

Map Reduce : Input (k-v), map (병렬가능), reduce (병렬불가)

문제상황: Top 10 url 찾고 count 해라

process: Input k-v  $\rightarrow$  초기 k-v  $\rightarrow$  Groupby  $\rightarrow$  output k-v

문제상황 2: 각 host가 접속하는 url의 총 유효성  $\Rightarrow$  JOIN 으로 해결

$\Rightarrow$  hash function 사용

Map에서 나온 값들을 중간 file에 저장하여 Reduce와 송수신

구성: Master node, map node, reduce node  
 $\hookrightarrow$  노드관리, Health check  
if, map 죽으면 다른 worker, reduce 죽으면 재시작  
master 죽으면 관리자 알림

Backup task: 처리량 너무 많아서 느려진 경우  $\rightarrow$  가까운 마스터에 할당

Combiners: 중간 결과 저장 및 Reducer에 전송

Partition function: Hash 통해 key 비교

Fixed proportion: 비율 고정, uniform distribution 가정.

문제상황: 한 user가 한 날에 같은 쿼리를 얼마나 검색 하는지

해결: user sampling (hash 이용)

Fixed-size tuple: tuple size 고정

Problem: 새로운 데이터 들어올 때 샘플 확률이 다 달라진다.

Solution: 저장 공간  $S$ 보다 time step  $n$ 이 크면  $\frac{S}{n}$  확률로 다음 item을 선택 여부 결정

$\therefore \frac{S}{n}$  확률로 데이터가 샘플 확률 유지

Sliding Window: 가장 최근  $N$ 개만 본다

문제상황: 0과 1이 들어올 때 최근  $k$  bit에 대해 1의 개수는? (단, window size  $k$  bit 작음)

Solution Method ① Exponential Windows

마지막 사각형 전까지 1의 개수 완벽하게 측정 가능, but 나머지 측정 (최대 50% error)

사각형 안엔  $2^k$  개의 1의 개수 저장공간  $O(\log^k N)$  1의 균등하게 분배된다는 가정

DGIM Method

$\Rightarrow$  window size 가변적 (1의 개수에 따라 window size 변화)

Timestamps  $\Rightarrow$  위치에 대한 정보,  $N$ 에 비례  $O(\log_2 N)$  bit 필요

Buckets  $\Rightarrow$  end 지점, 1의 개수 저장  $O(\log N)$  bit

$2^n$  저장 아니냐  $n$ 만 저장  $\therefore O(\log \log N)$

제약조건: 1개와 2개 같은 사이즈 bucket 존재, overlap  $X$ ,  $t \uparrow \rightarrow size \downarrow$ , end-time  $> N$  일때 사라짐

update: 0 이 들어오면 무시, 1 들어오면 새로운 bucket 생성

개수 세는 법: 마지막 bucket 빼고 모두 더함. 마지막 bucket의 반만 더함

예제: 원도유개수  $r$  개일때 Error  $O(\frac{1}{r})$

Bloom Filter: spam msg 판단

특징: FP (긍정 예측, 사실 부정)은 있지만 FN 없다.

method: hash function 이용하여 Bit 값 채움

ex)  $h_1(x) = x \% 11$  (odd),  $h_2(x) = x \% 11$  (even)  $\leftarrow x = \text{binary}$ , 오른쪽부터 곱함

input = 159 = 10011111 짝: 0111  $\% 11 = 7$  홀: 1011  $\% 11 = 0$  1000001000

$B[h(a)] = 1 \Rightarrow$  setting process

118 = 1110110  $h_1(y) = 3, h_2(y) = 5$  둘다 0  $\therefore$  리스트 안에 없다 = spam

Analysis: Hash function 하나 넣을때 Bit에 1이 채워질 확률

$= 1 - e^{-\frac{m}{n}} \approx$  FP 확률  $m = \text{Input 아이템}, n = \text{Target} = \text{bits} = \text{bucket}$

if hash function  $k \rightarrow$  검증도  $k \rightarrow$  1이 들어갈 확률  $1 - e^{-\frac{km}{n}}$

$FP = (1 - e^{-\frac{km}{n}})^k \therefore FP \downarrow$

Counting Distinct Elements

input  $x \rightarrow$  hash  $\rightarrow$  binary (ex. 1100)  $\Rightarrow r(a) = 2$  (오른쪽부터 0의 개수)

$\max(r(a)) = R \quad 2^R = \text{unique value counts}$

$e^{-m 2^{-R}} \rightarrow 2^R \approx m$

Counting Moments:  $\sum (m_i)^k \leftarrow$  적률

$k=2$  이면 얼마나 분포가 불균등한지 나타냄  $\rightarrow$  uneven

AMS Method: 기대값으로 추정

$\Rightarrow n(2C-1) \therefore n(C^k - (C-1)^k) \quad n = \text{stream 개수}, C = \text{count}$

$\Rightarrow$  실제론  $n$ 을 모르기 때문에 factor로 둔다.

Counting items  $\Rightarrow$  Exponentially Decaying Windows

$\Rightarrow$  element 당 점수 매긴다.  $\sum_{i=1}^{\infty} \delta (1-C)^{t-i} \Rightarrow$  데이터가 새로 들어올수록 먼저 들어온 데이터에 대한 중요도  $\downarrow$ , 새롭게 값이 들어오면  $+1$

$\Rightarrow$  가장 많음이 들어온 것은 weight가 가장 큰 것.  $\text{sum}(\text{weight}) = \frac{1}{c}$   
element 개수  $< \frac{2}{c}$

$$\text{Bayes} = P(x|y) = \frac{P(x)P(y|x)}{P(y)}$$

$$\text{chain rule} = P(a,b,c) = P(A|B \cap C) \times P(B|C) \times P(C)$$

$$\text{Entropy} = H(x) = -p \log p - (1-p) \log (1-p)$$

$$\begin{aligned} \text{Cross Entropy} &= H(P) + D_{KL}(P||Q) \\ &= H(P) + E_{x \sim P} [\log P(x) - \log Q(x)] \end{aligned}$$