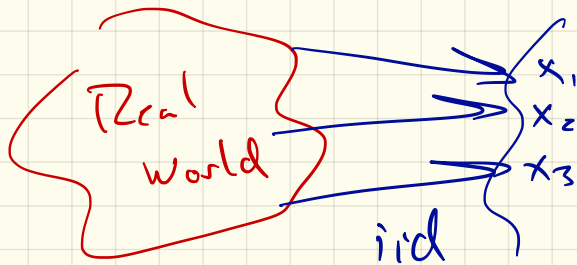


Data Mining LZ

Statistical Principles + Hashing

IID Data

Independent and Identically Distributed



input data

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

Set

This lecture

Assume iid set $X = \{x_1, x_2, \dots, x_m\}$

$x_i \in [n] = \{1, 2, 3, \dots, n\}$
in

Represent $[n]$: All possible IP addresses

All words in dictionary

All possible birthdays

$n = 365$

Assume each x_i uniform
in $[n]$

$$P_r [x_i = j] = \frac{1}{n} \quad \text{if } j \in [n]$$

0 o.w.

Hash Table and Hash Functions (Random)

Family Hash Functions \mathcal{H}

$h_a \in \mathcal{H} \leftarrow$ choice random

$h_a : \Sigma^k \rightarrow [n]$ deterministic

$$\Pr_{h_a \in \mathcal{H}} [h_a(\text{string1}) = h_a(\text{string2})] = \frac{1}{n}$$

as long as $\text{string1} \neq \text{string2}$

1. Built in Hash Function

ex. SHA-1

$$h_a(x) = \text{SHA-1}(\text{concat}(a, x))$$

↑
"salt"

Some string.

2. Multiplicative Hashing

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

↑
salt

$$\text{frac}(11.278) = 0.278$$

$$h_a(x) = \frac{xa}{2^q} \bmod m$$

large int

with binary representation

mix 0s 1s

3. Modular Hashing

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

Do Not Use

Input : sequence distinct strings

each hash with $h_a \in \mathbb{C}$

$$h_a(\text{str}) \rightarrow [n]$$

Qt: How many until collision.

"Birthday Paradox"

18 trials

Jan	1, 9
Feb	8, 15, 27
Mar	
Apr	24, 19
May	(17)
Jun	25, 8
Jul	
Aug	25, 20
Sep	17
Oct	
Nov	3, 30
Dec	3, 27

$$Pr[s_1 = s_2] = \frac{1}{n}$$

$$k \text{ people} \rightarrow \binom{k}{2} = \frac{k(k-1)}{2} \text{ pairs}$$

$$Pr. [\text{no coll } k \text{ pairs}] = \left(1 - \frac{1}{n}\right)^{\binom{k}{2}}$$

$$n = 365 \quad k = 23 \Rightarrow 0.467$$

pigeon hole
principal

?

$$k = 366$$

$$\left(1 - \frac{1}{n}\right)^{\binom{366}{2}} > 0$$

• Assume uniform prob.

• True prob after k steps

$$1 - \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \left(\frac{n-3}{n}\right) \dots \left(\frac{n-k}{n}\right)$$
$$= 1 - \prod_{i=1}^k \left(\frac{n-i}{n}\right)$$

Prob [coll > 50] after $\Rightarrow k = \sqrt{2n}$ steps

$$k = 18 \quad \text{Prob (coll)} = 0.34$$

$$k = 28 \quad \text{Prob (coll)} = 0.64$$

$$27 \approx \sqrt{2 \cdot 365}$$

QZ: When do we see all birthdays?

$\geq n$

$10 n^2 ?$

$5 n^{1.5} ?$

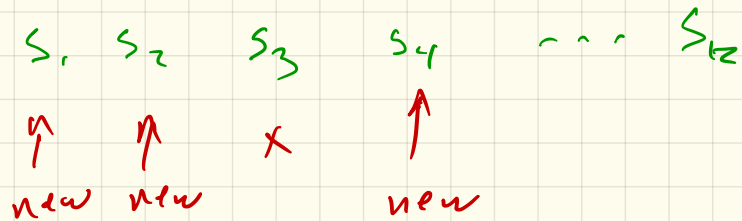
$25 \Theta(n) ?$

$40 n \log n ?$

$(n!) ? \equiv n \text{ factorial}$

Coupon collectors

Analyze Sequence



Let $T_i = \#$ steps until
ith distinct observation

$$t_i = T_i - T_{i-1}$$

$\equiv \#$ steps between (i-1)th distinct
and ith distinct.

$$\text{Expected total } \# = E\left[\sum_{i=1}^n t_i\right] = \sum_{i=1}^n E[t_i]$$

$$E[t_i] = \frac{n - (i-1)}{n} = \frac{n}{n-i+1} \quad H_n = n \text{th Harmonic number}$$

$$\sum_{i=1}^n E[t_i] = \sum_{i=1}^n \frac{n}{n-i+1} = n \sum_{i=1}^n \frac{1}{n-i+1} = n \sum_{i=1}^n \frac{1}{i}$$

$\frac{1}{n}, \frac{1}{n-1}, \frac{1}{n-2}, \dots, 1$

Prob [seeing ith distinct, given i-1 distinct]

$$= \frac{1}{n - (i-1)} = \frac{n - (i-1)}{n} = \frac{n-i+1}{n}$$

$$H_n = \gamma + \ln(n) + o(1/n) \approx 0.577$$

$$T_n = n H_n = n (0.577 + \ln n)$$

Probably Approximately Correct (PAC)

want to estimate μ

estimate \hat{x}

$$\delta \in [0, 1]$$

$$\hat{x}, \mu \in [0, 1]$$

$$\Pr[|\hat{x} - \mu| > \epsilon] < \delta$$

acceptable error

prob. of failure

$$0 \leq \epsilon < 1$$

$$\Pr[|\hat{x} - \mu| < \epsilon] \geq 1 - \delta$$

trials

$$\approx \frac{1}{\epsilon^2} \log \frac{1}{\delta}$$

