

Today:

- Wrap up LSH (see previous notes)
- New theme: representative subsets / coresets

Will cover:

- quantiles / median ← Today
- clustering (k-median)

## Representative subsets / coresets

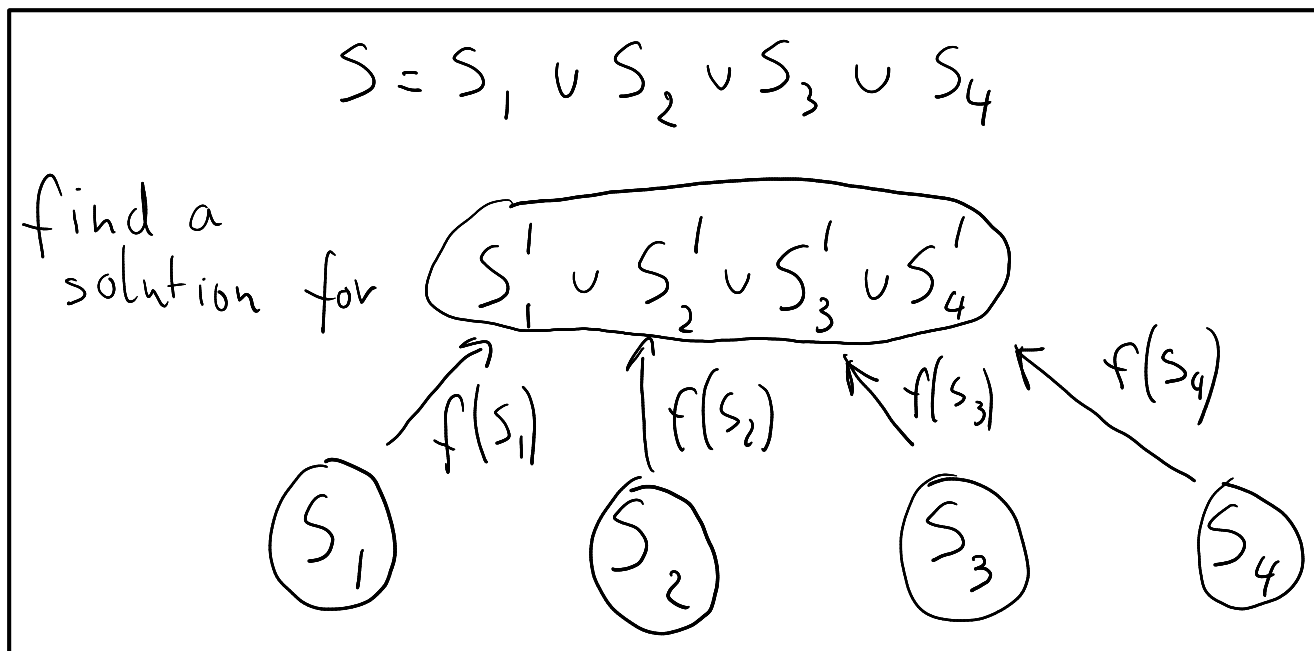
General goal:

- retain a small subset
- sufficient for solving the original problem(s)

Nice bonus properties (sometimes):

- replacing any subset with its coreset still leads to good approximation

- can combine coresets for separate subsets and still find a good solution for the original problem



Exact median: have to keep everything

$$S = S_1 \cup S_2$$

$\uparrow$              $\uparrow$   
 size  $n$     size  $n-1$

$$S_1 = \{s_1 < s_2 < s_3 < \dots < s_n\}$$

$n-i$  elements  $<$  coreset( $S_1$ )  $<$   $i-1$  elements

$S_2$  ← median of  $S = i$ -th element of  $S_1$

$S_2$  can select any element of  $S_1$

Instead:  $\epsilon$ -approximate median

output an element of  $S$  that is:

- greater than or equal to at least  $\frac{n}{2}(1-\epsilon)$  elements

and

- less than or equal to at least  $\frac{n}{2}(1+\epsilon)$  elements

---

Notation

multiset  $S \overset{\text{finite}}{\subseteq} \mathbb{R}$

$$\text{rank}_S(x) = |\{y \in S : y \leq x\}|$$

= number of elements in  $S$

that are at most  $x$

Fact:  $\text{median}(S) = \text{smallest } x \text{ s.t. } \text{rank}_S(x) \geq \frac{|S|}{2}$

| 13 - 3 |

# Coresets for quantiles ( $S, S'$ are multisets)

$S'$  is an  $\epsilon$ -coreset for  $S$  if

- same number of elements
- every element  $x \in S'$  is present in  $S$   
(possibly fewer copies of  $x$  in  $S$ )
- for all  $x \in \mathbb{R}$

$$\text{rank}_S(x) \leq \text{rank}_{S'}(x) \leq \text{rank}_S(x) + \epsilon |S|$$

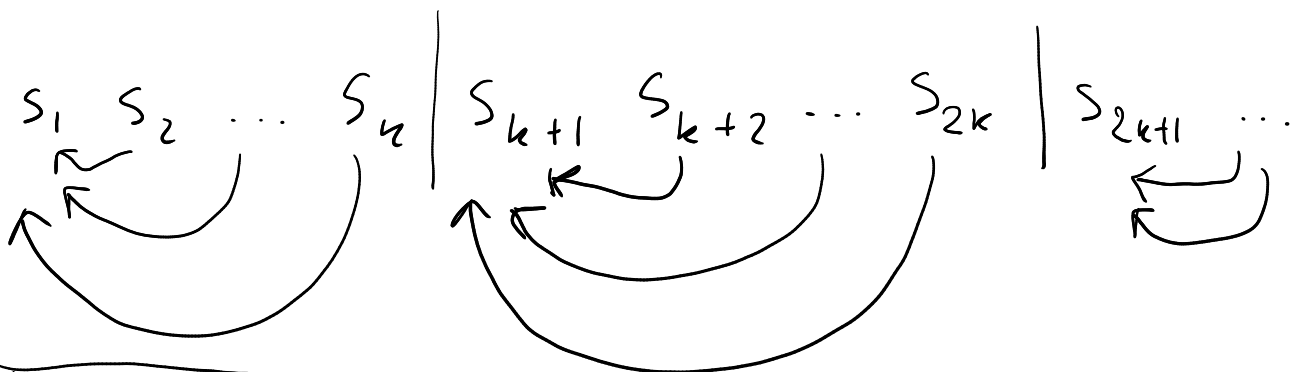
## Our construction

$$S = \{s_1 \leq s_2 \leq s_3 \leq \dots \leq s_n\}$$

$$\text{set } k = \lceil \epsilon n \rceil$$

Replace  $s_1, \dots, s_k$  with  $k$  copies of  $s_1$   
 $s_{k+1}, \dots, s_{2k}$  with  $k$  copies of  $s_{k+1}$

(last block might be shorter)



Output: multiset with at most  $1/\epsilon$  different elements

Storage:  $\leq 1/\epsilon \times (\text{element of } S + \text{count})$

---

### Why this works

$S'$  = output of our construction for  $S$

- first two properties of a coresets met trivially

- consider any  $x \in \mathbb{R}$ :

- each block  $\{s_1, \dots, s_k\}, \{s_{k+1}, \dots, s_{2k}\}, \dots$

spans a range  $[a_1, b_1], [a_2, b_2], \dots$   
 $\begin{array}{cccc} \parallel & \parallel & \parallel & \parallel \\ s_1 & s_k & s_{k+1} & s_{2k} \end{array}$

-  $a_1 \leq b_1 \leq a_2 \leq b_2 \leq a_3 \leq \dots$

- if  $x$  not in any range

$$\text{rank}_S(x) = \text{rank}_{S'}(x) \quad \checkmark$$

- otherwise let  $[a_j, b_j]$  be the last range to which  $x$  belongs

- elements move only within ranges and only decrease
- ranges  $[a_i, b_i]$  for  $i < j$  are such that  $b_i \leq x \Rightarrow$  no impact on rank of  $x$
- ranges  $[a_i, b_i]$  for  $i > j$  are such that  $x < a_i \Rightarrow$  no impact on rank of  $x$
- range  $[a_j, b_j] \ni x$ :  
at most  $k - 1 \leq \epsilon/|S|$  elements can decrease, and hence rank of  $x$  can grow by at most  $\epsilon/|S|$

$$\text{rank}_S(x) \leq \text{rank}_{S'}(x) \leq \text{rank}_S(x) + \epsilon/|S|$$


---

## Approximate median from a coresets

How: output the median, of an  $\epsilon$ -coreset  
 $s_*$

Why:  $S$ -original set,  $S'$ -coreset,  $s_* \in S$

①  $s_*$  less than or equal to  $|S|/2$  elements of  $S$

$\boxed{13 - 6}$

because  $s_* \leq \text{median}(S)$

$$\begin{aligned}
 (2) \quad & \text{rank}_{S'}(s_*) \leq \text{rank}_S(s_*) + \varepsilon |S| \\
 & s_* \text{ median of } S' \Rightarrow \text{rank}_{S'}(s_*) \geq \frac{|S'|}{2} = \frac{|S|}{2} \\
 & \boxed{\text{rank}_S(s_*) \geq \frac{|S|}{2} (1 - 2\varepsilon)}
 \end{aligned}$$

$s_*$  is an  $(2\varepsilon)$ -approximate median of  $S$

Nice property (1):

$$S = S_1 \cup S_2 \cup \dots \cup S_k$$

Union of  $\varepsilon$ -coresets for  $S_i$ 's is an  $\varepsilon$ -coreset for  $S$

Why:  $S_i$ -coreset for  $S_i$ ,  $S' = \bigcup_{i=1}^k S_i$   
 - first two properties of coresets met trivially

- consider any  $x \in R$

$$\text{For any } i \in [k]: \text{rank}_{S_i}(x) \leq \text{rank}_{S'_i}(x)$$

$$\underbrace{\sum_{i=1}^k \text{rank}_{S_i}(x)}_{\text{rank}_S(x)} \leq \underbrace{\sum_{i=1}^k \text{rank}_{S'_i}(x)}_{\text{rank}_{S'}(x)} \leq \underbrace{\sum_{i=1}^k (\text{rank}_{S_i}(x) + \varepsilon |S_i|)}_{\text{rank}_S(x) + \varepsilon |S|}$$

13 - 7

## Nice property (2)

$\varepsilon_1$ -coreset of  $\varepsilon_2$ -coreset is  $(\varepsilon_1 + \varepsilon_2)$ -coreset

Why:

$S_1$  is  $\varepsilon_2$ -coreset for  $S$

$S_2$  is  $\varepsilon_1$ -coreset for  $S_1$

- two first properties are trivial

- need to show two inequalities for any  $x \in \mathbb{R}$ :

$$\textcircled{1} \quad \text{rank}_S(x) \leq \text{rank}_{S_1}(x) \leq \text{rank}_{S_2}(x) \quad \checkmark$$

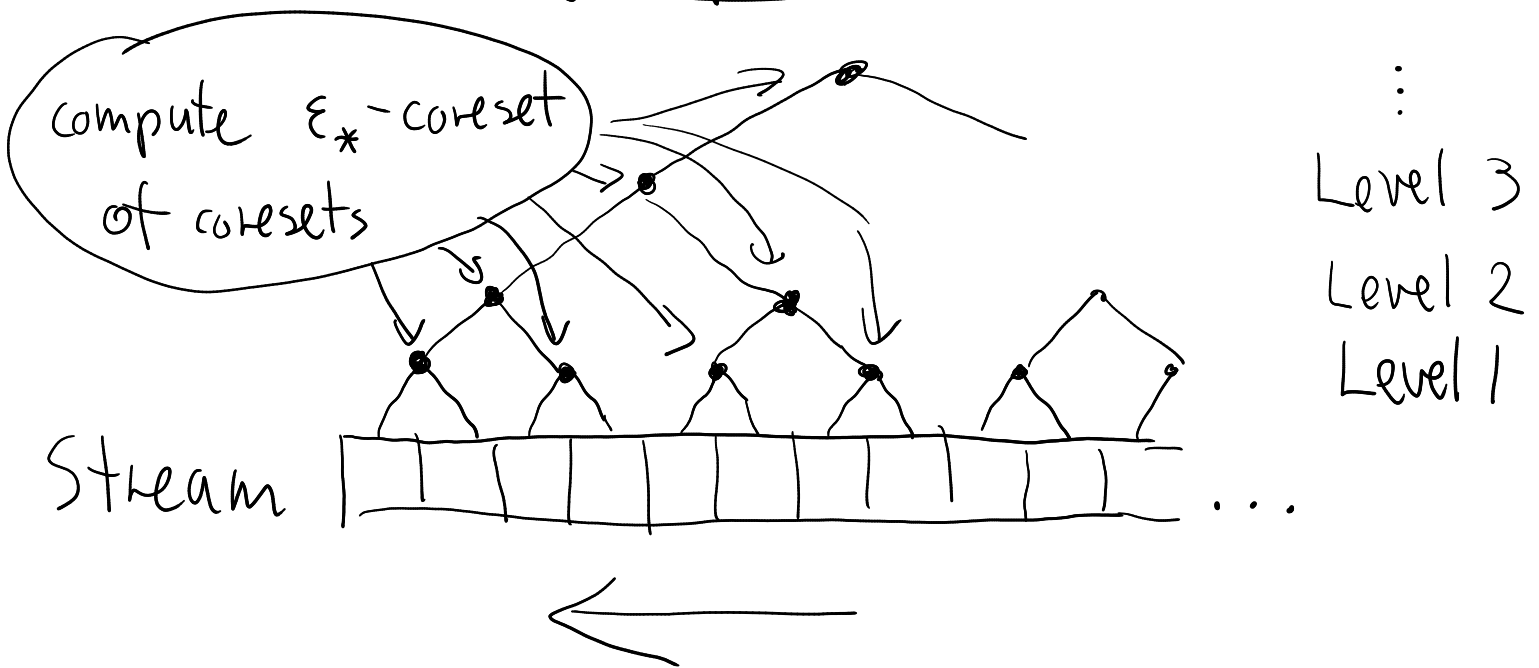
$$\textcircled{2} \quad \text{rank}_{S_2}(x) \leq \text{rank}_{S_1}(x) + \varepsilon_1 |S_1|$$

$$\leq \text{rank}_S(x) + \varepsilon_2 |S| + \varepsilon_1 |S|$$

$$= \text{rank}_S(x) + (\varepsilon_1 + \varepsilon_2) |S| \quad \checkmark$$



# Streaming application



- coresets computed on level  $i$  are  $(\epsilon_* \cdot i)$ -coresets
- $\log n$  levels
- to get  $\epsilon$ -approximate median, set  $\epsilon_* = \frac{\epsilon}{2 \log n}$

Storage: # levels  $\times$  (coreset size) =  
 $\log n \cdot \frac{2 \log n}{\epsilon} \times (\text{element} + \text{counter})$

$$O\left(\frac{\log^2 n}{\epsilon}\right) \text{ space}$$

## Alternate solution:

- Sample  $O(1/\epsilon^2)$  elements
- Output their median,  
which is  $\epsilon$ -approximate median with  
good probability
- Space:  $O(1/\epsilon^2)$

## Combine the two approaches:

- approximate median of the sample  
is still likely to be an approximate  
median of the original set  
(see Homework 2)
- Feed the sample into our algorithm
- Output approximate median of the sample
- Space:  $O\left(\frac{\log^2(1/\epsilon)}{\epsilon}\right)$