

Hashing for Data with some Distance Structure

(Near Neighbor Search and Locality Sensitive Hashing).

Often data has structure and hashing tends to remove it —

we want that changing ~~the~~ a single letter in a word causes it to hash to a very different place — look independent!

But at other times we may want to preserve some structure

— we're aiming for compression without losing the notion of similarity in the original space

Example: we're given strings and want to find pairs of "similar" strings
e.g. in document de-duplication.

or given a corpus S of strings, ~~find~~ and a query q ,
find strings in S that are "close to" q .

e.g. in web searching.

Note: exact problem is easier, using just standard hashing.

— hash all strings in the corpus into a table

— when query comes, look at strings that hash to location $h(q)$ and explicitly compare to them.

We're interested in approximate similarity.

Let's look at an example.

Suppose data are strings (maybe files of text)

need a notion of similarity (metric space)

(a) could imagine the underlying points belong to $\{A-Z a-z 0-9\}^*$
and distance is Edit distance.

+ Natural

- difficult to manipulate, do searches.

(b) strings are bit vectors (say in ASCII representation).

distance is Euclidean or Hamming distance

+ Easier to manipulate / understand

- no more intuition for what is similar or different.
- different length strings make it awkward

(c) Convert document into a "bag of words".

- Say 10^5 words in English, keep count of how many times each word appears in document / string.

+ each document $\in \mathbb{Z}^{10^5}$

- lost order information / proximity information.

- taking two copies of same doc gives very different point

- stop words can overwhelm all else

- Normalise, say by total # terms in doc (TF/IDF term frequency / inverse doc. freq.).
and ~~also~~ multiply by how common/rare word is in the corpus.

Say $\left(\frac{\text{freq of } w}{\text{total \# of words in doc}} \right) \times \log \left(\frac{\text{\# docs in corpus}}{\text{\# docs with word } w} \right)$

\uparrow TF \uparrow IDF

Also to keep some local structure :-

shingling — take k consecutive words in doc, call that a shingle
(eg. 3-shingles)

Now do some things for these shingles. ← also do for images

Maybe drop stop words.

+ hundreds of useful cool tricks.



But eventually have say

- (a) a set of items ← a bag of words, want to see if they are similar, or
- (b) a set of vectors ← bag put in some geometric representation.

What's the notion of similarity? or distance?

~~Let's~~ Let's look at the vector case, say.

• Each item is a vector v_i

• Distance is Hamming distance :- $\text{dist}(v_i, v_j) = \# \text{coordinates where they differ.}$

or Euclidean distance :- $\|v_i - v_j\| = \sqrt{\sum_{k=1}^d |v_{ik} - v_{jk}|^2}$

or L_1 -distance

(Manhattan/taxicab distance)

$$\|v_i - v_j\|_1 = \sum_{k=1}^d |v_{ik} - v_{jk}|$$

or millions of others.

Let's try to solve the (approximate) near neighbors problem

Given a metric space (distances, say \mathbb{R}^d)

and a corpus $S \subseteq \mathbb{R}^d$ of points

Create a data structure. such that it answers

queries of the form "return a near-neighbor of point q ".

i.e. ~~q~~ return a point $a \in S$ ^{answer}

$$\text{st } d(q, a) \leq c \times \min_{a^* \in S} d(q, a^*)$$

^{nearest neighbor distance}
^{approx factor. say 2}
^{say (1+ε).} ← today's lecture

Challenge: S is large. say $n = |S| = 10^{10}$ points

Strawman solution:

no data structure. Brute force search each time.

so $O(nd)$ time per query

^{# points in S} ^{$O(d)$ to compute distance, say}

Should do better.

Indeed, will use Locality Sensitive hashing

— similar items fall "close" to each other.

dissimilar items do not.

Want: data structure with small storage space, small preprocessing time and small query time

Solution #2 (Strouman++)

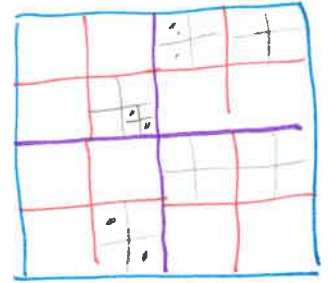
Use k-d trees (or variant)

Partition bounding box for all points in S into half along each coordinate

and so on, until each smallest box contains a single point.

For each box keep a representative point

Given query, go down the hierarchy, keeping closest representative at each level ~~find~~ look at all neighboring boxes as well.



So takes time $2^d \cdot \log\left(\frac{d_{\max}}{d_{\min}}\right)$ for queries

\uparrow # of neighboring boxes. \uparrow # levels of recursion

Curse of dimensionality!

- high dimensional point sets take $\exp(\text{dimension})$ query time

want to do better than that.

Note: have gone from $O(nd)$ query time, no preprocessing

to $O(2^d \cdot \log\left(\frac{d_{\max}}{d_{\min}}\right))$, $O(n \cdot \log\left(\frac{d_{\max}}{d_{\min}}\right))$ space

2 preprocessing.

Curse of dimensionality:

if n items in d -dimensional space, searching for near neighbors in this space typically takes $2^{\Omega(d)}$ time

Exponential dependence. 😞 if we do it naively.

(E.g. if we do it via quad-trees).
↑ again naively.

So let's be smarter.

Idea #1: Doubling Search.

Given S and parameter r , solve the problem

Ball search

— given query q , if $\exists a \in S$ at distance $d(q, a) \leq r$
return some a' at distance $\leq (1+\epsilon)r$.

if no $a \in S$ at distance $\leq r$ from query q ,
return whatever you want.

Now can try all values of r , get an answer a_r for that value

↑
powers of $(1+\epsilon)$
i.e. $r_i = (1+\epsilon)^i$

- check the distance (q, a_r)
- return the best of those.

if actual closest point is $a^* \in S$ at distance $\in [(1+\epsilon)^i, (1+\epsilon)^{i+1})$

then will get a good answer for $r = (1+\epsilon)^{i+1}$

that has distance from q being $\leq (1+\epsilon)^{i+2}$

$$\begin{aligned} &\leq (1+\epsilon)^2 \cdot d(q, a^*) \\ &\leq (1+2\epsilon) \text{dist}(q, a^*). \end{aligned}$$

OK: how to solve the doubling search version?

Given S , radius r , answer (approx) queries:

if \exists a point at distance $\leq r$ from q
return a point at dist $\leq 2r$ (say)

Here's idea:

Use hash function h : metric space \rightarrow buckets

such that

LSH. $\left[\begin{array}{l} \text{if } \text{dist}(x, y) \leq r \\ \text{dist}(x, y) \geq 2r \end{array} \right. \Rightarrow \begin{array}{l} P_c[h(x) = h(y)] \geq P_{\text{close}} \\ P_r[h(x) = h(y)] \leq P_{\text{far}} \end{array}$

Ideally: $P_{\text{close}} \gg P_{\text{far}}$

Use these to hash into buckets.

So $E[\text{#junk points in } h(q) \text{ bucket}] \leq \sum_{x: \text{dist}(q, x) \geq 2r} P_{\text{far}} \cdot n$

≤ 1 if $P_{\text{far}} = 1/n$.

But then maybe P_{close} is small as well.

So repeat $\cong 1/P_{\text{close}}$ times, i.e. use $1/P_{\text{close}}$

independent
many hash functions
to increase
probability we see
a close point.

Let's do a concrete example:

Want to distinguish points in $\{0,1\}^d$ ← Boolean cube with Hamming distance.
that are at distance $\leq r$ or at distance $\geq 2r$.

Try #0: $h(x) =$ pick a random coordinate i
and output x_i

$$\begin{aligned} \text{if } \text{dist}(x,y) \leq r &\Rightarrow \Pr[h(x) = h(y)] \geq 1 - \frac{r}{d} = P_{\text{close}} \\ \text{dist}(x,y) \geq 2r &\Rightarrow \Pr[h(x) = h(y)] \leq 1 - \frac{2r}{d} = P_{\text{far}} \end{aligned}$$

Very small gap between P_{close} & P_{far} .

Try #1: Parallel Repetition

$h(x) =$ pick l random coordinates with repetition
output those l coordinate values

$$\begin{aligned} \text{if } \text{dist}(x,y) \leq r &\Rightarrow \Pr[h(x) = h(y)] \geq \left(1 - \frac{r}{d}\right)^l = P_{\text{close}} \\ &\geq 2r \Rightarrow && \leq \left(1 - \frac{2r}{d}\right)^l = P_{\text{far}} \end{aligned}$$

As always: if $r \ll d$ (say) then $1 - \frac{r}{d} \approx e^{-r/d}$ etc.

$$\Rightarrow P_{\text{close}} \approx e^{-r/d \cdot l}$$

$$P_{\text{far}} \approx e^{-2r/d \cdot l} = \frac{1}{n} \text{ (say)}$$

$$\Rightarrow \frac{2r \cdot l}{d} = \log_2 n \Rightarrow l = \frac{d \cdot \log_2 n}{2r}$$

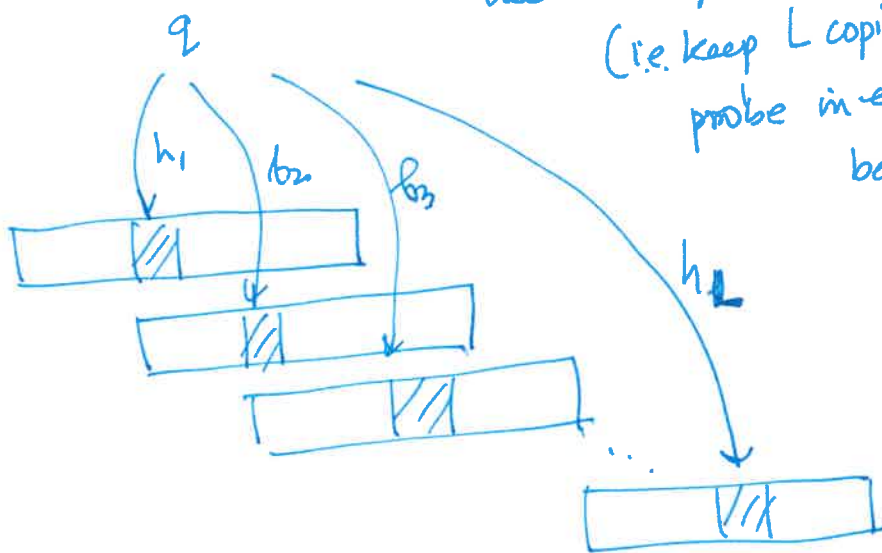
$$\Rightarrow P_{\text{close}} \cong e^{-r/d \cdot L} = e^{-\frac{1}{2} \log_2 n} = \frac{1}{\sqrt{n}}.$$

\Rightarrow if x and q are indeed at distance $\leq r$

then $\Pr[h(x) = h(q)] \geq \frac{1}{\sqrt{n}}$ \leftarrow small probability, but much better than P_{far} .

So final ingredient: ~~Parallel~~ Serial Repetition

use L independent hash functions
(i.e. keep L copies of data structures,
probe in each and return the best)



if we set $L = \sqrt{n} \log n$

$$\begin{aligned} \Rightarrow \Pr[\exists \text{ table } i \text{ st } h_i(q) = h_i(x)] &\geq 1 - \left(1 - \frac{1}{\sqrt{n}}\right)^L \\ &\cong 1 - \left(e^{-1/\sqrt{n}}\right)^L \\ &= 1 - e^{-\log n} = 1 - \frac{1}{n}. \end{aligned}$$

😊

So in summary:

used a weak hash function h with $P_{\text{close}}, P_{\text{far}}$ not very different

then reduced P_{far} by parallel repetition

and handled the "small P_{close} value" problem

by serial repetition.

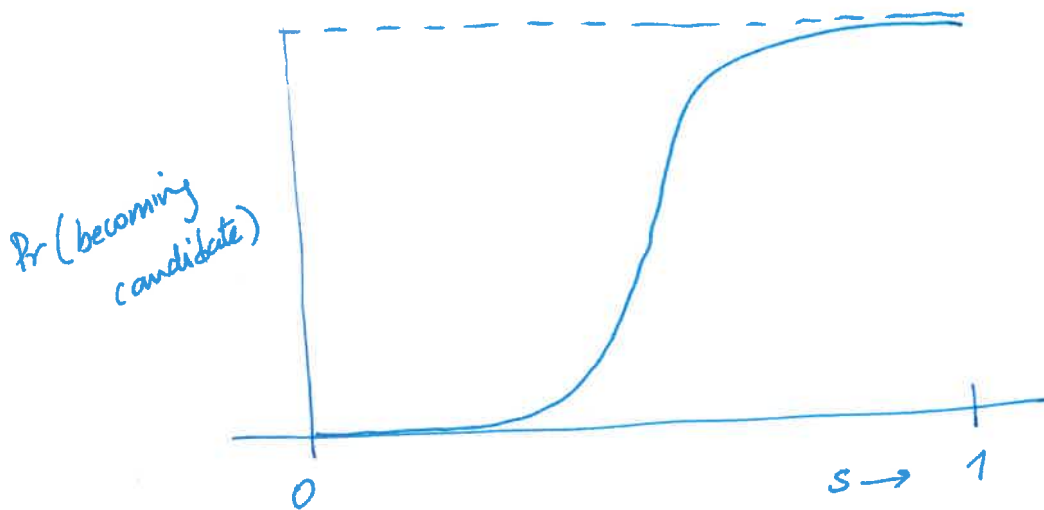
Note: if $s = \Pr[h(x) = h(q)]$ for a single hash function.

then $\Pr(h^t(x) = h^t(q))$ for the function after parallel rep]
~~is~~ $= s^t$.

$$\Rightarrow \Pr[h^t(x) \neq h^t(q)] = (1-s^t)$$

$$\Rightarrow \Pr[x \text{ is not seen in all } t \text{ copies of the data structure}] = (1-s^t)^t$$

$$\Rightarrow \Pr[x \text{ is one of the candidates for query } q] = 1 - (1-s^t)^t$$



$$1 - (1-s^t)^t \cong \frac{1}{2} \text{ at } (1-s^t)^t \cong \frac{1}{2} \text{ or } s \cong \left(\frac{1}{t}\right)^{1/2}$$

~~is~~ \Rightarrow

by choosing t, ϵ carefully

- can get space usage $n^{2-\epsilon}$ and query time n^{ϵ} in worst case
- in practice does quite a bit better. ↶ for ϵ vs $(1-\epsilon)r$

Observe: steps II (parallel repetition)

and III (serial repetition) are completely generic.


So: if we can get a small gap between collision probability in the close & far cases \Rightarrow can amplify the success.

Saw Step I for bit strings and Hamming distance.

Can do this for other metric spaces:-

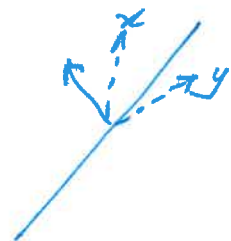
• Unit vectors and cosine distances

Elements are unit vectors.

$d(x, y) = \theta$ θ is angle between them 

hash function: pick a random hyperplane.

defines a halfspace. $x \in \text{Halfspace} \Rightarrow h(x) = 1$
else $h(x) = 0$.



$$\Pr[h(x) \neq h(y)] = \theta/\pi$$

• Sets and Jacquard distance

$$d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

"minhash"
↓

hash function: pick a random permutation on elements.

$h(S)$ = first element (in this permutation) in S .

$$\Pr[h(A) = h(B)] = \frac{|A \cap B|}{|A \cup B|} \Rightarrow \Pr[h(A) \neq h(B)] = d(A, B).$$

Take-aways

- Can easily extend the basic ideas behind hashing to structured data.

→ Call two items similar if $d(x, y) \leq r$
dissimilar if $d(x, y) \geq cr$.

Grey zone in the middle (don't care either way)

LSH idea

then can ask for

P_s [similar items collide] \geq large

P_d [dissimilar items collide] \leq small

And then amplify ~~the~~ correctness by repetition.

- Can do this for many metric spaces

- Hamming Cube (in lecture)

- Euclidean distances (embed into Ham Cube)

- Cosine distances

- Jacquard distances

← omitted, see MMDS book.

↑ depends on dimension to some extent
(not in same way as k-d tree)

- Popular and commonly used heuristic

(good in theory, even better in practice).

Next time: Reduce the ambient dimension of your data.