

# lecture 5 - Similarity Hashing

## Hash function

- $h: X \rightarrow [1..N]$  (typically, sometimes  $\mathbb{R}$ )
- useful as random number generator replacer  
(instead of rand, rand, rand use hash(1), hash(2)...) )
- alternatively also seed(x), rand()
- (but this is very expensive)

## Useful things to do with a hash fun

- 1) random value given  $\frac{h(x)}{N} \in [0, 1]$
- 2) Randomly hashed Parity ( $h(x)$ )  
is  $\{0, 1\}$  or  $\{\pm 1\}$  with prob.  $\frac{1}{2}$
- 3) Basic requirement  $\Pr_h \{h(x) = c\} = \frac{1}{N}$   
k-wise pairwise indep  $\Pr_h \{h(a_1) = c_1, \dots, h(a_n) = c_n\} = \frac{1}{N^k}$   
for  $\{a_i\}$  distinct
- 4) for simplicity (unrealistic) use ideal hash fun.

## Permutations from hashes (aka Feistel hash)

$(x, y) \rightarrow (y, x \text{ XOR } h(y))$   
repeat 3x to get crypt. strong hash

ubyg - Rack off Pseudo random Permutation Generator

$\Rightarrow x \rightarrow \text{Crypt}(x)$  is 1:1

## Locality Sensitive Hash func

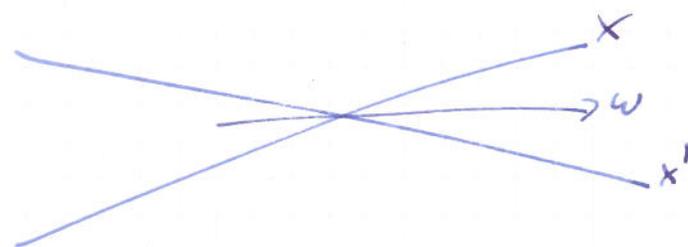
map  $x \rightarrow h(x)$  such that  
similarity between  $(x, x')$  is preserved.  
e.g. whp.

$$\left| \frac{d(h(x), h(x'))}{D(x, x')} - 1 \right| \leq \epsilon$$

not possible for all  $D(x, x')$  but  
available e.g. for Euclidean metrics

example for sets:  $M \times M$  for  $\|ij - i'j'\|$   
(this is the Indyk & Motwani Example)

## Similarity Estimation Techniques from Random Algorithms (N. Charlier '02)



- 1) draw  $w_i \sim \mathcal{N}(0, \sigma^2 \cdot \mathbb{I}_d)$
- 2) map  $x \rightarrow \text{sgn}(\langle w_1, x \rangle, \dots, \langle w_n, x \rangle)$   
 $=: h(x)$

Theorem:  $\mathbb{E}_h \frac{\pi}{2m} \|h(x) - h(x')\|_1$

$$= \arccos \langle x, x' \rangle$$

for  $\|x\| = \|x'\| = 1$

## Lecture 5 - Similarity Hashing

Estimator for inner product:

$$\langle x, x' \rangle \approx \|x\| \cdot \|x'\| \cdot \frac{1}{2^m} \|h(x) - h(x')\|$$

This is much cheaper:

→ modern CPUs can do 64-256 bit

Hamming distance caps. in 1 clock cycle ( $\frac{1}{2}$  ns)

→ only need to store 1 float ( $\|x\| \|x'\|$ ), 64 bit hash regardless of the dimensionality of the data

Theorem:

$$\Pr \left\{ \left| \langle x, x' \rangle - \cos \frac{\pi}{2^m} \|h(x) - h(x')\| \right| > \epsilon \right\} \leq 2 \exp(-\epsilon^2 / \pi^2)$$

Proof: need the Chernoff inequality (like Hoeffding)

→ replace 1 hash function

→ argument doesn't change by more than  $\frac{\pi}{2^m}$

→ cos doesn't change more either

↓ but  $\frac{\pi}{2^m} \cos \frac{\pi}{2^m} \|h(x) - h(x')\| \neq \langle x, x' \rangle$

→ but by inverting the condition we can get it

Definition: LSH scheme with disto:

$$\text{Sim}(x, x') := \Pr_{h \in \mathcal{H}} \left\{ h(x) = h(x') \right\}$$

(2)

Lemma:  $\text{Sim}(x, x')$  is a kernel

Proof:  $\delta(h(x), h(x'))$  is a kernel

Lemma:  $d(x, x') := 1 - \text{Sim}(x, x')$

satisfies the triangle inequality

Proof: For any given  $h$  we have

$$\left[ 1 - \delta(h(x), h(x')) \right] + \left[ 1 - \delta(h(x'), h(x'')) \right] - \left[ 1 - \delta(h(x), h(x'')) \right]$$

$$= 1 + \delta(h(x), h(x'')) - \delta(h(x), h(x')) - \delta(h(x'), h(x'')) \geq 0$$

Taking expectation over  $h$

Example: Sim Hash (aka dense)

$$x \rightarrow \text{sgn}(\langle w, x \rangle)$$

Example: Minwise Hash (Broder)

$$\text{may set } \pi^{-1}(\min(\pi(S))) =: h_\pi(S)$$

⇒ selects a random element from  $S$

$$\Pr_{\pi} \left\{ h_\pi(S) = h_\pi(S') \right\} = \frac{|S \cap S'|}{|S \cup S'|}$$

# lecture 5 - Similarity Hashing

Proof:



- Smallest element in  $S \cup S'$  is with prob  $\frac{|S \cap S'|}{|S \cup S'|}$  in  $|S \cap S'|$
- Prob that  $h(S) = h(S')$  only if the element is in both sets

## Detour: Shingles (Broder et al.)

- instead of picking single element pick  $k$  smallest entries
- not quite correct but good enough

Reber: Conditional Random Sampling (Li et al) (use the hash ID, too)

## Application: Duplicate detection on webpages

- 1)  $(h_i(x), docid) \rightarrow$  sort by  $h_i(x)$
- 2) query  $(docid, docid)$  pairs to find collision candidate
- 3) Count  $(docid, docid)$  to confirm links  
 $\downarrow$   $O(\# \text{duplicates})$  for large clusters

## Back to hash functions

- we have good  $sim(x, x') = Pr_h \{h(x) = h(x')\}$
- can we get a binary  $h(x)$ ?

Defin  $sim'(x, x') := Pr_{h, b} \{b(h(x)) = b(h(x'))\}$

$$= \frac{1}{2} + \frac{1}{2} \frac{\delta(h(x), h(x'))}{sim(x, x')}$$

↑  
random match for diff. hashes

Corollary: we can embed  $sim(x, x')$  on cube

## Application: Fast EX Sampling (Alford, Paris, Sankar)

The problem: in cluster we have to eval

$$p(y|x) \propto p(x|\theta_y) \cdot p(y|\theta)$$

$$= \exp(\langle \ell(x), \theta_y \rangle - g(\theta_y)) \cdot p(y|\theta)$$

Computing  $\langle \ell(x), \theta_y \rangle$  is fatal:

$$O(d \cdot k)$$

↑      ↖

$10^4$  docs for  $10^4 - 10^5$  clusters  
 (ag. docs)

CPU cost is  $\sim 10^9$  FLOPS for 1 dot  
per sampling round. This is fatal  $\leq$

Variant A) approximate by Sim hash  
and use Chernoff bound for  
rejection sampler

$\Rightarrow$  This is terrible since  $\exp(\epsilon)$   
where  $\epsilon = O\left(\frac{1}{\alpha}\right)$



Variant B) Metropolis-Hastings scheme

have  $p(x)$  to sample

but can only afford  $q(x'|x)$

$\Rightarrow$  accept with prob  $\min\left(1, \frac{p(x') q(x|x')}{p(x) q(x'|x)}\right)$

$\Rightarrow$  Special case  $q(x'|x) = q(x')$  to obtain

$$\min\left(1, \frac{p(x') q(x)}{p(x) q(x')}\right)$$

$\rightarrow$  So use approx  $\frac{\|x\| \|x'\| |h(x) - h(x')|}{3}$  as proposal

3 FLOPS  $\Rightarrow 3 \cdot 10^5$  FLOPS  
per proposal

$\rightarrow$  good acceptance ratio  $\Rightarrow$  fast mixing

Application: Similarity search

$\rightarrow$  hashes  
 $\rightarrow$  permute  
 $\rightarrow$  sort  
 $\rightarrow$  find neighbors

} repeat and find  
best match