

UU: Data Mining

L2: Statistical Phenomena

Jan 8, 2025

Data $X = \{x_1, x_2, \dots, x_n\}$

$X \stackrel{\text{iid}}{\sim} D(\theta)$

distribution, parameter θ



observations of the world

each x_i is one observation

$X \stackrel{\text{iid}}{\sim} D(\theta)$

drawn from

iid: independently and identically distributed

Assume domain of $\mathcal{D}(\theta)$ and data X

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

models: IP addresses $m \approx 16^{10}$

words

$m = 100,000$

all people in USA = 3.8 million

~~Assume~~

$x_i \sim \mathcal{D}(\theta)$

$P_i[x_i = j] = \frac{1}{m}$ if $j \in [m]$
 > 0 o.w.

random hash function

$$h: \Sigma^t \rightarrow [m]$$

h : deterministic

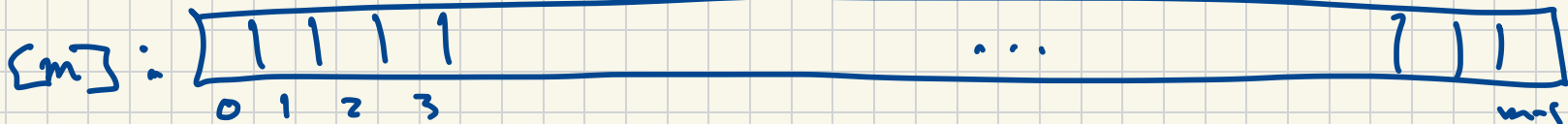
$$\Pr[h(i) = h(j)] = \frac{1}{m}$$

$h \sim \mathcal{H}$ $i \neq j$

$h_a \sim \mathcal{H}$ a : salt \leftarrow chosen at random

$$h_a = f(\text{string}, a)$$

complex $\in \Sigma^t$



Complex Deterministic Function

$$f: (\mathbb{Z}_1^t, \text{salt}) \rightarrow [m]$$

- SHA-1: hard to invert f

$$\boxed{f^{-1} ?}$$

$$\rightarrow [m]$$

$$m = 2^{160}$$

- Multiplicative Hashing $x \in \mathbb{Z}^* = \mathbb{R}$

$$h_a(x) = \lfloor m \cdot \text{frac}(x \cdot a) \rfloor$$

$$\lfloor (x \cdot a / 2^8) \rfloor$$

$$\text{frac}(17.23) = 0.23$$

$$\lfloor 17.23 \rfloor = 17$$

~~Do~~ not do modular hashing

$$h(x) = x \bmod m$$

Birth day Paradox

$m = 365$

$n = 23$

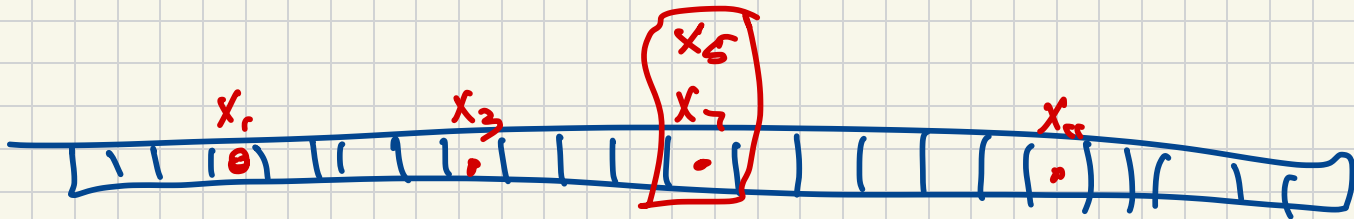
Draw $x_i \sim X$

$x_i \in [m]$ $P_i(x_i = j) = \frac{1}{m}$

x_1, x_2, \dots, x_n

When do I expect that two $x_i = x_j$
if $i \neq j$?

Answer: $n \approx \sqrt{2m}$



Jan
 Feb 14
 Mar 20
 Apr
 May (23)
 Jun 3, 26
 Jul, 25
 Aug
 Sep
 Oct
 Nov (y6)
 Dec

21, 20² collision of n=6
 19, 2, 18, 12
 31, 19, 24
 29, 29, 23
 24, 14, 22, (30), 29, 7, 21
 20, 17
 25, 1, 15, 5, 11
 17, 31, 5
 4, 21, 9, 24
 25, 11, 2, 7, 21, 12
 22, 20
 8, 14, 12

(x46)

$$Pr[\text{coll}, n=1] = 0$$

$$Pr[\text{coll}, n=2] = \frac{1}{m}$$

$$Pr[\text{coll}, n=3] \approx \left(1 - \left(1 - \frac{1}{m}\right)\right)^3$$

$$Pr[\text{coll}, n] \approx 1 - \left(1 - \frac{1}{m}\right)^{\binom{n}{2}}$$

$$\approx 1 - \left(1 - \frac{1}{m}\right)^{n^2/2}$$

$$n=23$$

$$\begin{aligned} \binom{n}{2} &= \# \text{ pairs} \\ &= \frac{n(n-1)}{2} \approx \frac{n^2}{2} \end{aligned}$$

$$\approx 0.048$$

$$n=6$$

$$\approx 1 - 0.997^{253} = 0.532$$

$$n=23$$

What is wrong w/ this analysis?

• bias : more birthdays in Spring.

↳ set of pairs ^{leap year} ~~not~~ ^{m=366 (Feb 29)} ~~real~~ ^{not} ~~birthdays~~ ^{not} iid.

What happens when $n = m + 1$

$$1 - \left(1 - \frac{1}{m}\right)^{\binom{m+1}{2}} \approx 1 - 0.0000007$$

$$1 - \left(\frac{m-1}{m}\right) \left(\frac{m-2}{m}\right) \left(\frac{m-3}{m}\right) \dots (0) = 1$$

Coupon Collector's

$x_1, \dots, x_n \stackrel{iid}{\sim} \text{Unif}([m])$

How many draws n until
I see all $j \in [m]$ so
some $x_i = j$?

$$n \approx 4m^2 ?$$

$$\approx 5m$$

$$m \log m = m(0.577 + \ln m)$$

$$E[n] = m \sum_{i=1}^m \frac{1}{i}$$

$\Rightarrow H_m$ harmonic number

$$\approx (0.577 + \ln m)$$