Lecture III Finding Nearest Neighbors in High Dimensions

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

CSE 547/STAT 548 Winter 2022

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Winter 2022 1/

э

イロト イヨト イヨト イヨト



- Distance functions
- Locality Sensitive Hashing
 - Hash functions and hash tables
 - What is Locality Sensitive Hashing
 - LSH functions from random projections
 - Approximate r-neighbor retrival by LSH
- 4 K-D trees, Ball trees etc.
- Big data and the curse of dimensionality
- Finding similar documents
 Min-Hash

Reading MMDS Ch.: 3. Finding similar items HTF Ch.:, Murphy Ch.: **Reading:** Lecture 16 notes by Moses Charikar, section 3.2; optionally Cormen, Leiserson, Rivest, Stein "Introduction to Algorithms", chapter on hashing.

Thanks to mmds.com (Leskovec, Rajaraman and Ullman) and randorithms.com (Ben Coleman)

2/

イロト イボト イヨト イヨト



Marina Meila (UW)

CSE 547/STAT 548 Winter 2022

3/



Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Winter 2022

4/























10 nearest neighbors from a collection of 20,000 images Marina Melia (UV) III NN in High Dimensions CSE 547/STAT 548 Winter 2022























Marina Meila (UW) III NN in High Dimensions CSE 547/STAT 548 Winter 2022

A Common Metaphor

- Many problems can be expressed as finding "similar" sets:
 - Find near-neighbors in <u>high-dimensional</u> space
 Examples:
- Examples:
 - Pages with similar words
 - For duplicate detection, classification by topic
 - Customers who purchased similar products
 - Products with similar customer sets
 - Images with similar features
 - Users who visited similar websites



The problem: finding neighbors in high dimensions

- Given \mathcal{D} of size *n* in \mathbb{R}^d , and given a query point x find the neighbors of x in \mathcal{D}
 - here: all neighbors in radius r
 - sometimes the k nearest-neighbors
 - sometimes just 1 neighbor
- query point can be in D, e.g. in clustering, dimension reduction, or not (e.g. retrieval, image completion)
- $n \ll 10^6$ and $d > 10^2$
- Brute force (suppose we need neighbors of all $x_i \in \mathcal{D}$)
 - compute time \$\mathcal{O}(n^2d)\$ Too large!
- Can we do it exactly in subquadratic time? Probably NO
 - [if the SETH (Strong Exponential Time Conjecture) holds]
- Rephrased problem: find approximate nearest neighbors
 - e.g. if x has neighbor $x' \in \mathcal{D}$ at distance r, return an $x'' \in \mathcal{D}$ at distance $\leq cr$
 - with c>1 some constant, and w.h.p.¹, usually measured by a confidence δ
 - we measure performance of algorithm as function of (c, r, δ)

¹with high probability

イロト イボト イヨト イヨト

Distance and similarity functions

Distances

- Euclidean $x, x' \in \mathbb{R}^d$, $d_{Euclid}(x, x') = ||x x'|| = \sqrt{x^T x + (x')^T x' 2x^T x'}$ • L1 (Manhattan) $x, x' \in \mathbb{R}^d$ $d_{11}(x, x') = ||x - x'||_1$
- E1 (Wainfattan) $x, x \in \mathbb{R}$ $a_{L1}(x, x) = ||x x||_1$ • Hamming $x, x' \in \{0, 1\}^d$ $d_H(x, x') = x^T x + (x')^T x' - 2x^T x' = \#x + \#x' - 2\#(x \cap x')$

Similarities

• cosine
$$x, x' \in \mathbb{R}^d$$
 or $\{0, 1\}^d \cos(x, x') = \frac{x^T x'}{\sqrt{(x^T x)((x')^T x')}}$
• Jaccard $x, x' \in \{0, 1\}^d$ $J(x, x') = \frac{\#(x \cap x')}{\#(x \cup x')} = \frac{x^T x'}{x^T x + (x')^T x' - x^T x'}$

- Note that if $x, x' \in \{0, 1\}^d$ they can be seen as indicator functions for subsets of 1 : n.
- Hence $x^T x' = \#(x \cap x')$ represents the cardinality of the intersection of sets given by x, x'• All distances above are metrics.

《日》《御》《日》《日》 - 日

Hash functions and hash codes

Let the data space be \mathbb{R}^d , and assume some fixed probability measure on this space.

- A family of hash functions is a set $\mathcal{H} = \{h : \mathbb{R}^d \to \{0,1\}\}$ with the following properties
 - For each h, $Pr[h(x) = 1] \approx \frac{1}{2}$
 - The binary random variables defined by the functions in *H* are mutually independent. (Or, if *H* is not finite, a "not too large" random sample of such random variables is mutually independent.)
- Let $h_{1:k}$ be a mutually independent subset of \mathcal{H} . We call

$$g(x) = [h_1(x) h_2(x) \dots h_k(x)] \in \{0,1\}^k$$
(1)

the hash code of x.

- Note that the codes g(x) are (approximately) uniformly distributed; the probability of any $g \in \{0,1\}^k$ is about $\frac{1}{2^k}$.
- Useful hash functions must be fast to compute.

イロト 不得 トイヨト イヨト

Hash tables

- A hash table \mathcal{T} is a data structure in which points in \mathbb{R}^d can be stored in such a way that
 - All points with the same code g are in the same bin denoted by T_g. The table need not use space for empty bins.
 - **②** Given any value $g \in \{0, 1\}^k$, we can obtain a point in \mathcal{T}_g or find if $\mathcal{T}_g = \emptyset$ in constant time (independent of the number of points *n* stored in \mathcal{T}). Some variance of back tables rature all points in \mathcal{T}_g or g, and list in constant time
 - Some versions of hash tables return all points in $\mathcal{T}_{g},$ e.g., as a list, in constant time.
 - (a) It is usually assumed that storing a point x with given code g(x) in a hash table is also constant time.
- Hence, using a hash table to store an x or to retrieve something, involves computing k hash functions, then a constant-time access to \mathcal{T} .
- When x' ≠ x and g(x') = g(x) we call this a collision. In some applications (not of interest to us), collisions are to be avoided.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Hashing vs. Locality Sensitive Hashing (LSH)



by Ben Coleman randorithms.com

Marina Meila (UW)

III NN in High Dimensions

Locality Sensitive Hash Functions and Codes

• A hash function h is locality sensitive iff for any $x, x' \in \mathbb{R}^d$

$$Pr[h(x) = h(x')] \ge p_1$$
 when $||x - x'|| \le r$ (2)

$$Pr[h(x) = h(x')] \le p_2$$
 when $||x - x'|| \ge cr$ (3)

with p_1, p_2, r and c > 1 fixed parameters (of the family \mathcal{H}) and $p_1 > p_2$. • W.I.o.g., we set $p_1 = p_2^{\rho}$ for some $\rho < 1$.



by Ben Coleman randorithms.com

CSE 547/STAT 548 Winter 2022 13

イロト イヨト イヨト

LSH functions

• A locality sensitive *h* makes a weak distinction between points that are close in space vs. points that are far away. A hash code *g* from locality sensitive hash functions sharpens this distinction, in the sense that the probability of far away points colliding can be made arbitrarily small.

$$p_{bad} = Pr[g(x) = g(x') | ||x - x'|| > cr] \le p_2^k$$
 (4)

- Assume x is not in \mathcal{T} ; for any $x' \in \mathcal{D}$ which is far from x,the probability that x' collides with x is $\leq p_{bad}$.
- We construct ${\mathcal T}$ so that $p_{bad} \leq \frac{1}{n}$ for n the sample size. For this we need Exercise (in Homework 1)

$$k = \frac{\ln n}{-\ln p_2} \Rightarrow p_{bad} \le \frac{1}{n}$$
 (5)

• Suppose $x' \in \mathcal{T}$ is "close" to x. What is the probability that g(x') = g(x)?

$$p_{good} = p_1^k = p_2^{\rho k} = \frac{1}{n^{\rho}}$$
 (6)

This is the probability that the bin $\mathcal{T}_{g(x)}$ contains x'.

- h depends on the distance d
- h and g sometimes depend on r

3

イロト 不得 トイヨト イヨト

How to find good hash functions?

- We need large families of h functions
- that are easy to generate randomly
- and fast to compute for a given x
- Generic method to obtain them: random projections

イロト イヨト イヨト イヨト

LSH function for Hamming distance

• $\mathcal{H} = \{h_j = bit_j(x), j = 1 : d\}$

• a random $h \in \mathcal{H}$ samples a random bit of x

Collision probability

$$p_1(x,x') = 1 - \frac{d_H(x,x')}{d}$$
 (7)

To bit sample, randomly choose an index This is sensitive to Hamming distance



by Ben Coleman

randorithms.com

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Winter 2022 16

イロト イヨト イヨト イヨト

randorithms.com

LSH function for Euclidean and L1 distance

• project x on a random line, round to multiples of r

$$h_{w,b}(x) = \lfloor \frac{w^T x + b}{r} \rfloor \tag{8}$$

- If $w \sim Normal(0, I_d)$, hash function for Euclidean distance
- If $w \sim Cauchy(0,1)^d$, hash function for L1 distance
- Collision probability (p = 2 for Normal, p = 1 for Cauchy)

$$p_1(x, x') = \text{deterministic function of } \|x - x'\|_p$$
(9)

• Hash function space \mathcal{H}_r is infinite, and depends on r



by Ben Coleman

イロト イヨト イヨト

CSE 547/STAT 548 Winter 2022

Analysis of projection on a random vector

- Data are $x \in \mathbb{R}^d$ as usual.
- Define $h_{w,b}: \mathbb{R}^d \to \mathbb{Z}$ by

$$h_{w,b}(x) = \lfloor \frac{w^T x + b}{r} \rfloor$$
(10)

with r > 0 a width parameter, $w \in \mathbb{R}^d, b \in [0, r)$.

- Intuitively, x is "projected" on w^2 , then the result is quantized into bins of width r, with a grid origin given by b.
- The family of hash functions is $\mathcal{H}_r = \{h_{w,b}, w \in \mathbb{R}^d, b \in [0, r)\}.$
- Sampling \mathcal{H}_r : $w \sim Normal(0, I_d)$, $b \sim uniform[0, r)$.
 - Because the Normal distribution is a stable distribution, this ensures that $w^T x$ is distributed as $Normal(0, ||x||^2)$. Exercise Verify this
 - Hence $w^T x w^T x'$ is distributed as $Normal(0, ||x x'||^2)$. Exercise Verify this
 - Moreover, if hash functions are sampled independently from $\mathcal{H}_{r,}(and nothing is known about x)$ then $h_{w,b}(x)$, $h_{w',b'}(x)$ are independent random variables. Exercise Prove this

²w is not necessarily unit length

イロト 不得 トイヨト イヨト

LSH function for angles

• project x on a random line, take the sign

$$h_{w,b}(x) = \operatorname{sign}(w^T x) \tag{11}$$

Collision probability

$$p_1(x,x') = 1 - \frac{\theta(x,x')}{\pi}$$
 (12)

• Hash function space ${\mathcal H}$ is infinite



by Ben Coleman randorithms.com



by Ben Coleman

< □ > < □ > < □ > < □ > < □ >

Clustering LSH

- \$\mathcal{H} = \{h = k(x), for some clustering of data\}\$
- h takes values in 1 : K
- This is a data dependent hash function family
- Clustering can be K-means, min-diameter, hieararchical . . .
- No theoretical guarantees, but works well in practice



by Ben Coleman

randorithms.com

CSE 547/STAT 548 Winter 2022 20

イロト イヨト イヨト イヨト

Approximate *r*-neighbor retrival by LSH

Input \mathcal{D} set of *n* points, *L* mutually independent hash codes $g_{1:L}$ of dimension *k*. dexing Construct *L* hash tables $\mathcal{T}^{1:L}$, each storing \mathcal{D} . trieval Given *x*

- compute g(x)
 for j = 1, 2, ... L if the bin T^j_{g(x)} ≠ Ø
 return some (all) x' from it.
 g stop if a single neighbor is wanted.
- Some analysis. We set $L = n^{\rho}$
 - Indexing time $\propto kn^{\rho+1}$
 - Retrieval time $\propto kn^{
 ho}$
 - Space used $\propto kn^{
 ho+1}$
 - For each $x' \in \mathcal{D}$ close to x, the probability that x' is NOT returned for any $j \in 1 : L$ is

$$(1 - \frac{1}{n^{\rho}})^{n^{\rho}} \approx \frac{1}{e} \tag{13}$$

This can be made arbitrarily small by multiplying L with a constant.

• For each $x' \in \mathcal{D}$ far from x, the probability that x' is NOT returned for any $j \in 1 : L$ is

$$\left(1-\frac{1}{N}\right)^{n^{\rho}} \approx \left(\frac{1}{e}\right)^{1/n^{1-\rho}} \approx \frac{1}{e^0} = 1 \tag{14}$$

• Hence, we are almost sure not to return a far point, and have a significant probability to return a close point when one exists, if no points neither far nor close are in the data. This is

Marina Meila (UW)

CSE 547/STAT 548 Winter 2022

²¹

Heuristics for neighbors in high-dimensions

- typically a form of hierarchical clustering
- K-D tree for low dimensions (but observed to work well in high dimensions too)
- Ball tree for high dimensions





< □ > < □ > < □ > < □ > < □ >





Marina Meila (UW)

III NN in High Dimensions









Marina Meila (UW)

III NN in High Dimensions



Marina Meila (UW)

III NN in High Dimensions





Marina Meila (UW)

III NN in High Dimensions

K-D Tree construction

node k:

- $b_{1:d}^{\min}$, $b_{1:d}^{\max}$ min, max of box in each dimension
- j_{\max} , Δ_{\max} = argmax, $\max_{j} \{ b_{i}^{\max} b_{i}^{\min}, j = 1 : d \}$ the largest dimension of the box
- n_k, \bar{x}_k, \ldots number of points in node, mean, other statistics
- if k is **leaf** then \mathcal{D}_k an array of the data under this node
- pointers p_k, l_k, r_k to parent and children nodes

Algorithm Split-Node(k)

- It is assumed that k is leaf, hence $I_k, r_k =$ null
- **()** Create new leaf nodes I_k , r_k children of k and set k as their parent
- 2 Let $b^* = (b_{i_{max}}^{max} + b_{i_{max}}^{min})/2$
- Create empty sets $\mathcal{D}_{l_k}, \mathcal{D}_{r_k}$
- For $i = 1 : n_k$
 - if $x_{i,j_{\text{max}}} < b^*$ then move x_i from \mathcal{D}_k to \mathcal{D}_{l_k} ; else move x_i to \mathcal{D}_{r_k}
 - update n_{l_k} , n_{r_k} and the other statistics as needed
 - update $b_{l_{1},r_{1}}^{\max,\min}$
- **5** Update $\Delta_{l_k, \max}, j_{l_k, \max}$ and $\Delta_{r_k, \max}, j_{r_k, \max}$

3

イロト イヨト イヨト イヨト

Searching for *r*-neighbors with K-D Tree

- Denote by Node_k the *d*-dimensional box $[b_1^{\min}, b_1^{\max}] \times \dots [b_d^{\min}, b_d^{\max}]$
- When is $B_r(x) \cap \operatorname{Node}_k \neq \emptyset$?
 - x close to a corner: closest corner is $c = [\min\{|b_j^{\min} x_j|, |b_j^{\max} x_j|\}]_{j=1:d}$
 - x is interior or close to a face: $x_j \in [b_i^{\min}, b_j^{\max}]$ if $j \neq j_0$, and $x_j \in [b_i^{\min} r, b_j^{\max} + r]$ for $j = j_0$
- When is $Node_k \subset B_r(x)$?
 - furthest corner is $c' = [\max\{|b_j^{\min} x_j|, |b_j^{\max} x_j|\}]_{j=1:d}$
 - if $||x c'|| \le r$ then all $Node_k \subset B_r(x)$

Retrieving all points in $\mathcal{D} \cap B_r(x)$

- Recursively from the root, examine Node_k
- If $B_r(x) \cap Node_k = \emptyset$, return with no output
- Else
- If $Node_k \subset B_r(x)$ output all \mathcal{D}_k and return
- Else examine children of k



イロト 不得 トイヨト イヨト

Task: Finding Similar Documents

- Goal: Given a large number (N in the millions or billions) of documents, find "near duplicate" pairs
 Applications:
 - Mirror websites, or approximate mirrors
 - Don't want to show both in search results
 - Similar news articles at many news sites
 - Cluster articles by "same story"

Problems:

- Many small pieces of one document can appear out of order in another
- Too many documents to compare all pairs
- Documents are so large or so many that they cannot fit in main memory

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Winter 2022

3 Essential Steps for Similar Docs

- 1. Shingling: Convert documents to sets
- 2. *Min-Hashing:* Convert large sets to short signatures, while preserving similarity
- Locality-Sensitive Hashing: Focus on pairs of signatures likely to be from similar documents
 - Candidate pairs!

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

CSE 547/STAT 548 Winter 2022

The Big Picture



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Marina Meila (UW)

III NN in High Dimensions

- no

29

CSE 547/STAT 548 Winter 2022



Shingling

Step 1: Shingling: Convert documents to sets

Marina Meila (UW)

III NN in High Dimensions

Documents as High-Dim Data

- Step 1: Shingling: Convert documents to sets
- Simple approaches:
 - Document = set of words appearing in document
 - Document = set of "important" words
 - Don't work well for this application. Why?

Need to account for ordering of words! A different way: Shingles!

Define: Shingles

- A k-shingle (or k-gram) for a document is a sequence of k tokens that appears in the doc
 - Tokens can be characters, words or something else, depending on the application
 - Assume tokens = characters for examples
- Example: k=2; document D₁ = abcab Set of 2-shingles: S(D₁) = {ab, bc, ca}
 - Option: Shingles as a bag (multiset), count ab twice: S'(D₁) = {ab, bc, ca, ab}

Compressing Shingles

- To compress long shingles, we can hash them to (say) 4 bytes
- Represent a document by the set of hash values of its k-shingles
 - Idea: Two documents could (rarely) appear to have shingles in common, when in fact only the hashvalues were shared

Example: k=2; document D₁= abcab
 Set of 2-shingles: S(D₁) = {ab, bc, ca}
 Hash the singles: h(D₁) = {1, 5, 7}

Similarity Metric for Shingles

- Document D₁ is a set of its k-shingles C₁=S(D₁)
- Equivalently, each document is a 0/1 vector in the space of k-shingles
 - Each unique shingle is a dimension
 - Vectors are very sparse
- A natural similarity measure is the

Jaccard similarity:

 $sim(D_1, D_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

III NN in High Dimensions

34

CSE 547/STAT 548 Winter 2022

Motivation for Minhash/LSH

- Suppose we need to find near-duplicate documents among N = 1 million documents
- Naïvely, we would have to compute pairwise Jaccard similarities for every pair of docs
 - N(N − 1)/2 ≈ 5*10¹¹ comparisons
 - At 10⁵ secs/day and 10⁶ comparisons/sec, it would take 5 days
- For N = 10 million, it takes more than a year...

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Motivation for Minhash/LSH

- Suppose we need to find near-duplicate documents among N = 1 million documents
- Naïvely, we would have to compute pairwise Jaccard similarities for every pair of docs
 - N(N − 1)/2 ≈ 5*10¹¹ comparisons
 - At 10⁵ secs/day and 10⁶ comparisons/sec, it would take 5 days
- For N = 10 million, it takes more than a year...

Min-Hash – Motivation

- Denote $S = \{ \text{ space of } k \text{-shingles } (k \text{-grams}) \}$
- $|S| = |alphabet|^k$ HUGE!
- document $\rightarrow c \in \{0,1\}^{|\mathcal{S}|}$ sparse!
- Similarity(document, document') = J(c, c') Jaccard

$$J(c,c') = rac{\#(c\cap c')}{\#(c\cup c')}$$

- Wanted compress $c \rightarrow x$, so that
 - $x \in \mathbb{Z}_+^L$ with $L \ll |\mathcal{S}|$
 - Jaccard is preserved (approximately), i.e.

$$J(c, c') \approx \frac{\#\{x_l = x'_l\}}{L}$$
 (15)

(fraction of equal elements in signatures approximates Jaccard)

- x is called signature of c
- How? Min-Hash
- Why not random bit hashing?

Min-Hash

Min-Hashing Example

Note: Another (equivalent) way is to

store row indexes:





CSE 547/STAT 548 Winter 2022

Min-Hash

The Min-Hash Property

- Choose a random permutation π
- <u>Claim</u>: Pr[h_π(C₁) = h_π(C₂)] = sim(C₁, C₂)
 Why?
 - Let **X** be a doc (set of shingles), **y** ∈ **X** is a shingle
 - Then: Pr[π(y) = min(π(X))] = 1/|X|

- Let \boldsymbol{y} be s.t. $\pi(\boldsymbol{y}) = \min(\pi(C_1 \cup C_2))$
- Then either: $\pi(y) = \min(\pi(C_1))$ if $y \in C_1$, or $\pi(y) = \min(\pi(C_2))$ if $y \in C_2$
- So the prob. that **both** are true is the prob. $\textbf{y} \in C_1 \cap C_2$
- $\Pr[\min(\pi(C_1))=\min(\pi(C_2))]=|C_1 \cap C_2|/|C_1 \cup C_2|= sim(C_1, C_2)$

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

27

0

0

1

0

0

1

One of the two cols had to have

1 at position v

0

0

1

0

1

0

Min-Hash high-level summary

- Choose a family of hash functions $\mathcal{H} = \{h_{\pi}\}$
 - where π are permutatons of ${\cal S}$
 - $h_{\pi}(c) \in \{0, 1, \dots |\mathcal{S}| 1\}$
 - $h_{\pi}(c)$ = number 0's at the beginning of $\pi(c)$ = location of 1st 1 in $\pi(c)$ (zero-indexed)
- so that

$$Pr[h_{\pi}(c) = h_{\pi}(c')] = J(c,c')$$
 for all π, c, c' (Min-Hash Property)

- Choose L random permutations $\pi_{1:L}$
- Map c vectors to x by

$$x(c) = [h_{\pi_1}(c), h_{\pi_2}(c), \dots h_{\pi_L}(c)]$$

• Approximate J(c, c') by averaging

$$J(c,c') = \frac{1}{L} \sum_{l=1}^{L} \mathbb{1}_{[x_l = x'_l]}$$

Marina Meila (UW)

CSE 547/STAT 548 Winter 2022 40

э

イロト イヨト イヨト イヨト

Input matrix (Shingles x Documents)

Min-Hashing Example

Permutation π

2

3

1

5

2 4 6 3 2 6 6

1	0	1	0
1	0	0	1
0	1	0	1
0	1	0	1
0	1	0	1
1	0	1	0
1	0	1	0

Signature matrix M



Similarities:

	1-3	2-4	1-2	3-4
Col/Col	0.75	0.75	0	0
Sig/Sig	0.67	1.00	0	ο

J. Leskovec, A. Raiaraman, J. Ullman: Mining of Massive Datasets. http://www.mmds.org

CSE 547/STAT 548 Winter 2022 41

Finding similar documents: Summary

3

イロト 不得 トイヨト イヨト