CS-GY 6763: Lecture 6 Near-neighbor search in high dimensions

NYU Tandon School of Engineering, Prof. Christopher Musco

Nrg

Dimensionality reduction: Given vectors \mathbf{x} , \mathbf{y} , compute small space compressions $C(\mathbf{x})$ and $C(\mathbf{y})$ that can be used to estimate the distance or similarity between \mathbf{x} and \mathbf{y} .

EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Distributional JL Lemma)

Let Π be a random matrix that compresses to $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ rows. Then with probability $(1 - \delta)$:

$$(1-\epsilon)\|\mathbf{x}-\mathbf{y}\|_2^2 \le \|\mathbf{\Pi}\mathbf{x}-\mathbf{\Pi}\mathbf{y}\|_2^2 \le (1+\epsilon)\|\mathbf{x}-\mathbf{y}\|_2^2$$



DIMENSIONALITY REDUCTION FOR JACCARD SIMILARITY



KEY APPLICATION: MODERN VECTOR SEA SEARCH



KEY APPLICATION: MODERN VECTOR SEA SEARCH



All modern vector search systems use "fancier" versions of methods studied in this class:

- Product Quantization
 PCA-based methods

Dimensionality reduction methods are typically paired with vector indexing methods.

Goal of Dimensionality Reduction: Reduce **dependence on** *d* in *O*(*nd*) search cost. Reduce space complexity.

Goal of Vector Indexing: Reduce **dependence on** *n* in *O*(<u>n</u>*d*) search cost. Often at the cost of added space complexity.

This problem can already be solved in low-dimensions using space partitioning approaches (namely kd-trees)

0(m)



ISSUE WITH KD-TREES



ISSUE WITH KD-TREES



HIGH DIMENSIONAL NEAR NEIGHBOR SEARCH

Only been attacked much more recently:

- Locality-sensitive hashing [Indyk, Motwani, 1998]
- Spectral hashing [Weiss, Torralba, and Fergus, 2008]
- Vector quantization [Jégou, Douze, Schmid, 2009]
- Graph-based vector search [Malkov, Yashunin, 2016, Subramanya et al., 2019]

Key ideas behind all of these methods:

- 1. Trade worse space-complexity +(preprocessing time) for better time-complexity. I.e., preprocess database in data structure that uses $(\Omega(\underline{n}))$ space.
- 2. Allow for approximation.

INTUITIVELY WHY DO PREPROCESSING AND SPACE HELP?

Question: Suppose you want to search over points in $[-1, 1]^d$ and to achieve accuracy ϵ . I.e., for a given $\mathbf{y} \in [-1, 1]^d$, you want to find $\tilde{\mathbf{q}}$ with $\|\mathbf{y} - \tilde{\mathbf{q}}\|_2 \le \min_i \|\mathbf{y} - \mathbf{q}_i\|_2 + \epsilon$.

Can you construct a data structure that supports *we* time search but uses <u>exponential space</u>?



LOCALITY SENSITIVE HASH FUNCTIONS

Let $h : \mathbb{R}^d \to \{1, \dots, m\}$ be a random hash function. We call h <u>locality sensitive</u> for similarity function $\underline{s}(\underline{q}, \underline{y})$ if $\Pr[h(\mathbf{q}) == h(\underline{y})]$ is:

• Higher when **q** and **y** are more similar, i.e. $\underline{s(q, y)}$ is higher.

Lower when **q** and **y** are more dissimilar, i.e. s(q, y) is lower.





LOCALITY SENSITIVE HASH FUNCTIONS

LSH for s(q, y) equal to Jaccard similarity: $\int e \frac{1}{20} \frac{1}{13} d$

- Let $c: \{0,1\}^d \rightarrow [0,1]$ be a single instantiation of MinHash.
- Let $g: (0,1) \rightarrow \{1,\ldots,m\}$ be a uniform random hash function.
- Let $h(\mathbf{q}) = g(c(\mathbf{q}))$.



$$Pr\left[c(q) := c(2)\right] = J(q, 3)$$

LSH for Jaccard similarity:

- Let $c: \{0,1\}^d \rightarrow [0,1]$ be a single instantiation of MinHash.
- Let $g : [0, 1] \rightarrow \{1, \dots, m\}$ be a uniform random hash function.

$$Let h(\mathbf{x}) = g(c(\mathbf{x})). \qquad g(c(\mathbf{x})) \cdot g(c(\mathbf{x})) \cdot f(c(\mathbf{x})) \cdot f(c(\mathbf{x})) \cdot f(c(\mathbf{x})) \cdot f(\mathbf{x}) + f(\mathbf{x}) +$$

Basic approach for LSH-based near neighbor search in a database.

Pre-processing:

- Select random LSH function $h: \{\underline{0,1}\}^d \to 1, \dots, m$.
- Create table T with m = O(n) slots.¹
- For i = 1, ..., n, insert (\mathbf{q}) into $T(h(\mathbf{q}_i))$.



Query:

- Want to find near neighbors of input $\mathbf{y} \in \{0,1\}^d$.
- Linear scan through all vectors $\mathbf{g} \in T(\underline{h}(\mathbf{y}))$ and return any that are close to \mathbf{y} . Time required is $O(\underline{d} \cdot |T(\underline{h}(\mathbf{y})|)$.

¹Enough to make the O(1/m) term negligible.

NEAR NEIGHBOR SEARCH



Two main considerations:

- **False Negative Rate**: What's the probability we do not find a vector that is close to **y**?
- **False Positive Rate**: What's the probability that a vector in $T(h(\mathbf{y}))$ is not close to \mathbf{y} ?

A higher false negative rate means we miss near neighbors.

A higher false positive rate means (increased runtime) – we need to compute $S(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to \mathbf{y} .

Note: The meaning of "close" and "not close" is application dependent. E.g. we might specify that we want to find anything with Jaccard similarity > .4, but not with Jaccard similarity < .2.

Let's use Jaccard similarity as a running example. We will discuss LSH for inner product/Euclidean distance as well. Suppose the nearest database point **q** has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we do not find q?

$$Pr\left(\mathcal{G}\in\mathcal{T}(h(\gamma))\right)$$

$$= Pr\left(h(q) = h(\gamma)\right) = J(\gamma, \beta) = .\Upsilon$$
Folze myotron prod. = 1-. $\Upsilon = .\mathcal{L}$

REDUCING FALSE NEGATIVE RATE



 $h_1(3) \dots h_r(3)$

Pre-processing:

- Select t independent LSH's $(h_1, \ldots, h_t : \{0, 1\}^d \rightarrow 1, \ldots, m.$
- Create tables T_1, \ldots, T_t , each with *m* slots.
- For i = 1, ..., n, j = 1, ..., t,
 - Insert \mathbf{q}_i into $T_j(h_j(\mathbf{q}_i))$.

REDUCING FALSE NEGATIVE RATE

Query:

- Want to find near neighbors of input $\mathbf{y} \in \{0, 1\}^d$.
- Linear scan through all vectors in $T_1(h_1(\mathbf{y})) \cup T_2(h_2(\mathbf{y})) \cup \dots, T_t(h_t(\mathbf{y})).$

J- 10

Suppose the nearest database point **q** has
$$J(\mathbf{y}, \mathbf{q}) = .4$$
.

What's the probability we find q?

$$-(1-.4)^{+} - (1-.6^{+})^{+}$$



Suppose there is some other database poin z with J(y, z) = .2. What is the probability we will need to compute J(z, y) in our

hashing scheme with one table? I.e. the probability that **y** hashes into at least one bucket containing **z**.

In the new scheme with t = 10 tables?

$$|-(1.2)^{+} = |-.8^{+} = \frac{99\%}{1}$$

$$\frac{1}{12} = 1 - .8^{+} = \frac{99\%}{12}$$

(89%)

Change our locality sensitive hash function. Tunable LSH for Jaccard similarity: • Let $c_1, \ldots, c_r : \{0, 1\}^d \rightarrow [0, 1]$ be independent random MinHash's. • Let $g: [0,1]^r \to \{\underline{1,\ldots,m}\}$ be a uniform random hash function. • Let $h(\mathbf{x}) = g(c_1(\mathbf{x}), \ldots, c_r(\mathbf{x})).$.12 .27 .5 ... r "bands" $c_1(\mathbf{g}) \in \mathbf{g}$ q C,(9) Juwitus un rend - un $g(c_1(\mathbf{q}), c_2(\mathbf{q}), ..., c_r(\mathbf{q}))$ 23 2 m 1

<u>Tunable</u> LSH for Jaccard similarity: $\int \frac{p}{r} \int \frac{p}{r} \int \frac{1}{r} \int \frac{p}{r} \int \frac{1}{r} \frac{1}$

- Choose parameter $r \in \mathbb{Z}^+$.
- Let c_1 , \mathcal{W}_{A} : $\{0,1\}^d \rightarrow [0,1]$ be random MinHash. ((3) (3))
- Let $g:[0,1]^r \to \{1,\ldots,m\}$ be a uniform random hash function.
- Let $h(\mathbf{x}) = g(\underline{c_1}(\mathbf{x}), \dots, \underline{c_r}(\mathbf{x})).$ $f: J(\mathfrak{g}, \mathfrak{Y}) \subseteq \mathcal{V}$
- If $J(\underline{q}, \underline{y}) = \underline{v}$, then $\Pr[h(\mathbf{q}) == h(\underline{y})] = \underbrace{\underline{v}}_{-} \cdot \cdot ((-\underline{v}) \cdot \underline{y})$

verligible.

(2(8)= (2()) and

TUNABLE LSH





Full LSH cheme has two parameters to tune:



Effect of **increasing number of bands** *r* on:



Probability we check **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:



Probability we check **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$\approx 1 - (1 - v^r)^t$$



r = 5, t = 40



Probability we check **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^r)^t$$



Increasing both *r* and *t* gives a steeper curve.

Better for search, but worse space complexity.

FIXED THRESHOLD



Space complexity: 40 hash tables $\approx 40 \cdot O(n)$. $\downarrow O(uc)$ Directly trade space for fast search.

Possible to prove concrete worst-case results for distance functions that satisfy triangle inequality.

Theorem (Indyk, Motwani, 1998. Point Location in Ball) Fix a distance R. If there exists some q with $\|\mathbf{q} - \mathbf{y}\|_0 \le R$, return a vector $\mathbf{\tilde{q}}$ with $\|\mathbf{\tilde{q}} - \mathbf{y}\|_0 \le C \cdot R$ in:

- Time: $O(n^{1/C})$.
- Space: $O(n^{1+1/C} + nd)$.

 $\|\boldsymbol{q}-\boldsymbol{y}\|_0=$ "hamming distance" = number of elements that differ between \boldsymbol{q} and $\boldsymbol{y}.$

If there is no point at distance *R*, algorithm does not need to return anything.

To obtain a nearest-neighbor search algorithm build multiple data structures for exponentially growing distances:

R 2R 4R 8R ...

Search from most accurate level to least accurate.



To obtain a nearest-neighbor search algorithm build multiple data structures for exponentially growing distances:

R 2R 4R 8R ...

Search from most accurate level to least accurate.



To obtain a nearest-neighbor search algorithm build multiple data structures for exponentially growing distances:

R 2R 4R 8R ...

Search from most accurate level to least accurate.



Total number of levels = $O(\log(d_{\max}/d_{\min}))$, where $d_{\max} = \max_{i,j} ||\mathbf{q}_i - \mathbf{q}_j||$ and $d_{\min} = \min_{i,j} ||\mathbf{q}_i - \mathbf{q}_j||$. d_{\max}/d_{\min} is called the **dynamic range**.

Theorem (Indyk, Motwani, 1998)

Let \mathbf{q} be the closest database vector to \mathbf{y} . Return a vector $\mathbf{\tilde{q}}$ with $\|\mathbf{\tilde{q}} - \mathbf{y}\|_0 \le C \cdot \|\mathbf{q} - \mathbf{y}\|_0$ in:

• Time:
$$\tilde{O}(n^{1/C})$$
. $O(\overline{C})$ $C = 2$

• Space:
$$\tilde{O}(n^{1+1/c} + nd)$$
. $O(\eta^{1.5} + nd)$

Similar results can be proven for other metrics, including Euclidean distance. But you need a good LSH function.

Good locality sensitive hash functions exists for other similarity measures.

Cosine similarity $\cos(\theta(\mathbf{x}, \mathbf{y})) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$:



Cosine similarity is natural "inverse" for Euclidean distance when $\|\mathbf{x}\|_2^2 = \|\mathbf{y}\|_2^2 = 1$ (often the case for ML-based embeddings).

$$(o_{2}(\theta(x, z)) = \langle x, y \rangle = (-\frac{|x-y||_{r}}{2})$$

$$(|x-y||_{r} = ||x||_{r}^{2} + ||_{r}||_{r}^{2} - 2\langle x, y \rangle = 2 - 2\langle x, y \rangle$$

LSH functions also exist for Euclidean distance, but are a bit more complex to describe/analyze. See [Andoni, Indyk, 2006] if you are interested. Locality sensitive hash for cosine similarity:

- Let $g \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $\underline{f}: \{-1,1\} \rightarrow \{\underline{1},\ldots,m\}$ be a uniformly random hash function.
- $\underline{h} : \mathbb{R}^d \to \{1, \dots, m\}$ is defined $h(\mathbf{x}) = \underline{f(\operatorname{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))}.$

If $cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $Pr[h(\mathbf{x}) == h(\mathbf{y})]$?

SIMHASH ANALYSIS IN 2D



SimHash can be banded, just like our MinHash based LSH function for Jaccard similarity:

- Let $\underline{g}_1, \dots, \underline{g}_r \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $f: \{-1,1\}^r \to \{1,\ldots,m\}$ be a uniformly random hash function.

•
$$h : \mathbb{R}^d \to \{1, \dots, m\}$$
 is defined
 $h(\mathbf{x}) = \underbrace{f([\operatorname{sign}(\langle \mathbf{g}_1, \mathbf{x} \rangle), \dots, \operatorname{sign}(\langle \mathbf{g}_r, \mathbf{x} \rangle)])}_{\Pr[h(\mathbf{x}) == h(\mathbf{y})]} \approx \underbrace{\left(1 - \frac{\theta}{\pi}\right)^r}_{\Pr[h(\mathbf{x}) == h(\mathbf{y})]}$

SIMHASH ANALYSIS IN 2D



SIMHASH ANALYSIS 2D



 $\Pr[sign(\langle g, x \rangle) == sign(\langle g, y \rangle)] = probability x and y are on the same side of hyperplane orthogonal to g.$

SIMHASH ANALYSIS HIGHER DIMENSIONS



There is always some <u>rotation matrix</u> **U** such that **Ux**, **Uy** are spanned by the first two-standard basis vectors and have the same cosine similarity as **x** and **y**.

SIMHASH ANALYSIS HIGHER DIMENSIONS



There is always some <u>rotation matrix</u> **U** such that **x**, **y** are spanned by the first two-standard basis vectors.

Note: A rotation matrix U has the property that $U^T U = I$. I.e., U^T is a rotation matrix itself, which reverses the rotation of U.

Claim:

$$\begin{aligned} \Pr[\operatorname{sign}(\langle g, \mathbf{x} \rangle) &== \operatorname{sign}(\langle g, \mathbf{y} \rangle) = \Pr[\operatorname{sign}(\langle g, \mathbf{U}\mathbf{x} \rangle) == \operatorname{sign}(\langle g, \mathbf{U}\mathbf{y} \rangle)] \\ &= \Pr[\operatorname{sign}(\langle g[1, 2], (\mathbf{U}\mathbf{x})[1, 2] \rangle) == \operatorname{sign}(\langle g[1, 2], (\mathbf{U}\mathbf{y}[1, 2] \rangle)] \\ &= 1 - \frac{\theta}{\pi}. \end{aligned}$$

The first step is the trickiest here. Why does it hold?

LSH is widely used in practice, but is starting to get replaced by other methods. Most of these are <u>data dependent</u> in some way.

Starting point: Think of LSH as a randomized space-partitioning method.



In practice, we can often get partitions with better <u>margin</u> but partitioning in a data-dependent way.

Common approach: Split data using k-means clustering.



NEAREST-NEIGHBOR SEARCH IN PRACTICE

Common approach: Split data using k-means clustering.



Main approach behind "k-means tree" and "inverted file index" based near-neighbor search methods like Meta's FAISS library and Google's SCANN.

NEAREST-NEIGHBOR SEARCH IN PRACTICE

New kid on the block: Graph-based nearest neighbor search.



Idea behind methods like NSG, HNSW, DiskANN, etc. Inspired by Milgram's famous "small-world" experiments from the 1960's.

Can we better explain the success of data-dependent nearest-neighbor search methods?

