

L2: Statistical Phenomena

Aug 20, 2025

Data Mining



Jeff Phillips

Data $X = \{x_1, x_2, \dots, x_n\} \sim D(\theta)$
iid
 \uparrow param.
independent identically distributed.

Central Limit Theorem

as $n \rightarrow \infty$, get better estimates

$$X \subset [m] = \{0, 1, 2, \dots, m-1\}$$

- IP addresses $m = 10^{16}$
- words in English $m = 100,000$
- people in USA $m = 360 \text{ M.}$

$$D(\alpha) = \text{Unif}([m])$$

$$\Pr(x_i = j) = \frac{1}{m} \quad \text{if } j \in [m]$$

0 otherwise

Hash Function

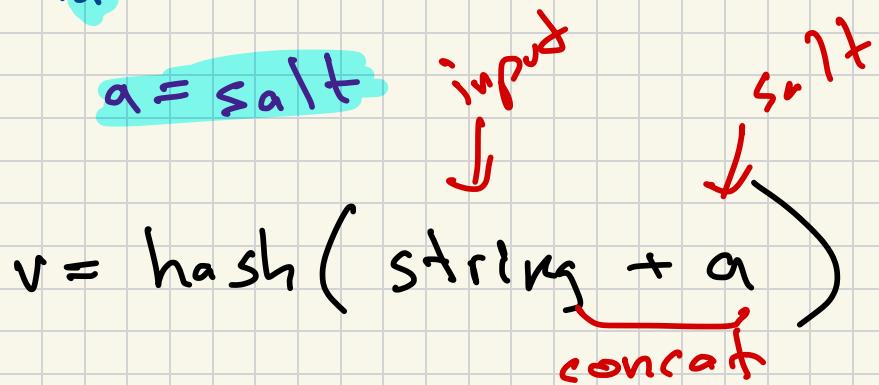
choose $h_{a \sim \mathcal{H}}$

$$\underline{h} : \Sigma^* \rightarrow [m] \quad \text{Deterministic}$$

SHA1

$$\Pr_{\substack{h_{a \sim \mathcal{H}}}} [h_a(\text{string}) = j] = \frac{1}{m} \quad \text{if } j \in [m]$$

0 otherwise.



Birthday Paradox

$m = \text{unique days year}$
365

Jan	26, 12, 1
Feb	26
Mar	6, 19
Apr	30, 5
May	8, 14, 25, 11
June	26, 21, 23
July	10
Aug	9, 16, 22, 30, 29, 15, 24
Sep	21, 15, 29, 11
Oct	17, 24
Nov	
Dec	

$\Pr(\text{collision } n \text{ trials})$

$$\Pr(\text{coll } n=1) = 0$$

$$\Pr(\text{coll } n=2) = \frac{1}{365}$$

$$\Pr(\text{coll } n=3) = 1 - \left(1 - \frac{1}{365}\right)^3$$

$$\Pr(\text{coll } n) = 1 - \underbrace{\left(1 - \frac{1}{365}\right)}_{0.997}^{\binom{n}{2}} \approx \frac{n^2}{2}$$

if $n=23$ $\binom{23}{2} = 253$

$$1 - (0.997)^{253} \approx 0.532$$

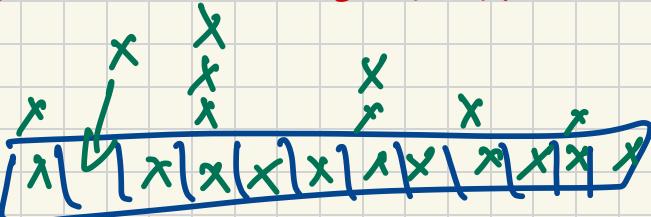
When first collision
 $n \approx \sqrt{2m}$

What is wrong?

- birthday are not uniform?
- leap year $m = 365$ Feb 29
also not uniform
- math correlation
in pairs
- samples independent?
(twins)

Coupon Collectors

Size m domain \rightarrow See all distinct elements $j \in [m]$



How many trials until see all $[m]$

- at least m .

0.577...

$$\frac{2^m}{m!} \approx m!$$

$$m^3$$

$$2^m$$

$$E[\text{\#trials}] = m(\gamma + \ln m)$$

t_j = # trials between j th new coupon
and the $(j+1)$ th new coupon

$$t_0 = 1$$

$$\mathbb{E}[t_j] = \frac{m}{m-j}$$

$$j = \frac{m}{2}$$

$$\frac{m}{m-j} = \frac{m}{m-\frac{m}{2}} = \frac{m}{\frac{m}{2}} = 2$$

$$j = m-1 \quad \frac{m}{m-(m-1)} = m$$

m -th Harmonic
num

$$T = \sum_{j=0}^{m-1} t_j$$

$$\begin{aligned} \mathbb{E}[T] &= \sum_{j=0}^{m-1} \mathbb{E}[t_j] = \sum_{j=0}^{m-1} \frac{m}{m-j} \\ &= m \cdot \sum_{j=0}^{m-1} \frac{1}{m-j} = m \cdot \sum_{j=1}^m \frac{1}{j} \end{aligned}$$

H_m

Summarize

- To use math we need simple models.
- Models are wrong, but sometimes OK.
- Non-trivial phenomena
 - collision $\approx \sqrt{2m}$
 - see all $\approx m \log n$