Quantile estimation,
Reservoir sampling,
Glivenko—Cantelli

If tokens in a stream come from an ordered set, e.g. $[m]$, we can define the quantile function a.k.a. the __empirical CDF__

$$q(a) = \frac{\#\{i \mid a_i \leq a\}}{n}$$

# stream tokens $\leq a$

total length of stream.

Two goals for quantile sketching.
The algorithm estimates $q(a)$ using estimator $\hat{q}(a)$.

$\underline{(\varepsilon, \delta)\text{-PAC}}$: $\forall a \quad \Pr\left(|\hat{q}(a) - q(a)| > \varepsilon\right) < \delta$.

$\underline{\text{Uniformly } (\varepsilon, \delta)\text{-PAC}}$:
$$\Pr\left(\|\hat{q} - q\|_\infty > \varepsilon\right) < \delta$$

where $\|\hat{q} - q\|_\infty$ denotes $\max_{a \in [m]} |\hat{q}(a) - q(a)|$.

This lecture in a nutshell:

1. **Downsampling** the stream — maintaining in memory a random sample of $t \approx \frac{1}{\varepsilon^2}$ elements — is a good way to maintain a quantile sketch.

2. **Reservoir sampling** achieves downsampling in $O(t + \log_m(n))$ space.

3. Given $t \approx \frac{1}{\varepsilon^2}$ uniformly random elements of the stream, the empirical CDF of those $t$ elements is an $\varepsilon$-approx to the empirical CDF of the entire stream.

<span style="color:green">Glivenko–Cantelli Theorem</span>

<span style="color:red">DKW Inequality</span>

Dvoretzky – Kiefer – Wolfowitz

## Reservoir Sampling for $t=1$ random elements

Store $a_1$ in memory

for $s = 2, 3, \ldots, n$:

with probability $\frac{1}{s}$ overwrite stored elt with $a_s$

Probability that $a_i$ is stored at the end?

- pick $a_i$ and store it when we first see it. $\left(\frac{1}{i}\right)$

- don't pick $a_s$ for all $s > i$.
$$\left(1 - \frac{1}{i+1}\right) \cdot \left(1 - \frac{1}{i+2}\right) \cdot \cdots \cdot \left(1 - \frac{1}{n}\right).$$

Product of all these probabilities:
$$\frac{1}{i} \cdot \left(\frac{i}{i+1}\right) \cdot \left(\frac{i+1}{i+2}\right) \cdots \cdot \left(\frac{n-1}{n}\right) = \frac{1}{n}.$$

Generalization to selecting $t$ out of $n$.

1. Initialize empty buffer of size $t$.
$$b = (b_1, \ldots, b_t) = (\bot, \bot, \ldots, \bot).$$

2. for $s = 1, \ldots, n$:
   if $s \leq t$ store $b_s = a_s$

   if $s > t$:
   with probability $\frac{t}{s}$:
   sample $i \in [t]$ uniformly random
   overwrite $b_i \leftarrow a_s$
   else: // probability $1 - \frac{t}{s}$
   do nothing

Fact. If $n \geq t$ then the buffer contents after processing stream of length $n$

are a uniformly random $t$-element subset.

## Quantile estimation using reservoir sampling.

Maintain reservoir sample $b_1, \ldots, b_t$ in memory.
When queried about $q(a)$ for $a \in [m]$ report

$$\hat{q}(a) = \frac{\#\{i \mid b_i \leq a\}}{t}$$

a.k.a. $\hat{q}$ is the empirical CDF of the reservoir sample.

## Analyzing the error $\left| \hat{q}(a) - q(a) \right|$

We'll analyze an algorithm where $b_1, \ldots, b_t$ are independent, each is uniformly distributed.

Space: $O\left(t + \log_m(n)\right)$

↑ storing $b_1, \ldots, b_t$

↑ storing the counter $s \in [n]$ when one element of $[m]$ takes $O(1)$ storage.

$$\Pr\left( \left| \hat{q}(a) - q(a) \right| > \varepsilon \right) < 2 e^{-2\varepsilon^2 t}$$

by Hoeffding Bound.

$$X_i = \begin{cases} 1 & \text{if} \quad b_i \leq a \\ 0 & \text{o.w.} \end{cases} \qquad i \in [t]$$

$b_i$ is unif random

$$\mathbb{E}[X_i] = Pr(b_i \leq a) \underset{\downarrow}{=} \frac{\#\{j \mid a_j \leq a\}}{n} = q(a)$$

Let $\quad X = X_1 + \cdots + X_t = t \cdot \hat{q}(a)$

the event $\quad |\hat{q}(a) - q(a)| > \varepsilon \quad$ is equiv't to

$$|X - \mathbb{E}X| > \varepsilon t.$$

$b_1, \ldots, b_t$ indep't $\implies X_1, \ldots, X_t$ indep't

(Hoeffding)
$$\implies Pr(|X - \mathbb{E}X| > \varepsilon t) < 2e^{-2\varepsilon^2 t^2 / t}$$

Summary: For $(\varepsilon, \delta)-$PAC quantile est
maintain $\quad t \quad$ indep. unif. random samples
from the stream where $\quad t \quad$ is chosen
s.t.
$$2e^{-2\varepsilon^2 t} < \delta$$

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

## Upgrading to a uniform $(\varepsilon, \delta)-$PAC estimate

Keep same algorithm. Vary the value of $t$
according to the analysis technique.

## Plan 1. Union bound

There are $m$ events of the form

$$\left| \hat{q}(a) - q(a) \right| > \varepsilon$$

one for each $a \in [m]$.

We can make each have probability less than $\frac{\delta}{m}$ by setting

$$t > \frac{1}{2} \varepsilon^{-2} \ln\left(\frac{2m}{\delta}\right).$$

## Plan 2. More careful union bound

**Lemma.** (weak DKW) If $b_1, \ldots, b_t$ are independent ident distrib samples from any distribution on $\mathbb{R}$ with CDF $F(a)$ and if we define the empirical CDF

$$\hat{q}(a) = \frac{\#\{i \mid b_i \le a\}}{t}$$

then

$$\Pr\left( \|\hat{q} - F\|_\infty > \varepsilon \right) < \frac{4}{\varepsilon} e^{-\varepsilon^2 t/2}.$$

Compare with

$$\Pr\left( \| \hat{q} - F \|_{\infty} > \varepsilon \right) < 2m\, e^{-2\varepsilon^2 t}$$
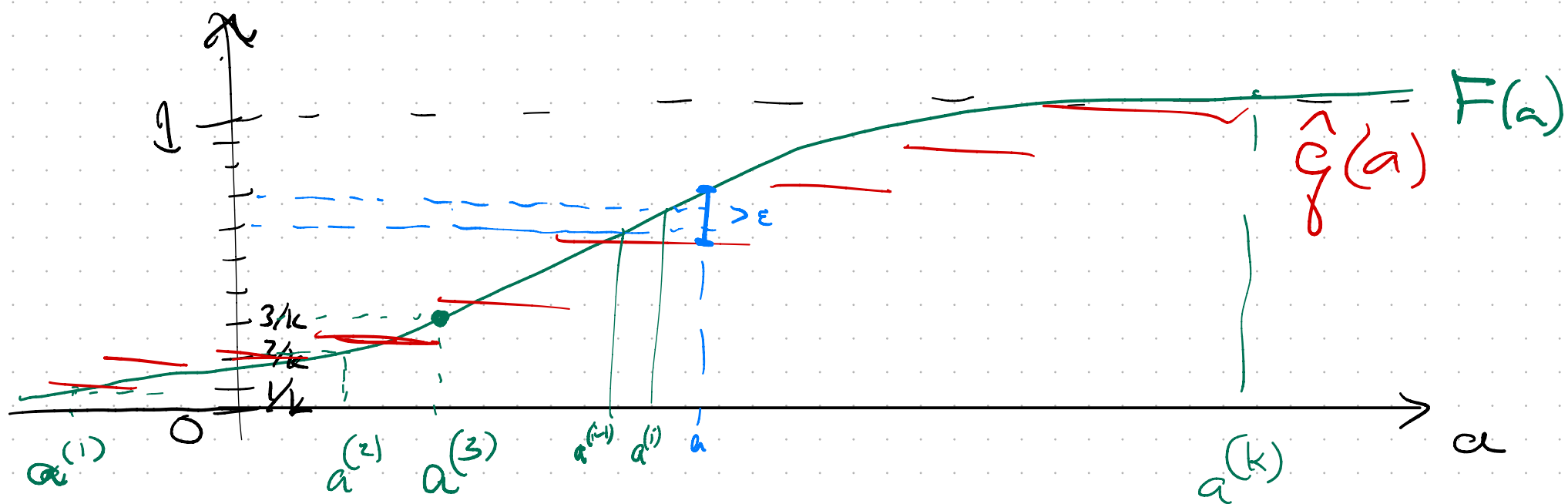
where $F$ supported on $[m]$ and we use naive union bound a la Plan 1.

## Proof of weak DKW.

Let $k = \lceil \frac{2}{\varepsilon} \rceil$.

We'll find a set of $2(k-1) < \frac{4}{\varepsilon}$ bad events, such that if none of them happen,

$$\| \hat{q} - F \|_{\infty} \leq \varepsilon.$$



$$a^{(i)} = \sup \left\{ a \mid F(a) < \frac{i}{k} \right\}.$$

Proof shows that if $\exists a$ s.t.

$$|\hat{q}(a) - F(a)| > \varepsilon \qquad \text{then}$$

$$\exists i \in [k-1] \quad \text{s.t.}$$

$$F(a^{(i)}) - \hat{q}(a^{(i)}) > \varepsilon/2 \qquad \text{or}$$

$$\check{q}(a^{(i)}) - F(a^{(i)}_-) > \varepsilon/2$$

$$\overset{\|}{\phantom{x}} \qquad\qquad \overset{\|}{\phantom{x}}$$

$$\sup\{\hat{q}(a) \mid a < a^{(i)}\} \qquad \sup\{F(a) \mid a < a^{(i)}\}$$