

2 Mar 2026

Sketching element frequencies.

Recap: $N = 2^b$ $[N] \cong \{0, 1\}^b$ "space of tokens"

x_1, x_2, \dots, x_m stream of tokens

$\vec{f} \in \mathbb{R}^N$ frequency vector

$$f_x = \# \{i \mid x_i = x\}$$

= number of times x occurs in stream
(a number between 0 and m)

Mista-Ones outputs a list of k tokens guaranteed to include every token x s.t. $f_x \geq \frac{m}{k+1}$.

Sketching \vec{f} means maintaining a data struct C that occupies small space ($\text{poly}(b, \log m)$)

and enables approximately answers queries

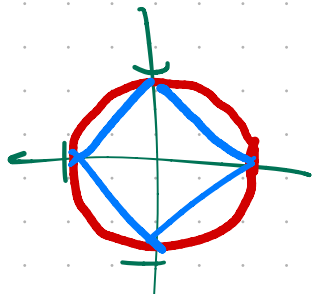
$$x \mapsto \hat{f}_x \quad (\text{estimate of } f_x)$$

(ϵ, δ) -PAC

$$\Pr(\text{normalized error} > \epsilon) < \delta.$$

Count-Min Sketch: $\Pr(|\hat{f}_x - f_x| > \epsilon \|\vec{f}\|_1)$

Count Sketch: $\Pr(|\hat{f}_x - f_x| > \epsilon \|\vec{f}\|_2)$



2-D. $\|x\|_1^2 = (|x| + |y|)^2 = |x|^2 + 2|x||y| + |y|^2 \geq |x|^2 + |y|^2 = \|x\|_2^2$

(tougher to satisfy)

In \mathbb{R}^N , $\|x\|_1 \geq \|x\|_2$ always
 $\|x\|_1$ can exceed $\|x\|_2$ by factor \sqrt{N} .
 E.g. $x = \vec{1}$.

Count-Min Sketch

Sample random $h_1, \dots, h_t: [N] \rightarrow [B]$
 from \mathcal{Z} -universal family, where $t = \lceil \log_2(\frac{1}{\delta}) \rceil$.
 initialize $C = 0$. (dimension $B \times t$)

for $i = 1, 2, \dots, m$:

let $(b_1, \dots, b_t) = (h_1(x), \dots, h_t(x))$
 $\forall j \ C[b_j, j] \leftarrow C[b_j, j] + 1$ // $C[b]$ counts how
 many tokens hashed
 to bucket b .
 end

Query(x): return $\hat{f}_x = \min_j C[h_j(x), j]$

as. \hat{f}_x never underestimates f_x .

$\hat{f}_x \geq f_x$ with prob. 1.

Q: $\Pr(\hat{f}_x - f_x > \epsilon)$?

$\hat{f}_x - f_x = \# \text{ tokens } x_i \neq x \text{ s.t.}$
 $h(x_i) = h(x)$.

$$E[\hat{f}_x - f_x] = ??$$

Suppose $x_i \neq x$. What is

$$P_r(h(x_i) \neq h(x)) = \frac{B-1}{B}$$

$$P_r(h(x_i) = h(x)) = \frac{1}{B}$$

$$E[\hat{f}_x - f_x] = \sum_{i: x_i \neq x} P_r(h(x_i) = h(x))$$

$$= \frac{1}{B} \cdot \#\{i \mid x_i \neq x\}$$

$$= \frac{1}{B} \cdot (m - f_x) \leq \frac{1}{B} \cdot m = \frac{1}{B} \cdot \|f\|_1$$

Idea #1, Set $B = \frac{1}{\epsilon \delta}$.

$$E[\hat{f}_x - f_x] \leq \epsilon \delta \|f\|_1$$

$$P_r(\hat{f}_x - f_x > \epsilon \|f\|_1) \leq \delta \text{ by Markov.}$$

Space requirement:

• C stores B counters in range $0 \dots m$

$$B \cdot \log(m) = \frac{1}{\epsilon \delta} \log(m) \text{ bits.}$$

• description of h : $O(\log B \cdot \log N) = O(\log(\frac{1}{\epsilon \delta}) \cdot b)$.

Now with $B = \left\lceil \frac{2}{\epsilon} \right\rceil$, $t = \left\lceil \log_2\left(\frac{1}{\delta}\right) \right\rceil$

we have

- Never underestimate $\Pr(\hat{f}_x \geq f_x) = 1$

- Probably overshoot f_x by at most $\epsilon \cdot \|f\|_1$

$$\Pr(\hat{f}_x - f_x > \epsilon \|f\|_1) < \delta$$

- Space requirement:

C needs $\underbrace{\frac{2}{\epsilon}}_{\text{rows}} \times \underbrace{\log_2\left(\frac{1}{\delta}\right)}_{\text{cols}} \times \underbrace{\log m}_{\text{bits per entry}}$

h_1, \dots, h_t needs $O\left(\log\left(\frac{1}{\epsilon}\right) \cdot b \cdot \log\left(\frac{1}{\delta}\right)\right)$

Count-Sketch. Use random signs to bring about cancellations in the "noise" without cancelling "signal."

Algorithm:

$$t = ??$$

$$B = ??$$

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 $i = 1, \dots, m$:

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

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

// could increment or decrement

end for

Query(x):

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

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

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

$$\text{output } \hat{f}_x = \text{median}(a_1, \dots, a_t).$$

What can we prove about this?

- Error is two-sided: $\hat{f}_x - f_x$ could be positive or negative.

- Error likely to be small:

$$P(|\hat{f}_x - f_x| > \epsilon \|f\|_2) < \begin{array}{l} \text{small when} \\ B \gg \frac{1}{\epsilon^2} \\ \text{and } t \gg \log_2\left(\frac{1}{\epsilon}\right). \end{array}$$

- Space

C requires $B \times t \times \log(2m)$ bits

$h_1, \dots, h_t, g_1, \dots, g_t$

$$O(\log B \cdot t \cdot b).$$

$$\text{Total space} = O(\text{Greek} \cdot (b + \log m)).$$

Analysis of count sketch

Focus on column j with

hash functions $h = h_j, g = g_j,$

estimate $a_j = C[h(x), j] \cdot g(x).$

Claim 1. $E[a_j] = f_x.$

Proof. let $X_{by} = \begin{cases} 1 & \text{if } h(y) = b \\ 0 & \text{if } h(y) \neq b \end{cases}$

$$Z_y = g(y)$$

$$C[b, j] = \sum_{y \in [N]} f_y \cdot X_{by} \cdot Z_y$$

$$a_j = \sum_{b \in [B]} X_{bx} \cdot C[b, j] \cdot Z_x$$

$$E[a_j] = \sum_b \sum_y f_y \cdot E[X_{bx} X_{by} Z_x Z_y]$$

$\hookrightarrow 0$
except if
 $x = y$

$$= \sum_b f_x E[X_{bx}]$$

$$= f_x$$