

4 March 2026

Finishing Count Sketch  
+ Reservoir Sampling

## COUNT SKETCH

initialize array  $C$  of dimension  $B \times t$   
 $C \equiv 0$ .

sample random hash functions (2-univ)

$$h_1, h_2, \dots, h_t: [N] \rightarrow [B]$$

$$g_1, g_2, \dots, g_t: [N] \rightarrow \{\pm 1\}$$

for  $i = 1, \dots, n$ :

$$(b_1, \dots, b_t) = (h_1(x_i), \dots, h_t(x_i))$$

$$\forall j \quad C[b_j, j] = C[b_j, j] + g_j(x_i)$$

// could increment or decrement

end for

QUERY( $x$ ):

for  $j = 1, \dots, t$  let

$$v_j = C[h_j(x), j] \cdot g_j(x)$$

// in absence of hash collisions  
this would be exactly  $f_x$ .

output  $\hat{f}_x = \text{median}(v_1, \dots, v_t)$ .

Analysis focuses on one

estimator  $v_l$ ,  $l \in [t]$ .

Fix  $x$  (taken whose

frequency is queried).

$$X_{yl} = \begin{cases} 1 & \text{if } h_l(y) = h_l(x) \\ 0 & \text{otherwise} \end{cases}$$

$$Z_{yl} = g_l(x) g_l(y) X_{yl} f_y$$

the "noise term" in  $v_l$   
coming from token  $y$ .

Fact:  $v_l = \sum_{y \in [N]} Z_{yl}$

$$E(v_e) = \sum_y E(z_{ye})$$

$$= \sum_y E[g_e(x) g_e(y) X_{ye} f_y]$$

$$= f_x + \sum_{y \neq x} E[g_e(x) g_e(y) X_{ye}] f_y$$

PW IND  $\underbrace{E[g_e(x) g_e(y)]}_{\text{only depends on } g_e} \cdot \underbrace{E[X_{ye}]}_{\text{only depends on } h_e}$

~~$E[g_e(x)] \cdot E[g_e(y)]$~~

$$E(v_e) = f_x$$

Now focus on  $\text{Var}(v_e)$ , with an eye toward using Chebyshev.

$$\text{Var}(v_x) = \mathbb{E} \left[ (v_x - F_x)^2 \right]$$

$$= \mathbb{E} \left[ \left( \sum_{y \neq x} z_{yl} \right)^2 \right]$$

$$= \sum_{y \neq x} \mathbb{E} [z_{yl}^2] + \sum_{y \neq x} \sum_{w \notin \{x, y\}} \mathbb{E} [z_{yl} z_{wl}]$$

$$z_{yl}^2 = \cancel{g_l(x)} \cancel{g_l(y)} \overset{+1}{X_{yl}} \overset{+1}{f_y^2}$$

$$= X_{yl} f_y^2$$

$$\mathbb{E} [z_{yl}^2] = \mathbb{E} [X_{yl}] f_y^2 = \frac{1}{B} f_y^2$$

$$z_{yl} z_{wl} = \cancel{g_l(x)} \cancel{g_l(y)} \overset{+1}{X_{yl}} \overset{+1}{f_y} \cancel{g_l(x)} \cancel{g_l(w)} \overset{+1}{X_{wl}} \overset{+1}{f_w}$$

$$\mathbb{E} [z_{yl} z_{wl}] = \mathbb{E} \left[ \underbrace{g_l(y) g_l(w)}_{\text{dep on } g} \underbrace{X_{yl} X_{wl}}_{\text{dep on } h} \right] f_y f_w$$

$$\begin{aligned}
&= E(g_\mu(y)g_\mu(w)) \cdot E(x_y x_w) f_y f_w \\
&= E(g_\mu(y)) E(g_\mu(w)) E(x_y x_w) f_y f_w \\
&= 0 \cdot 0 \cdot E(x_y x_w) f_y f_w \\
&= 0
\end{aligned}$$

$$E(v_\ell) = f_x$$

$$\text{Var}(v_\ell) = \frac{1}{B} \sum_{y \neq x} f_y^2 \leq \frac{1}{B} \|f\|_2^2$$

$$P\left(|v_\ell - f_x| > \varepsilon \|f\|_2\right) < \frac{\text{Var}(v_\ell)}{\varepsilon^2 \|f\|_2^2}$$

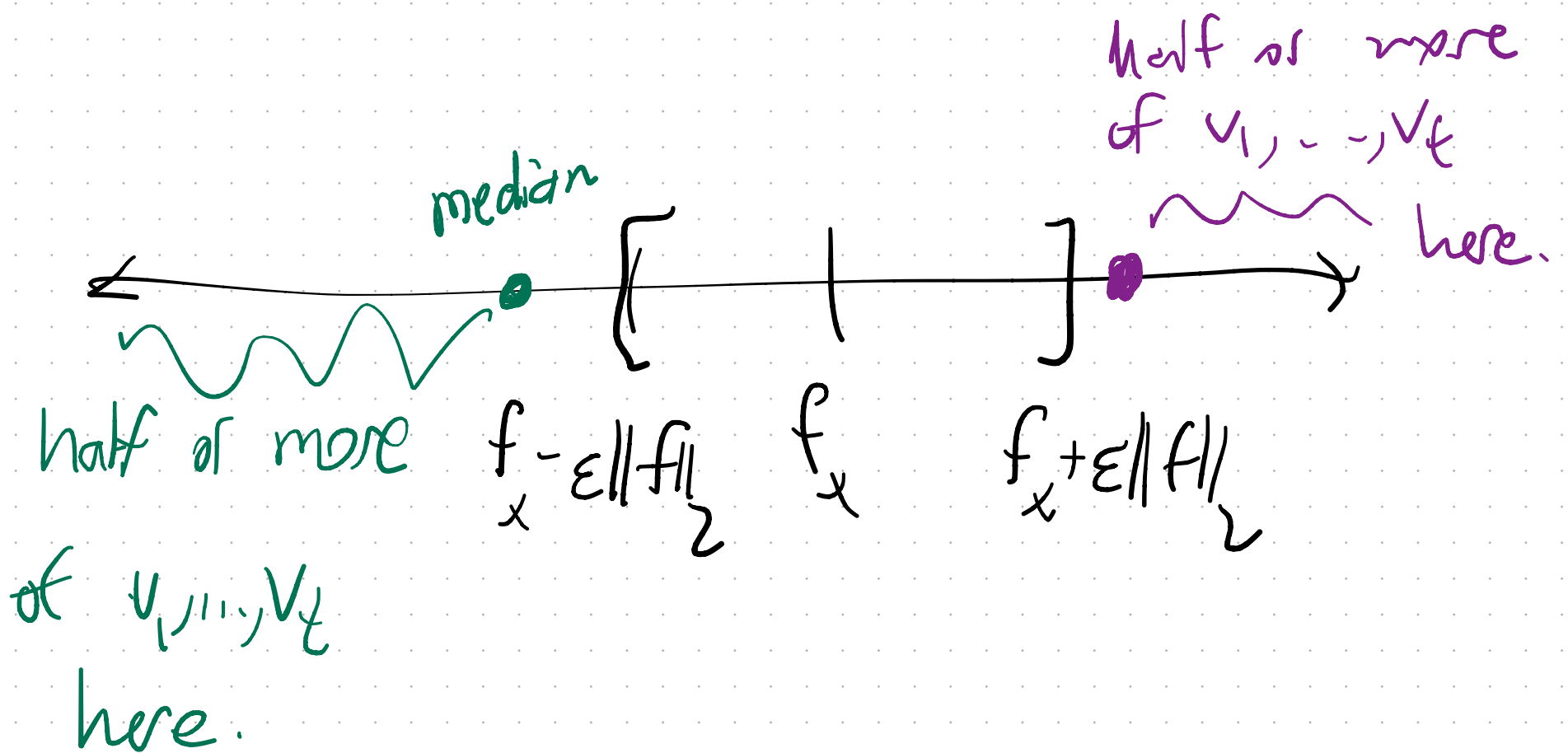
$$< \frac{1}{\varepsilon^2 B} < \frac{1}{3}$$

... by setting  $B > 3/\varepsilon^2$ .

If median  $(v_1, \dots, v_t)$  differs from  $f_x$   
 by more than  $\varepsilon \|f\|_2$

it means

$$\#\{l \mid |v_l - f_x| > \varepsilon \|f\|_2\} \geq \frac{t}{2}.$$



$$\mathbb{E} \left[ \#\{l \mid |v_l - f_x| > \varepsilon \|f\|_2 \right] < \frac{t}{3}.$$

Hoeffding  $\Pr(\text{this}) < e^{-\frac{2 \cdot (t/6)^2}{t}} = e^{-t/18}.$

To make  $e^{-t/18} < \delta$  set

$$t > 18 \ln(1/\delta).$$

Conclusion: Count sketch

with  $B = \Theta\left(\frac{1}{\varepsilon^2}\right)$

$$t = \Theta\left(\ln \frac{1}{\delta}\right)$$

satisfies  $\varepsilon \cdot \|f\|_2$  with

$$\text{prob} \geq 1 - \delta.$$

Space requirement:

$$N = 2^b$$

$$O\left(\frac{1}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right) \underbrace{\ln(mN)}_{\ln(m) + b}\right) \text{ bits,}$$

2-universal  
Storing one hash function from

$$[N] \rightarrow [B]$$

takes how much space?

$$B = p \quad \text{prime}$$

$$[N] \subseteq \mathbb{F}_p^d \quad d = \lceil \log_p(N) \rceil.$$

Use inner product  $h_{a,b}(x) = \langle a, x \rangle + b \pmod p$ .

Description of  $h$ :  $d+1$  values mod  $p$ .

coords of  $a$   $b$

$$(d+1) \cdot \log_2(p) \quad \text{bits}$$

$$\leq (\log_p(N) + 2) \cdot \log_2(p) \quad \text{bits}$$

$$= O(\log_2(N))$$

$$B = 2^k \quad k \in \mathbb{N}$$

use  $k$  indep hash functs with  $\{0,1\}$  output

$$k \cdot \log_2(N) \quad \text{bits}$$