

APPROXIMATION ALGORITHMS

RANDOM SAMPLING & PRIORITY SAMPLING

RASMUS PAGH

UNIVERSITY OF COPENHAGEN



TODAY

- WARM-UP: POLLS
- RESERVOIR SAMPLING
- TAIL INEQUALITIES AND ESTIMATORS
- WEIGHTED RANDOM SAMPLING
- PRIORITY SAMPLING

POLLS

- CASE 1: - 2000 CALLS ARE MADE TO RANDOM PHONE NOS
- 1000 PEOPLE RESPOND
 - 60% SAY THEY WILL VOTE FOR CURRENT GOVT.
 - 40% SAY THEY WILL VOTE FOR OPPOSITION

HOW MANY WILL VOTE FOR GOVT. IN POPULATION AT LARGE?

- CASE 2: A NEWSPAPER CALLS 100 COMPANIES AND ASK ABOUT THE PROBABILITY THEY WILL HIRE IN 2021
- PROBABILITIES p_1, \dots, p_{100} ARE REPORTED

WHAT FRACTION OF COMPANIES DO WE EXPECT WILL HIRE?

RE. CASE 1 NEED TO KNOW RESPONSE RATE
 (SAMPLING PROBABILITY) FOR THE SETS
 $S_{\text{Govt.}}$ OF GOVT. VOTERS AND S_{Opp} OF OPP. VOTERS

		RESPONSE	
		YES	NO
VOTE	GOVT.	30% 30%	20% 30%
	OPP.	20% 20%	30% 20%

RANDOM VOTER X
 SET A: VOTERS WHO ANSWER
 $\Pr[X \in A | X \in S_{\text{Govt.}}] = \Pr[X \in A \cap S_{\text{Govt.}}] / \Pr[X \in S_{\text{Govt.}}]$

RATIO BETWEEN OBS. AND ANSWER

CAN BE OBSERVED

WHAT WE WANT

CAN DISTINGUISH IF WE KNOW HOW TO WEIGH EACH ANSWER

POSSIBLE DISTRIBUTION
 ANOTHER POSSIBLE DISTRIBUTION

RE. CASE 2

$$X_i = \begin{cases} 1 & \text{if COMPANY } i \text{ HIRES} \\ 0 & \text{OTHERWISE} \end{cases}$$

$$Y_i = \begin{cases} 1 & \text{if COMPANY} \\ & i \text{ IS SAMPLED} \\ 0 & \text{OTHERWISE} \end{cases}$$

$$E[X_i] = \text{Pr}[\text{COMPANY } i \text{ HIRES}] = p_i$$

$$\begin{aligned} E[\text{\# HIRES AMONG 100}] &= E\left[\sum_{i=1}^{100} X_i\right] \\ &= \sum_{i=1}^{100} E[X_i] \\ &= \sum_{i=1}^{100} p_i \end{aligned}$$

SUPPOSE n COMPANIES
IN TOTAL:

$$\begin{aligned} E[\text{TOTAL \# HIRES}] &= E\left[\sum_{i=1}^n X_i\right] \\ &= \sum_{i=1}^n \text{Pr}[X_i Y_i = 1] / \text{Pr}[Y_i = 1] \\ &= \left(\sum_{i=1}^{100} p_i\right) / \left(\frac{100}{n}\right) = \frac{n}{100} \sum_{i=1}^{100} p_i \end{aligned}$$

FISCHER-YATES SHUFFLE

RESERVOIR SAMPLING

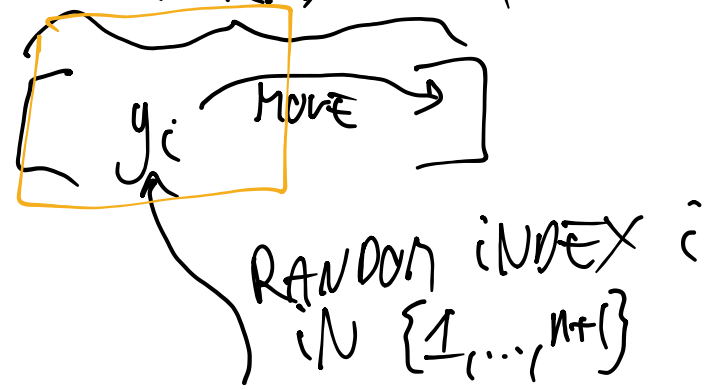
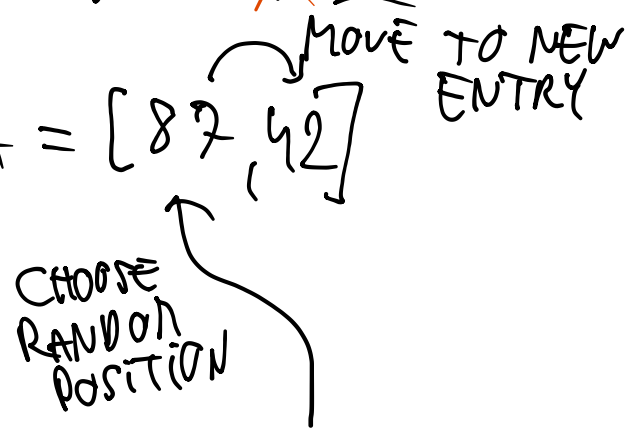
MAINTAIN: - COUNTER n
- RANDOM PERMUTATION OF n ELEMENTS

SPACE S

UNDER INSERTION OF NEW ELEMENTS

INITIALLY $n = \cancel{2}$
ARRAY $A = [87, 42]$

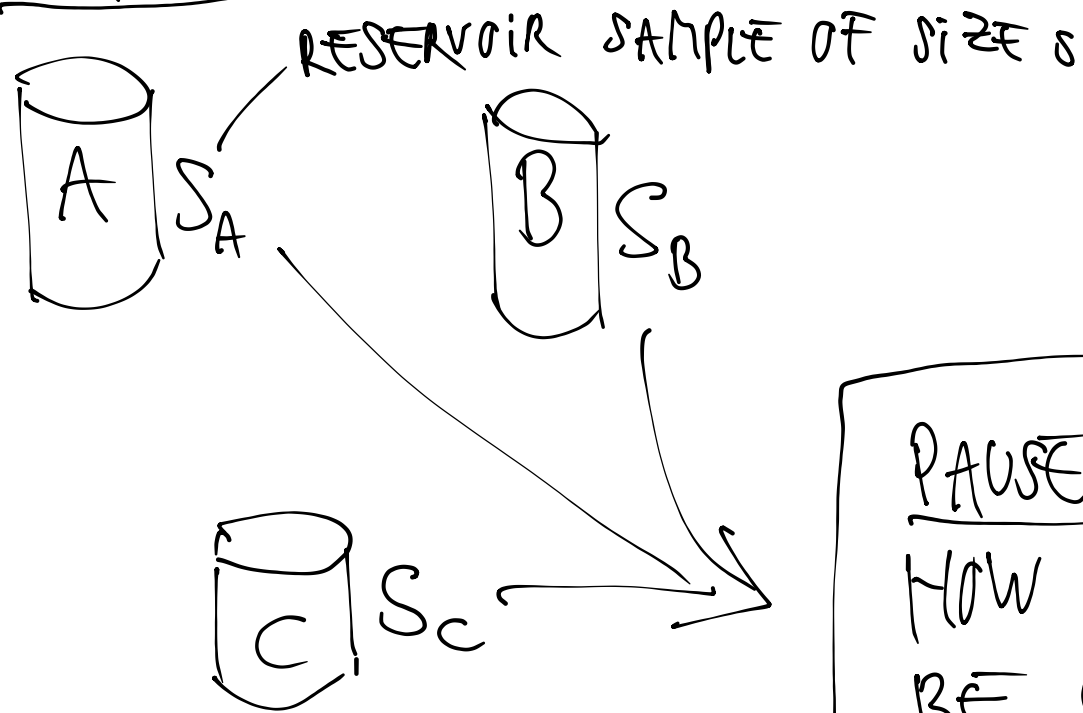
KEEP ONLY FIRST S ENTRIES
IN SHUFFLE.
EXTEND TO $n+1$



INSERTIONS: 42, 87, x_{n+1}

INVARIANT: A IS A RANDOM PERMUTATION OF x_1, \dots, x_n

SAMPLING DISTRIBUTED DATA



PAUSE AND THINK:
HOW CAN S_A, S_B, S_C
BE COMBINED INTO
A RESERVOIR SAMPLE
OF $A \cup B \cup C$?

size s

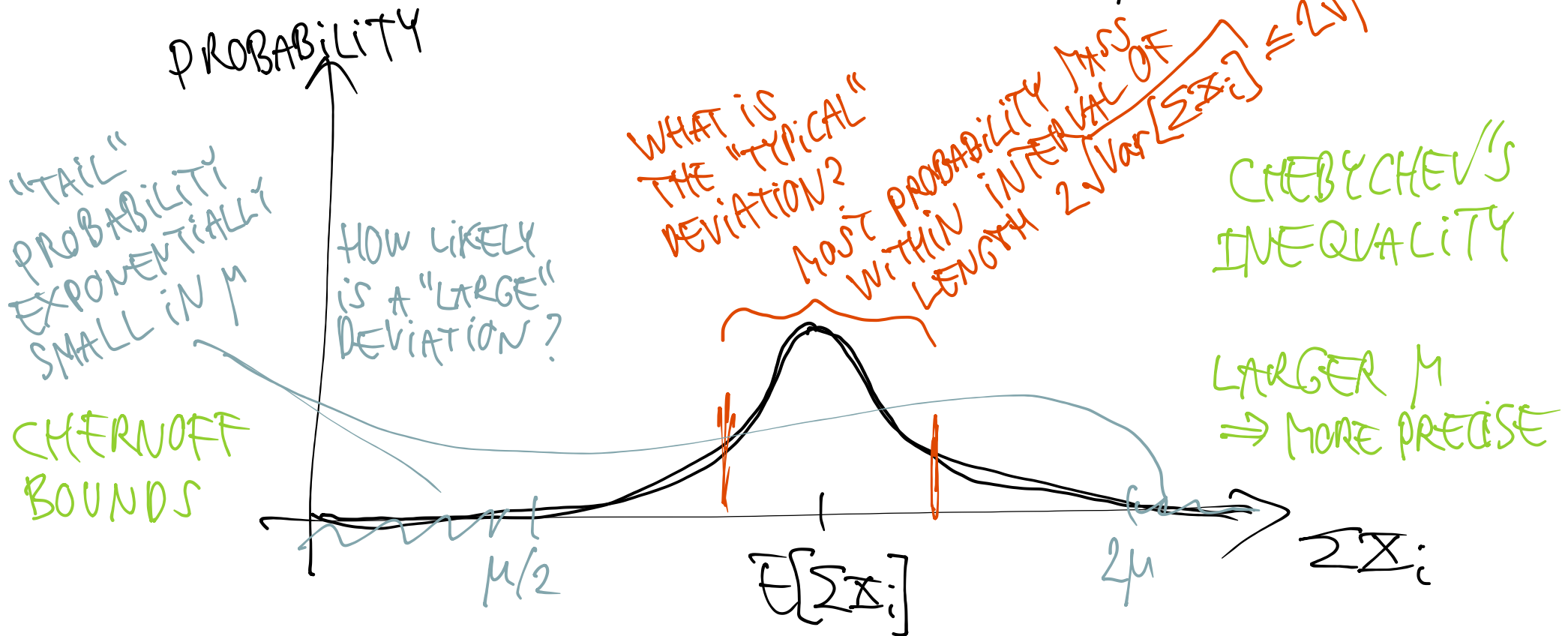
COMPUTERS
WITH DATA SETS
 A, B, C . (DISJOINT)

$$\text{WANT } \Pr[x \in S_{A \cup B \cup C}] = \frac{s}{|A \cup B \cup C|}.$$

TAIL INEQUALITIES INDEPENDENT

SETTING: HAVE RANDOM $X_1, \dots, X_n \in \{0, 1\}$, WANT TO
KNOW $\mu = \mathbb{E}[\sum_i X_i]$

HOW CLOSE IS THE OBSERVED SUM TO μ ?



TAIL BOUNDS & APPROXIMATION

CHEBYCHEV'S INEQUALITY:

$$\Pr[|X - \mathbb{E}[X]| > k] \leq \text{Var}(X)/k^2$$

(EXAMPLE)
APPROXIMATION GUARANTEE:

IF $\text{Var}(X) \leq \mathbb{E}[X]$, X IS BETWEEN $\frac{1}{2}\mathbb{E}[X]$ AND $2\mathbb{E}[X]$ WITH PROBABILITY $1 - 4/\mathbb{E}[X]$.

CHERNOFF BOUND:

IF $X = \sum_{i=1}^n X_i$, FOR INDEP. $X_i \in \{0, 1\}$

$$\Pr[X \leq (1-\epsilon)\mathbb{E}[X]] \leq \exp\left(-\frac{\epsilon^2}{2}\mathbb{E}[X]\right)$$

$$\Pr[X \geq (1+\epsilon)\mathbb{E}[X]] \leq \exp\left(-\frac{\epsilon^2}{4}\mathbb{E}[X]\right)$$

APPROXIMATION GUARANTEE:

IF $\mathbb{E}[X] \geq \frac{4 \ln(2/\delta)}{\epsilon^2}$ THEN

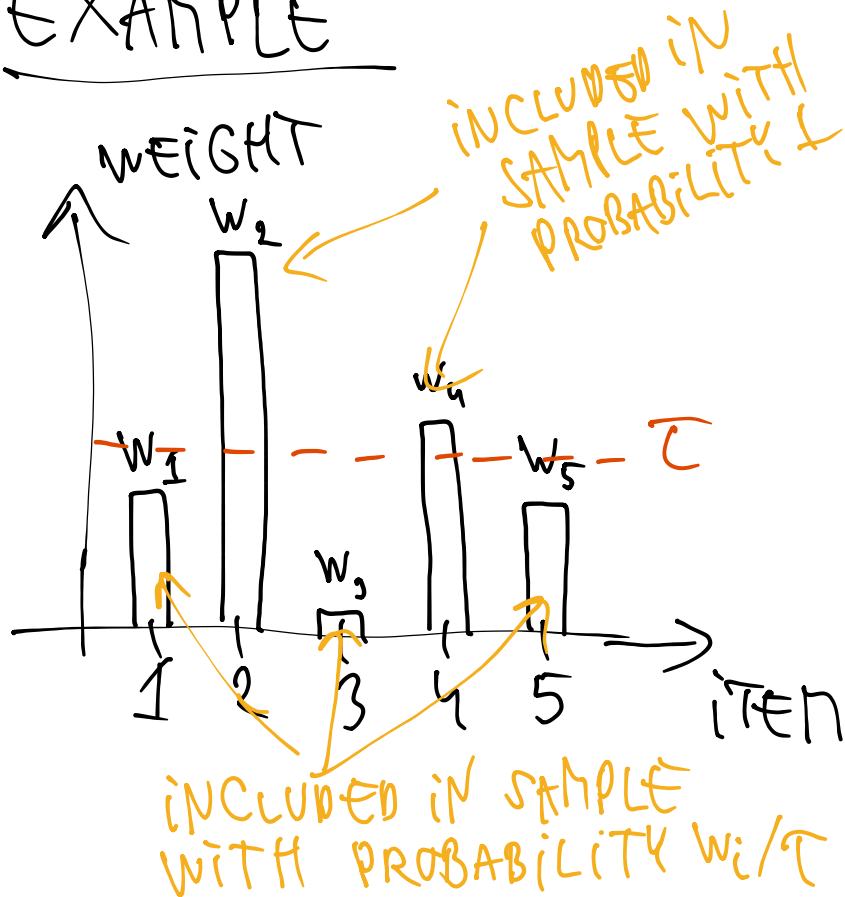
$$X \in [(1-\epsilon)\mathbb{E}[X], (1+\epsilon)\mathbb{E}[X]]$$

WITH PROBABILITY $\geq 1 - \delta$

WEIGHTED RANDOM SAMPLING

- ITEM i COMES WITH WEIGHT $w_i \geq 0$
- WANT TO SAMPLE HIGHER WEIGHT ITEMS MORE FREQUENTLY

EXAMPLE



$S=1$: SAMPLE i WITH PROB. $\frac{w_i}{\sum_{j=1}^n w_j}$

PROPORTIONAL TO WEIGHT

NORMALIZING FACTOR

$S=5$: SAMPLE ALL!
 $\tau=0$

GENERAL S :

FIND $\tau \geq 0$ SUCH THAT

$$\sum_{i=1}^n \min(1, w_i/\tau) = S$$

SAMPLE i WITH PROB. $p_i = \min(1, w_i/\tau)$

ESTIMATORS FROM WEIGHTED SAMPLES

INTERESTED IN $\mu = \sum_{i \in Q} w_i$

SET NOT KNOWN BEFORE SAMPLING

EFFICIENT UPDATES,
MEANING: SEE BOOK
(COMPLICATED)

CAN COMPUTE $\sum_{i \in S_n Q} w_i$ FROM SAMPLE.

HORVITZ-THOMPSON ESTIMATOR

UNBIASED ESTIMATOR $\hat{M} = \sum_{i \in S_n Q} w_i / p_i = \sum_{i \in Q} \mathbb{1}(i \in S) \cdot w_i / p_i$

$$E[\hat{M}] = \sum_{i \in Q} \Pr[i \in S] w_i / p_i = \sum_{i \in Q} p_i w_i / p_i = \mu$$

FACT: WEIGHTED SAMPLING PROPORTIONAL TO WEIGHT
GIVES THE SMALLEST POSSIBLE VARIANCE
AMONG ALL SAMPLING METHODS USING SPACES

PRIORITY SAMPLING

A SIMPLER WAY OF DOING WEIGHTED SAMPLING

PRIORITY OF ITEM x IS COMPUTED AS FOLLOWS:

- SAMPLE $\alpha_x \in (0, 1)$ UNIFORMLY

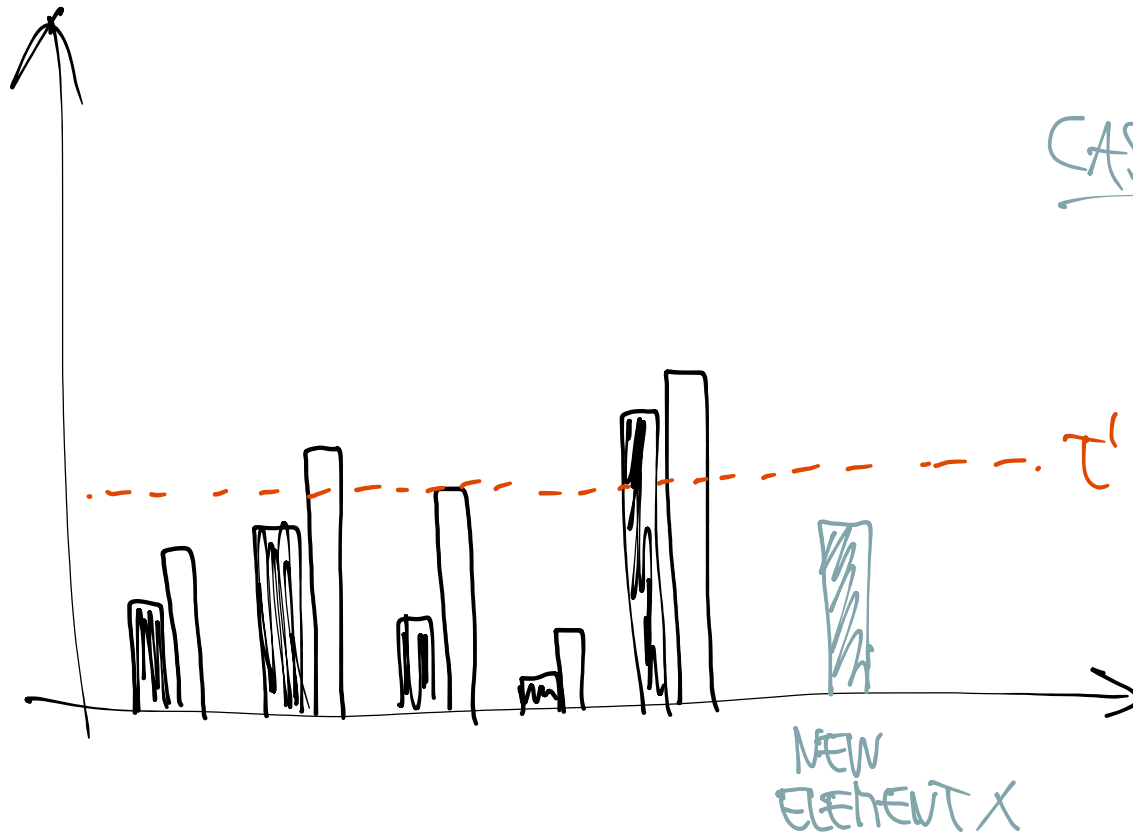
- LET $q_x = w_x / \alpha_x$

SAMPLE: ITEMS WITH s HIGHEST PRIORITIES } MAINTAIN
THRESHOLD τ EQUALS $(s+1)$ ST HIGHEST PRIORITY. } IN A
PRIORITY
QUEUE

ESTIMATOR: $\hat{\mu} = \sum_{x \in Q_{s+1}} \max(w_x, \tau)$

$$E[\hat{\mu}] = \sum_{x \in Q} w_x \leftarrow E[\max(w_x, \tau) \mid \text{sth HIGHEST PRIORITY FOR ELEM. } y \neq x \text{ IS } \tau] \\ = w_x$$

ILLUSTRATION



 WEIGHT

 PRIORITY (\geq WEIGHT)

CASE 1: $W_x \geq \tau'$: IN SAMPLE WITH PROB. 1

CASE 2: $W_x < \tau'$: IN SAMPLE IFF $W_x/\alpha_x > \tau'$
↑
PROBABILITY W_x/τ

PRIORITY SAMPLING IS DUE TO DUFFIELD, LUND & THORUP
KNOWN TO BE OPTIMAL UP TO A SINGLE SAMPLE.