





NN in high dimensions LSH

Project · published twa

Lecture III Finding Nearest Neighbors in High Dimensions

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

CSE 547/STAT 548 Spring 2025

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Spring 2025 1/

イロト イヨト イヨト イヨト

2 Big Data



🚺 Motivation – finding similar items 🛛 🗲 🛨

Locality Sensitive Hashing

- Hash functions and hash tables
- What is Locality Sensitive Hashing
- LSH functions from random projections
- Approximate r-neighbor retrival by LSH
- K-D trees, Ball trees etc.
- Big data and the curse of dimensionality
- 6 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)

イロト イヨト イヨト



Marina Meila (UW)

CSE 547/STAT 548 Spring 2025

3 /



quory (image) = x

wanted {x'} heighbors" of x 2

4 /















5 /









10 nearest neighbors from a collection of 20,000 images Marina Mella (UW) III NN in High Dimensions CSE 547/STAT 548 Spring 2025























Haring Meila (UW) Haring 2025

A Common Metaphor

- Many problems can be expressed as finding "similar" sets:
- Find near-neighbors in <u>high-dimensional</u> space
 Examples:
 - Bagos with similary
 - 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 \times find the neighbors of \times 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 <u>clustering</u>, 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 O(n²d) 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 $x' \in D$ at distance r, return an $x'' \in D$ 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)

Approx NN alps -> index D r usually fixed _____ distance x ∈ D ⇒ neigh(x) 3 X

¹with high probability



Note that if x, x' ∈ {0,1}^d they can be seen as indicator functions for subsets of 1 : n.
Hence x^Tx' = #(x ∩ x') represents the cardinality of the intersection of sets given by x, x'
All distances above are metrics.



$$\begin{array}{l} x \in \{0, 1\}^{d} = \lambda i \in [n], \ x_{i} = 1 \\ [n] = 1 : n \equiv \{1, 2, \cdots, n\} \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \alpha_{i} + \beta_{i}} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \beta_{i}} = \frac{1}{3} = \frac{1}{3} \Rightarrow d_{3} = 1 - J(x_{i}, x_{i}) \\ \hline J(x_{i}, x_{i}) = \frac{1}{1 + \alpha_{i} + \beta_{i}} = \frac{1}{3} = \frac$$

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Spring 2025

9/

Hash functions and hash codes

hj(x)€h0,1g hj_he +je ;=

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

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

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.





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 \hat{T} or g, as a 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.

э.

Hashing vs. Locality Sensitive Hashing (LSH)



by Ben Coleman randorithms.com

Marina Meila (UW)

III NN in High Dimensions

CSE 547/STAT 548 Spring 2025 12

< □ > < □ > < □ > < □ > < □ >

Locality Sensitive Hash Functions and Codes

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



III NN in High Dimensions

CSE 547/STAT 548 Spring 2025 13

э

Intuition $P_r[h(x)=h(x^{t}), ||x-x^{t}||=k] = p(r)$ about (2), (3)





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)

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

• 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

$$P_{n} = \Pr\left[h\left(x\right) = h\left(x'\right)\right] \qquad ||x - x'|| < \mathbb{R} \qquad \text{good}$$

$$P_{n} = \Pr\left[-u - 1\right] \qquad ||x - x'|| \ge cr \qquad c > 1 \quad bad$$

$$g \in \left\{0, 12^{k} \qquad \Pr\left[g(x) = g(x')\right] \mid ||x - x'|| > cr \qquad \right] = \Pr\left[2 = \Pr\left[bad\right]$$

$$Want \qquad P_{bad} = \frac{1}{n} \implies \frac{1}{n} - \Pr\left[2 \implies k = \frac{\log n}{\log \frac{1}{22}}\right]$$

$$P_{good} = \Pr\left[g(x) = g(x')\right] \mid ||x - x'|| < \mathbb{R} \qquad = \Pr\left[2 = \frac{1}{n}\right] = \frac{1}{n^{2}}$$

$$P_{good} = P_r \left[g(x) = g(x') \right| ||x - x'|| < |z] = p_1^k = p_2^{sk} = \left(\frac{1}{n}\right)^s = \frac{1}{n^s}$$

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

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)
$$\mathcal{P}_{L}\left[h(x) \neq h(x')\right]$$





by Ben Coleman

randorithms.com

x = 20,13

< □ > < □ > < □ > < □ > < □ >

LSH function for Euclidean and L1 distance



 $h \cdot \mathbb{R}^d \to \mathbb{Z}$

17

Sample $W_{i:e} \sim N(o_i I_d)$ $b_{i:e} \sim umif[o_i r] \implies h_{i:e}$ $g(x) = [h_i(x) \cdots h_e(x)]$ 1,51 W2,62 $g(x) = \left[h_{w_{1}b_{1}}^{(x)} h_{w_{2}b_{2}}^{(x)} \right]$