

Statistical Principles

Note Title

1/13/2016

Hashing

- Birthday Paradox
- Coupon collector's Problem
- Central Limit Theorem

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

$x_i \sim D(\theta)$

i.i.d



parameters

- independently and identically distributed

Big Data

complex / accurate estimates

combining many small observations

~ ~

Hash Function

$$h_a \in \mathcal{H}$$

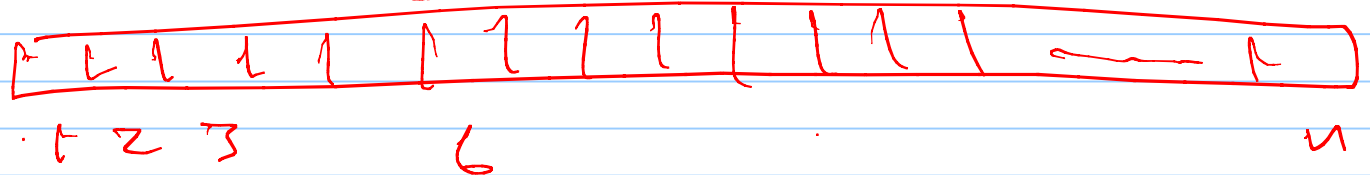
$$h: [m] \rightarrow [n]$$

all IP addresses

$$\{1, 2, 3, \dots, n\}$$

$$x \in [m]$$

$$h(x)$$



$$\Pr[h_a(x) = h_a(x')] = \frac{1}{n} \quad x, x' \in [m]$$

SHA-1(Σ)

$\rightarrow n = 2^{160}$

Σ = string of bits
 $\{0, 1\}^*$

add salt : string a

SHA-1(concat(a , x))

input x

Multiplicative Hashing

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

Modular Hashing : $h(x) = x \bmod n$

Birthday Paradox

$n: [\text{person}] \rightarrow [365]$

Jan: 24, 6, 1, 31, 9

Feb: 2, 29, 1, 9

Mar: 24, 6, 1, 31, 9

Apr: 5, 2, 17

May: 20, 21, 17

June: 20, 21, 17

July: 20, 21, 17

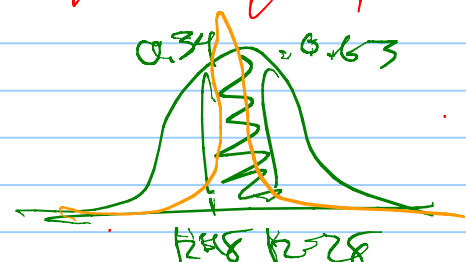
Aug: 20, 21, 17

Sep: 20, 21, 17

Oct: 12, 3

Nov: 23

Dec: 25



$P_r[\text{Collision } k \in [18, 28]]$

1st person

2nd person

1st collision

$365 \quad k \approx \sqrt{2 \cdot n}$

k people

$\left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \dots \left(\frac{n-k}{n}\right)$

$\binom{n}{2} \text{ pairs} \approx \frac{n^2}{2} \approx \frac{n^2}{2}$

$P_{\text{coll}}[\text{no coll}] \approx \left(1 - \frac{1}{365}\right)^k$

$k=23$

$(0.997)^{253} = 0.467$

Jan : ~~1111~~

Feb : 11

Mar : 11

Apr : ~~1111~~

May : ~~1111~~

June : 1111

Jul : 1111

Aug : 1111

Sep : 11

Oct : 1111

Nov : ~~1111~~

Dec : 1

$$n^2 = 1404$$

$r_0 = \#$ frocks before coupon = 36
Seeding with distinct n coupons

$$t_0 = r_0 - r_{0-1} = r_0 - r_{0-1}$$

$\#$ frocks between $(i-1)$ th dist coupon and i th

4 |

$$r_n = \sum_{i=1}^n t_i$$

$$P[\text{single}] = \frac{n-1}{n} \quad \left[\frac{t_i}{r_i} \right]$$

$$E[r_n] = \sum_{i=1}^n \frac{n}{n-i+1} = n \left(\sum_{i=1}^n \frac{1}{i} \right) \approx \ln n \approx n(0.6 + \ln n)$$

Coupon Collectors

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

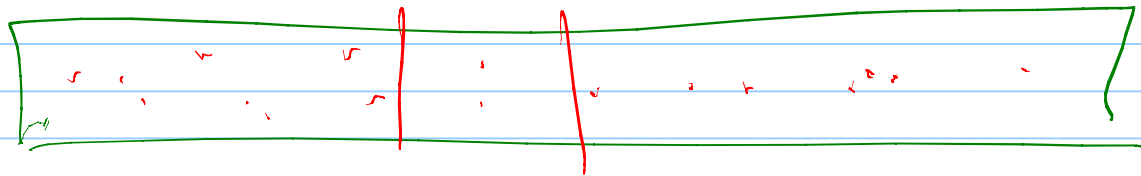
$h^+ \rightarrow [n]$ How many draws to get all
at least n draws

n events $\{P_1, P_2, \dots, P_n\}$ $\sum P_i = 1$

$$\mathbb{E}[R] = \frac{1}{P_*} (\gamma + \ln(n))$$

$$P_* = \min_i P_i$$

$$\varepsilon = \frac{1}{(n+1)} \Rightarrow \frac{1}{\varepsilon} \log \frac{1}{\varepsilon}$$



Central Limit Theorem / Chernoff-Hoeffding Bound

Random Variables X_1, X_2, \dots, X_n

$$-125 \leq X_i \leq 125$$

$$A = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\Pr[|A - \mathbb{E}[A]| \geq \alpha] \leq 2 \exp\left(\frac{-n \alpha^2}{2 \Delta^2}\right)$$

$$\exp(x) = e^x$$

