

9 Mar 2026

Reservoir Sampling

Quantile Estimation

Glivenko - Cantelli Theorem

Announcements. (1) Pset 1 & Quiz 4 grades released  
Regrades available for 2 weeks (PSet)  
1 week (quiz)

(2) Typos corrected on PSet 2.  
See Ed post.

Suzzle. You're going to see a stream of  $m$  tokens, each in  $\{0, 1\}^b$ .

You only have space to store  
one token and one counter in range  
 $0, 1, \dots, m$ .

How to output a uniformly random  
element of the token sequence?

Idea. Store one token,  $y$ , in memory.

When processing  $x_t$ , replace  $y$  with  $x_t$   
with "some probability" such that it remains  
uniformly random."

## RESERVOIR SAMPLING:

initialize  $y = x_1$ .

for  $t$  in  $2, 3, 4, \dots, m$ :

observe token  $x_t$ .

with probability  $1/t$ :

$y \leftarrow x_t$

else:

$y$  remains unchanged

endfor

output  $y$ .

Claim  $\forall t \geq 1$ , after  $t$  rounds of this loop,  
 $y$  is unlr distrib over  $\{x_1, \dots, x_t\}$ .

Proof. Induct on  $t$ .

Base case  $t=1$ :  $y = x_1$  ✓

Induct step: after loop iteration  $t+1$ ,

$$\Pr(y = x_{t+1}) = \frac{1}{t+1} \text{ by design.}$$

For any  $1 \leq s \leq t$ ,

$$\Pr(y = x_s) = \Pr(y \text{ unchanged in iteration } t+1)$$

$$\cdot \Pr(y = x_s \text{ at start of iter})$$

$$= \left(1 - \frac{1}{t+1}\right) \cdot \frac{1}{t}$$

$$= \frac{1}{t+1}$$

RESERVOIR SAMPLING for drawing  
 a  $k$ -element subset of  $\{x_1, \dots, x_m\}$   
 uniformly at random from all  
 $k$ -element subsets.

initialize  $(y_1, \dots, y_k) = (x_1, \dots, x_k)$

for  $t = k+1, k+2, \dots, m$ :

with probability  $\frac{k}{t}$ :

// Store  $x_t$  in  $\vec{y}$ , overwriting  
 a random stored element

Sample index  $i \in [k]$  unif. random

$y_i \leftarrow x_t$

else:

ignore  $x_t$ , leave  $\vec{y}$  unchanged.

Correctness proof: again by induction.

(omitted, but in book.)

Space requirement:  $k \cdot b + \log m$   
 $y_1, \dots, y_k$   $t$ .

## Quantile estimation.

Tokens are elements of  $[N] = \{1, 2, \dots, 2^b\}$ .

Goal. Store a sketch of stream  $x_1, \dots, x_m$   
and afterward answer quantile queries:

$$q(i) = \frac{\#\{j \mid x_j \leq i\}}{m}$$

$(\epsilon, \delta)$ -PAC quantile est.

On query  $i$ , output  $\hat{q}(i)$  s.t.

$$\Pr(|\hat{q}(i) - q(i)| > \epsilon) < \delta.$$

$(\epsilon, \delta)$ -PAC uniform quantile est.

Output a data structure encoding  
a vector of quantile estimates



$(\hat{g}(0), \hat{g}(1), \dots, \hat{g}(m))$  such that

$$\Pr(\exists i \text{ with } |\hat{g}(i) - g(i)| > \epsilon) < \delta.$$

Clearly  $(\epsilon, \delta)$ -PAC UQE  $\implies$   $(\epsilon, \delta)$ -PAC QE

$$(\epsilon, \frac{\delta}{N})\text{-PAC QE} \implies (\epsilon, \delta)\text{-PAC UQE.}$$

First let's talk about QE.

Claim: For appropriate choice of  $k$ , the

following algo works:

- use reservoir sample to downsample  $x_1, \dots, x_m$  to  $y_1, \dots, y_k$ .

- using  $y_1, \dots, y_k$  simulate  $k$  indep.

stream samples with replacement  $z_1, \dots, z_k$ .

- answer query  $g(i)$  with

$$\hat{g}(i) = \frac{\#\{l \mid \cancel{y_l} \leq i\}}{k}$$

How large must  $k$  be for this to work?

To show.  $|\hat{g}(i) - g(i)|$  likely small.

Idea. Use Hoeffding.

$$X_\ell = \begin{cases} 1 & \text{if } \cancel{y_\ell} \leq i \\ \emptyset & \text{otherwise} \end{cases}$$

Indeed  $E[X_\ell] = \frac{1}{m} \cdot \# \{j \mid y_j \leq i\}$

because  $y_\ell$  is unif distrib

over  $\{x_1, \dots, x_m\}$ .

$$\hat{g}(i) = \frac{1}{k} (X_1 + X_2 + \dots + X_k).$$

$$\Pr(|\hat{g}(i) - g(i)| > \varepsilon)$$

$$= \Pr(|k \cdot \hat{g}(i) - k \cdot g(i)| > k\varepsilon)$$

$$= \Pr(|(X_1 + \dots + X_k) - E[X_1 + \dots + X_k]| > k\varepsilon)$$

(Hoeff)

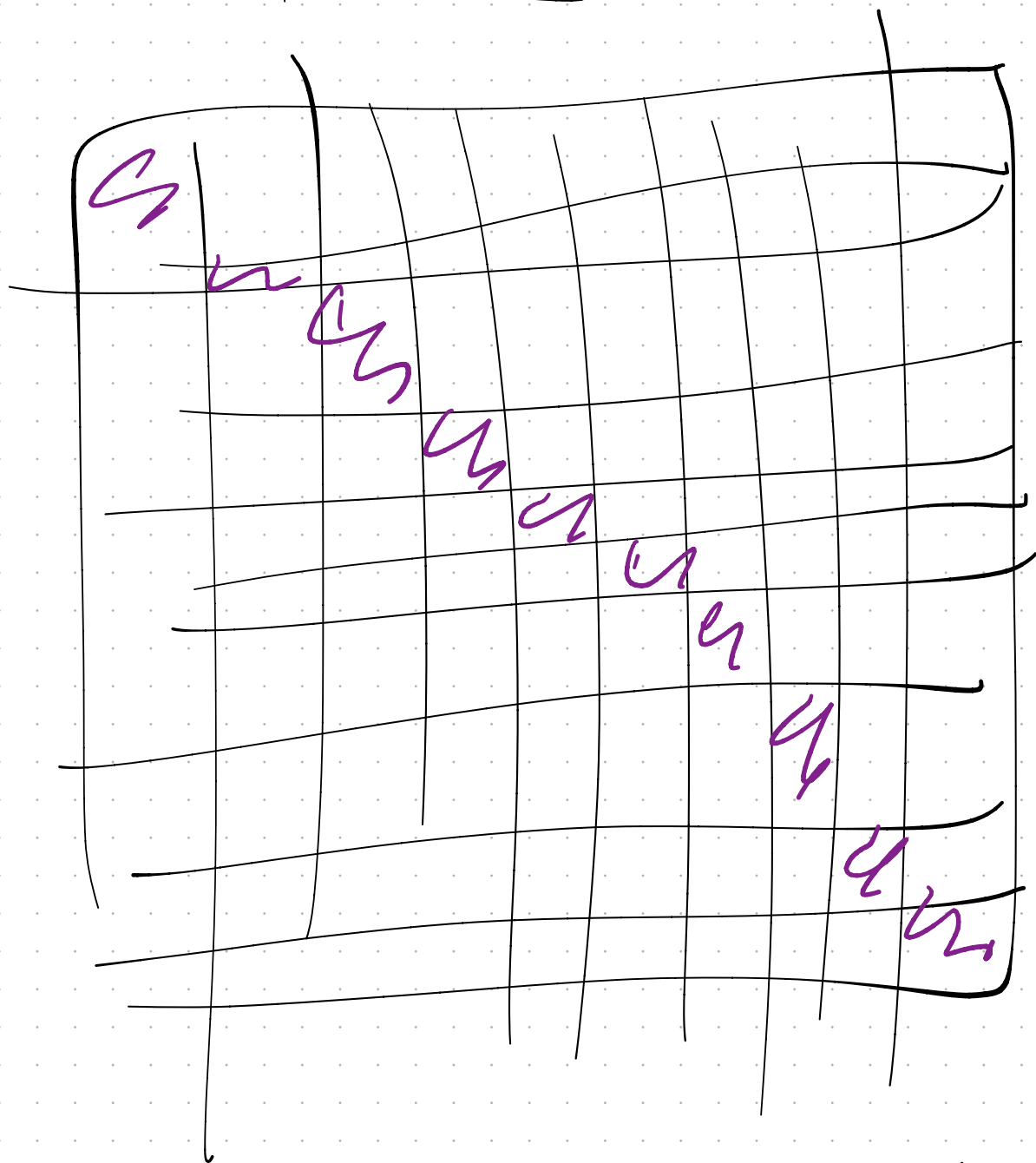
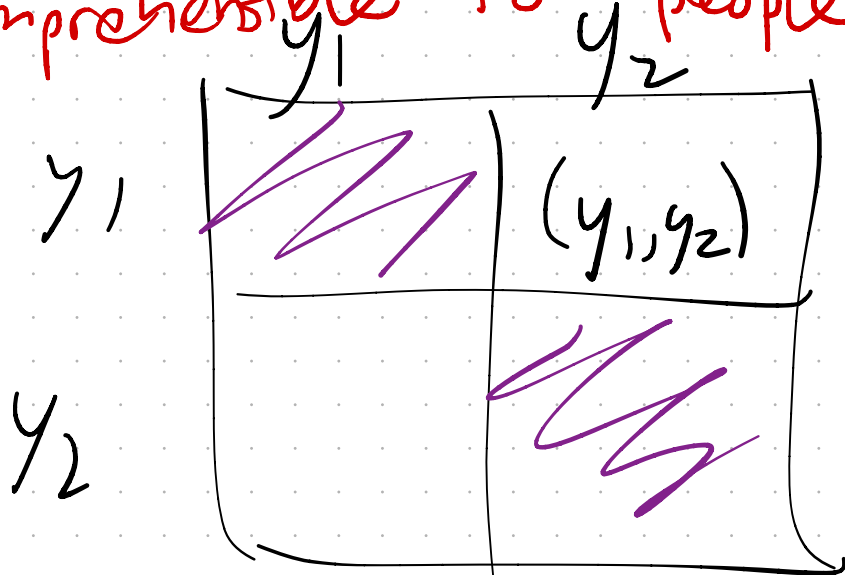
$$\leq 2e^{-\frac{2(k\varepsilon)^2}{k}} = 2e^{-2k\varepsilon^2}$$

To make this  $< \delta$ , set

$$k > \frac{1}{2\varepsilon^2} \ln\left(\frac{2}{\delta}\right)$$

# < BEGIN "SCRAP PAPER" DIGRESSION >

Readers of these notes after the lecture should ignore this part of the notes. They were used as a visual aid but aren't meant to be comprehensible to people who weren't present.



$k=2$ :

$$z_1 = \begin{cases} y_1 & \text{w. prob. } 1/2 \\ y_2 & \text{--- } 1/2 \end{cases}$$

$$z_2 = \begin{cases} z_1 & \text{w. prob. } 1/m \\ \{y_1, y_2\} - \{z_1\} & \text{w. prob. } \frac{m-1}{m} \end{cases}$$

More generally.

After sampling  $z_1, \dots, z_{k-1}$ ;

let  $d = \#$  distinct values in  
 $(z_1, \dots, z_{k-1})$ .

$z_k = \left\{ \begin{array}{l} \text{w prob } \frac{d}{m} \text{ repeat } \text{random} \text{ one of the} \\ \text{d distinct vals} \end{array} \right.$

$\text{w. prob. } \frac{m-d}{m} \text{ sample } \text{random} \text{ one of the}$   
unsampled elements of  
tuple  $(y_1, \dots, y_m)$ .

< END SCRAP PAPER DIGRESSION >

Better still.

Run  $k$  indep copies of reservoir

sampling (each sampling one

stream element) to yield

$(z_1, \dots, z_k)$  each i.i.d uniform

over  $\{x_1, \dots, x_m\}$ .

Then  $\hat{g}(i) = \frac{\# \{ \ell \mid z_\ell \leq i \}}{k}$ .

Summary. sample size

$$k = \frac{1}{2\epsilon^2} \ln\left(\frac{2}{\delta}\right)$$

suffices for  $(\epsilon, \delta)$ -PAC QLE.

$$\implies k = \frac{L}{2\epsilon^2} \ln\left(\frac{2N}{\delta}\right)$$

suffices for  $(\epsilon, \delta)$ -PAC UQE.

In fact

$$k = \frac{L}{2\epsilon^2} \ln\left(\frac{2N}{\delta}\right)$$

samples are sufficient even for  
 $(\epsilon, \delta)$ -PAC UQE.

DKW inequality

(Dvoretzky-Kiefer-Wolfowitz)

If  $z_1, \dots, z_k$  are iid draws from a distribution with CDF  $F$ , and  $\hat{F}_z$  denotes the empirical CDF

$$\hat{F}_z(a) = \frac{\#\{z_i \leq a\}}{k}$$

then  $\Pr\left(\|\hat{F}_z - F\|_\infty > \varepsilon\right) < 2e^{-2\varepsilon^2 k}$

N.B. For a single  $a \in \mathbb{R}$ ,

$$\Pr\left(|\hat{F}_z(a) - F(a)| > \varepsilon\right) < 2e^{-2\varepsilon^2 k}$$

(Like Hoeffding's  $bd$  + union  $bd$  combined without requiring union  $bd$ )

$\Rightarrow$   $(\varepsilon, \delta)$ -PAC UQE requires space

$$O\left(\log m + \frac{1}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right) b\right)$$

## → Glivenko-Cantelli Theorem

If  $z_1, z_2, \dots$  is an infinite sequence of iid samples from distn. with CDF  $F$  and

$$\hat{F}_t(a) = \frac{\#\{s \mid 1 \leq s \leq t, z_s \leq a\}}{t}$$

= empirical CDF of first  $t$  samples

then  $\Pr\left(\|\hat{F}_t - F\|_{\infty} \rightarrow 0 \text{ as } t \rightarrow \infty\right) = 1$

Reason. DKW implies  $\forall \epsilon > 0$

$\mathbb{P}\left(\#t : \|\hat{F}_t - F\|_{\infty} > \epsilon\right)$  is finite

→ w. prob. 1,

$\|\hat{F}_t - F\|_{\infty} \leq \epsilon$  for all suff large  $t$ .