_t|] \ge 1 \delta$
- Want space close to our previous $O(1/\epsilon)$ (log($|\Sigma|$) + log t) bits
- Let $h: \Sigma \to \{0,1,2,...,k-1\}$ be a hash function (will specify later)
- Maintain an array A[0, 1, ..., k-1] to store non-negative integers

```
when update a_t arrives: 

if a_t = (add, e) then A[h(e)] + +

else a_t = (del, e), and A[h(e)] - -
```

•
$$\operatorname{est}_{\mathsf{t}}(e) = A[h(e)]$$

•A[h(e)] =
$$\sum_{e' \in \Sigma} \text{count}_t(e')$$
 · $\mathbf{1}(\text{h}(e') = \text{h}(e))$, where $\mathbf{1}(\text{condition})$ evaluates to 1 if the condition is true, and evaluates to 0 otherwise

• $A[h(e)] = \sum_{e' \in \Sigma} count_t(e') \cdot \mathbf{1}(h(e') = h(e))$, where $\mathbf{1}(condition)$ evaluates to 1 if the condition is true, and evaluates to 0 otherwise

•
$$A[h(e)] = count_t(e) + \sum_{e'\neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e)),$$

• $A[h(e)] = \sum_{e' \in \Sigma} count_t(e') \cdot \mathbf{1}(h(e') = h(e))$, where $\mathbf{1}(condition)$ evaluates to 1 if the condition is true, and evaluates to 0 otherwise

•
$$A[h(e)] = count_t(e) + \sum_{e' \neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e)),$$

• $est_t(e) - count_t(e) = \sum_{e' \neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e))$

• $A[h(e)] = \sum_{e' \in \Sigma} count_t(e') \cdot \mathbf{1}(h(e') = h(e))$, where $\mathbf{1}(condition)$ evaluates to 1 if the condition is true, and evaluates to 0 otherwise

•
$$A[h(e)] = count_t(e) + \sum_{e'\neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e)),$$

•
$$\operatorname{est}_{\mathsf{t}}(\mathsf{e}) - \operatorname{count}_{\mathsf{t}}(\mathsf{e}) = \sum_{\mathsf{e}' \neq \mathsf{e}} \operatorname{count}_{\mathsf{t}}(\mathsf{e}') \cdot \mathbf{1}(\mathsf{h}(\mathsf{e}') = \mathsf{h}(\mathsf{e}))$$

• Since we have a small array A with k locations, there are likely many $e' \neq e$ with h(e') = h(e), but can we bound the expected error?

- Recall: Family H of hash functions h(U)> $\{0, 1, ..., k-1\}$ is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$
- There is a simple family where h can be specified using $O(\log |U|)$ bits. Here, $|U| = |\Sigma|$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

•
$$E[est_t(e) - count_t(e)] = E[\sum_{e'\neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e))]$$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

•
$$E[est_t(e) - count_t(e)] = E[\sum_{e' \neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e))]$$

= $\sum_{e' \neq e} count_t(e') \cdot E[\mathbf{1}(h(e') = h(e))]$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

•
$$E[est_t(e) - count_t(e)] = E[\sum_{e' \neq e} count_t(e') \cdot \mathbf{1}(h(e') = h(e))]$$

$$= \sum_{e' \neq e} count_t(e') \cdot E[\mathbf{1}(h(e') = h(e))]$$

$$= \sum_{e' \neq e} count_t(e') \cdot Pr[h(e') = h(e)]$$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

$$\begin{split} \bullet \ \mathsf{E}[\mathsf{est}_\mathsf{t}(\mathsf{e}) - \mathsf{count}_\mathsf{t}(\mathsf{e})] &= \mathsf{E}[\sum_{\mathsf{e}' \neq \mathsf{e}} \mathsf{count}_\mathsf{t}(\mathsf{e}') \cdot \mathbf{1}(\mathsf{h}(\mathsf{e}') = \mathsf{h}(\mathsf{e}))] \\ &= \sum_{\mathsf{e}' \neq \mathsf{e}} \mathsf{count}_\mathsf{t}(\mathsf{e}') \cdot \mathsf{E}[\ \mathbf{1}(\mathsf{h}(\mathsf{e}') = \mathsf{h}(\mathsf{e}))] \\ &= \sum_{\mathsf{e}' \neq \mathsf{e}} \mathsf{count}_\mathsf{t}(\mathsf{e}') \cdot \underbrace{\mathsf{Pr}[\mathsf{h}(\mathsf{e}') = \mathsf{h}(\mathsf{e})]} \\ &\leq \sum_{\mathsf{e}' \neq \mathsf{e}} \mathsf{count}_\mathsf{t}(\mathsf{e}') \cdot \underbrace{\left(\frac{1}{\mathsf{k}}\right)} \end{aligned}$$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

$$\begin{split} \bullet \; \mathsf{E}[\mathsf{est}_\mathsf{t}(e) - \mathsf{count}_\mathsf{t}(e)] &= \mathsf{E}[\sum_{e' \neq e} \mathsf{count}_\mathsf{t}(e') \cdot \mathbf{1}(\mathsf{h}(e') = \mathsf{h}(e))] \\ &= \sum_{e' \neq e} \mathsf{count}_\mathsf{t}(e') \cdot \mathsf{E}[\; \mathbf{1}(\mathsf{h}(e') = \mathsf{h}(e))] \\ &= \sum_{e' \neq e} \mathsf{count}_\mathsf{t}(e') \cdot \mathsf{Pr}[\mathsf{h}(e') = \mathsf{h}(e)] \\ &\leq \sum_{e' \neq e} \mathsf{count}_\mathsf{t}(e') \cdot \left(\frac{1}{\mathsf{k}}\right) \\ &= \underbrace{\left(S_\mathsf{t}\right) \cdot \mathsf{count}_\mathsf{t}(e)}_{\mathsf{k}} \leq \underbrace{\left(S_\mathsf{t}\right)}_{\mathsf{k}} \end{split}$$

• Recall: Family H of hash functions h: U -> {0, 1, ..., k-1} is universal if for all $x \neq y$, $\Pr_{h \leftarrow H}[h(x) = h(y)] \leq \frac{1}{k}$

• There is a simple family where h can be specified using O(log |U|) bits. Here, $|U| = |\Sigma|$

•
$$\mathbf{E}[\operatorname{est}_{t}(e) - \operatorname{count}_{t}(e)] = E[\sum_{e' \neq e} \operatorname{count}_{t}(e') \cdot \mathbf{1}(h(e') = h(e))]$$

$$= \sum_{e' \neq e} \operatorname{count}_{t}(e') \cdot E[\mathbf{1}(h(e') = h(e))]$$

$$= \sum_{e' \neq e} \operatorname{count}_{t}(e') \cdot \Pr[h(e') = h(e)]$$

$$\leq \sum_{e' \neq e} \operatorname{count}_{t}(e') \cdot \left(\frac{1}{k}\right)$$

$$= \frac{|S_{t}| - \operatorname{count}_{t}(e)}{k} \leq \frac{|S_{t}|}{k} \leq C$$

 $k = 1/\epsilon$ makes this at most $\epsilon \cdot |S_t|$. Space is $O(\frac{1}{\epsilon})$ counters plus storing hash function

High Probability Bounds for CountMin

- Have $0 \le \operatorname{est}_{\mathsf{t}}(\mathsf{e}) \operatorname{count}_{\mathsf{t}}(\mathsf{e}) \le (|S_{\mathsf{t}}|/k)$ n expectation from CountMin
 - With probability at least 1/2, $est_t(e) count_t(e) \le 2|S_t|/k$ Why?



High Probability Bounds for CountMin

- Have $0 \le \operatorname{est}_{\mathsf{t}}(\mathsf{e}) \operatorname{count}_{\mathsf{t}}(\mathsf{e}) \le |\mathsf{S}_{\mathsf{t}}|/k$ in expectation from CountMin
 - With probability at least 1/2, $est_t(e) count_t(e) \le 2|S_t|/k$ Why?

• Can we make the success probability 1- δ ?

High Probability Bounds for CountMin

- Have $0 \le \operatorname{est}_{\mathsf{t}}(\mathsf{e}) \operatorname{count}_{\mathsf{t}}(\mathsf{e}) \le |\mathsf{S}_{\mathsf{t}}|/k$ in expectation from CountMin
 - With probability at least 1/2, $est_t(e) count_t(e) \le 2|S_t|/k$ Why?
- Can we make the success probability 1-δ? Repeat in time

 the original scheme
 - Independent repetition: pick m hash functions $h_1, ..., h_m$ with $h_i : \Sigma \to \{0, 1, 2, ..., k-1\} \text{ independently from H. Create array } A_i \text{ for } h_i$ when update a_t arrives:

for each i from 1 to m

if
$$a_t = (add, e)$$
 then $A_i[h_i(e)] + +$
else $a_t = (del, e)$ and $A_i[h_i(e)] - -$

$$\mathtt{best}_t(e) := \min_{i=1}^m A_i[h_i(e)].$$

What is our new estimate of $count_t(e)$?

$$\underbrace{\left(\text{best}_t(e) \right)}_{i=1} := \min_{i=1}^m A_i[h_i(e)].$$

Each A_i[h_i(e)] is an overestimate to count_t(e)

$$\underline{\mathtt{best}_t(e)} := \min_{i=1}^m A_i[h_i(e)].$$

- Each A_i[h_i(e)] is an overestimate to count_t(e)
- By independence, $\Pr[\text{for all i}, A_i[h_i(e)] \text{count}_t(e) \ge 2|S_t|/k) \le 6$





$$\mathsf{best}_t(e) := \min_{i=1}^m A_i[h_i(e)].$$

- Each $A_i[\Pi_i(e)]$ is an overestimate to $count_t(e)$ $2 \in S_t$ By independence, $Pr[for\ all\ i,\ A_i[h_i(e)] count_t(e) \ge 2|S_t|/k] \le \left(\frac{1}{2}\right)^m$
- For $k = \frac{2}{\epsilon}$ and $m \neq \log_2(\frac{1}{\delta})$, the error is at most $\epsilon |S_t|$ with probability 1- δ

What is our new estimate of count_t(e)

$$\mathtt{best}_t(e) := \min_{i=1}^m A_i[h_i(e)].$$

- Each A_i[h_i(e)] is an overestimate to count_t(e)
- By independence, $\Pr[\text{for all i, } A_i[h_i(e)] \text{count}_t(e) \ge 2|S_t|/k] \le \left(\frac{1}{2}\right)^m$
- For $k = \frac{2}{\epsilon}$ and $m = \log_2(\frac{1}{\delta})$, the error is at most $\epsilon |S_t|$ with probability 1- δ
- Space: m) $k = O(\frac{\log(\frac{1}{\delta})}{\epsilon})$ counters each of $O(\lg t)$ bits m) $O(\log |\Sigma|) = O(\log(\frac{1}{\delta})\log|\Sigma|)$ bits to store hash functions

For fixed e. Pr[large emor] < o E-Heavy Hitters Pr[alle have small error]
Pr[ae, has large error] = []. S

Our new estimate best_t(e) satisfies

$$\Pr[|\mathsf{best}_t(e) - \mathsf{count}_t(e)| \leq \epsilon |S_t|] \geq 1 - \delta$$
 and uses
$$O(\frac{\log(\frac{1}{\delta})\log t}{\epsilon} + \log(\frac{1}{\delta})\log |\Sigma|) \text{ bits of space}$$

• What if we want with probability 9/10, simultaneously for all e, $|\text{best}_{t}(e) - \text{count}_{t}(e)| \le \epsilon |S_{t}|$?

Sofar, the analysis is for an arbitrary fixed e

€-Heavy Hitters

• Our new estimate best_t(e) satisfies

$$\Pr[|\mathsf{best}_{\mathsf{t}}(\mathsf{e}) - \mathsf{count}_{\mathsf{t}}(\mathsf{e})| \le \epsilon |\mathsf{S}_{\mathsf{t}}|] \ge 1 - \delta$$

and uses $O(\frac{\log(\frac{1}{\delta})\log t}{\log(\frac{1}{\delta})\log|\Sigma|})$ bits of space

$$\log(5)$$

$$= \log(10.1)$$

- What if we want with probability 9/10, simultaneously for all e, $|\text{best}_t(e) \text{count}_t(e)| \le \epsilon |S_t|$?
- Set $\delta = \frac{1}{10(\Sigma)}$ and apply a union bound over all $e \in \Sigma$