

L2 Statistical Principles & Hashing

Wednesday, January 8, 2020 2:58 PM

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

where $x_i \sim U$
↑
independent and identically distributed
← unknown distribution

To start, each $x_i \in [m] = \{1, 2, \dots, m\}$

IP Addresses
all possible words
of days in the year.

Here, U is a known distribution

$U = \text{uniform distribution}$

↳ $P(x_i = j) = 1/m \quad \forall j \in [m]$

↳ Appears in algos, data structures when Hashing (Hash table)

(Random) Hash Function:

$h: \text{Domain} \rightarrow \text{Range}$
 $\sum^k [m]$

→ deterministic
(same string maps to same hash)

• Randomly select $h_a \in \mathcal{H}$ (family of hash functions)

$P[h_a(x) = h_a(y)] = \frac{1}{m} \quad \text{where } x \neq y$

3 versions of hash functions.

• Built-in hash functions

↳ SHA-1: $(\Sigma = \{0, 1\})^k \rightarrow [m = 2^{160}]$

a = salt

↳ input $x \rightarrow \text{SHA-1}(\text{concat}(x, a))$

• Multiplicative Hashing

↳ $h_a(x) = \lfloor m \cdot \text{frac}(x \cdot a) \rfloor$ where $\text{frac}(17.32) = 0.32$

$= (x \cdot a / 2^q) \bmod m$ q large number

• Modular Hashing $h(x) = x \bmod m$

(DO NOT USE)

Ex.

1. $u \in \text{Unif}[0, 1]$

2. $\lfloor u \cdot m \rfloor \rightarrow j \in [m]$

Q: How many samples $x_1, x_2, \dots, x_n \in [m]$
so avg two $x_i = x_j$ (a collision)

A: $n \approx \Theta(\sqrt{m})$

Analyzing:

$P[\text{collision}, \text{domain}[m], n \text{ steps}]$

$$n=1 \rightarrow P=0$$

$$n=2 \rightarrow P = \frac{1}{m}$$

$$n=3 \rightarrow P = \left(\frac{1}{m}\right)^2 \approx \left(\frac{1}{m}\right)^{n^2/2}$$

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

$$P[\text{no collision}] = \left(1 - \frac{1}{m}\right)^{\binom{n^2}{2}}$$

$$P[\text{collision}] = 1 - \left(1 - \frac{1}{m}\right)^{\binom{n^2}{2}}$$

$$\text{set } n = \sqrt{2m}$$

$$\Rightarrow 1 - \left(1 - \frac{1}{m}\right)^m \approx 1 - \frac{1}{e}$$

$$P[\text{collision}] = 1 - \left(\frac{m-1}{m}\right)\left(\frac{m-2}{m}\right)\left(\frac{m-3}{m}\right) \cdots \left(\frac{m-(n-1)}{m}\right)$$

Problems with the analysis

1) Birthdays are not iid.

2) Some years have 366 days

3) Twins exist

Q: How many n until we observe all $j \in [m]$? (coupon collector's problem)

A: $n = m \cdot \ln(m)$

Simulation Analysis

$E[r_m]$ where r_i = # of trials until i th distinct item.

$$= E\left[\sum_{i=1}^m t_i\right]$$

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

$$r_1 = 1, r_2 \approx 2$$

and

$$\text{epoch } t_i = r_i - r_{i-1}$$

where

$$E[t_i] = \frac{1}{p_i} = \frac{1}{\left(\frac{m-i+1}{m}\right)} = \left(\frac{m}{m-i+1}\right)$$

$$\text{Now, } E[r_m] = \sum_{i=1}^m E[t_i] = \sum_{i=1}^m \frac{m}{m-i+1} = m \sum_{j=1}^m \frac{1}{j} \text{ where } j = m-i+1$$

H_m (Harmonic #)

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

$\text{truth}(M)$

$\text{Sample}(x)$

$$d(M, \text{Alg}(x)) \leq \epsilon \quad \leftarrow \text{error}$$

$$P[d(M, \text{Alg}(x)) > \epsilon] < \delta$$

$$P[a(\mu, \text{Alg}(x)) > \epsilon] < \delta$$

← prob. failure

Probably Approximately correct
(PAC)