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

NYU Tandon School of Engineering, Prof. Christopher Musco

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

Lemma (MinHash)

Let C be a length $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ MinHash sketch. Then with probability $(1 - \delta)$, we can return an estimate \tilde{J} based on C(**x**) and C(**y**) with:

$$J(\mathbf{x}, \mathbf{y}) - \epsilon \leq \tilde{J} \leq J(\mathbf{x}, \mathbf{y}) + \epsilon.$$



KEY APPLICATION: MODERN VECTOR SEA SEARCH



Cost of naive algorithm is O(nd).

Dimensionality reduction reduces search cost to *O*(*nk*) and reduces space requirements.



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

- Quantized JL/SimHash
- b-bit MinHash

- 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*(*nd*) 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).



Search time is roughly $O(d \cdot \log n \cdot 2^d)$), which is only sublinear for $d = o(\log n)$.

ISSUE WITH KD-TREES



ISSUE WITH KD-TREES



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:

- Trade worse space-complexity + preprocessing time for better time-complexity. I.e., preprocess database in data structure that uses Ω(n) space.
- 2. Allow for approximation.

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 *O*(1) time search but uses <u>exponential space</u>?

Let $h : \mathbb{R}^d \to \{1, \dots, m\}$ be a random hash function.

We call h <u>locality sensitive</u> for similarity function s(q, y) if Pr [h(q) == h(y)] is:

- Higher when \mathbf{q} and \mathbf{y} are more similar, i.e. $s(\mathbf{q}, \mathbf{y})$ is higher.
- Lower when **q** and **y** are more dissimilar, i.e. *s*(**q**, **y**) is lower.



LSH for *s*(**q**, **y**) equal to 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 $h(\mathbf{q}) = g(c(\mathbf{q}))$.



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 $h(\mathbf{x}) = g(c(\mathbf{x}))$.

 $\mathsf{lfJ}(\mathsf{q},\mathsf{y})=\mathsf{v}_{\mathsf{,}}$

 $\Pr[h(q) == h(y)] =$

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

Pre-processing:

- Select random LSH function $h: \{0,1\}^d \rightarrow 1, \dots, m$.
- Create table T with m = O(n) slots.¹
- For $i = 1, \ldots, n$, insert \mathbf{q}_i 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 $\mathbf{q} \in T(h(\mathbf{y}))$ and return any that are close to \mathbf{y} . Time required is $O(d \cdot |T(h(\mathbf{y})|)$.

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 <u>is close</u> to **y**?
- False Positive Rate: What's the probability that a vector in T(h(y)) is not close to 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 \mathbf{q} has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we do not find q?

REDUCING FALSE NEGATIVE RATE



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 $T_j(h_j(\mathbf{q}_i))$.

Query:

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

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

What's the probability we find q?

(10, 99%)

Suppose there is some other database point **z** with $J(\mathbf{y}, \mathbf{z}) = .2$. What is the probability we will need to compute $J(\mathbf{z}, \mathbf{y})$ 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?

(89%)

Change our locality sensitive hash function.

Tunable LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.
- Let $c_1, \ldots, c_r : \{0, 1\}^d \rightarrow [0, 1]$ be independnt random MinHash's.
- + Let $g: [0,1]^r \to \{1,\ldots,m\}$ be a uniform random hash function.

• Let
$$h(\mathbf{x}) = g(c_1(\mathbf{x}), \dots, c_r(\mathbf{x})).$$



Tunable LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.
- Let $c_1, \ldots, c_r : \{0, 1\}^d \rightarrow [0, 1]$ be random MinHash.
- + Let $g: [0,1]^r \to \{1,\ldots,m\}$ be a uniform random hash function.
- Let $h(\mathbf{x}) = g(c_1(\mathbf{x}), \dots, c_r(\mathbf{x})).$

If J(q, y) = v, then $\Pr[h(q) == h(y)] =$

TUNABLE LSH



Full LSH cheme has two parameters to tune:



Effect of **increasing number of tables** t on:

False Negatives

False Positives

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

False Negatives

False Positives

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



r = 5, t = 5

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 $J(\mathbf{q}, \mathbf{y}) = v$:

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



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.

Use Case 1: Fixed threshold.

- Shazam wants to find match to audio clip **y** in a database of 10 million clips.
- There are 10 true matches with J(y, q) > .9.
- There are 10,000 <u>near matches</u> with $J(y, q) \in [.7, .9]$.
- All other items have J(y, q) < .7.

With r = 25 and t = 40,

- + Hit probability for J(y,q) > .9 is $\gtrsim 1-(1-.9^{25})^{40}=.95$
- + Hit probability for J(y,q) \in [.7, .9] is $\lesssim 1-(1-.9^{25})^{40}=.95$
- + Hit probability for J(y,q) < .7 is $\lesssim 1-(1-.7^{25})^{40}=.005$

Upper bound on total number of items checked:

 $10 + .95 \cdot 10,000 + .005 \cdot 9,989,990 \approx 60,000 \ll 10,000,000.$

Space complexity: 40 hash tables $\approx 40 \cdot O(n)$. 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 $\tilde{\mathbf{q}}$ with $\|\tilde{\mathbf{q}} - \mathbf{y}\|_0 \le C \cdot R$ 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 q be the closest database vector to y. Return a vector \tilde{q} with $\|\tilde{q} - y\|_0 \le C \cdot \|q - y\|_0$ in:

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

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



 $-1 \leq \cos(\theta(\mathbf{x}, \mathbf{y})) \leq 1.$

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).

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

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

Theorem (to be proven): If $cos(\theta(x, y)) = v$, then $\Pr[h(\mathbf{x}) == h(\mathbf{y})] = 1 - \frac{\theta}{\pi} + \frac{\theta/\pi}{m} = 1 - \frac{\cos^{-1}(v)}{\pi} + \frac{\theta/\pi}{m}$ 0.9 0.8 collision probability 2 2 2 2 0.2 0.1 0 L -1 -0.8 -0.6 -0.4 -0.2 0 0.2 0.4 0.6 0.8 cosine similarity

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

- Let $\mathbf{g}_1, \ldots, \mathbf{g}_r \in \mathbb{R}^d$ 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 $h(\mathbf{x}) = f([sign(\langle \mathbf{g}_1, \mathbf{x} \rangle), \dots, sign(\langle \mathbf{g}_r, \mathbf{x} \rangle)]).$

$$\Pr[h(\mathbf{x}) == h(\mathbf{y})] \approx \left(1 - \frac{\theta}{\pi}\right)^r$$

SIMHASH ANALYSIS IN 2D

To prove: $\Pr[h(\mathbf{x}) == h(\mathbf{y})] \approx 1 - \frac{\theta}{\pi}$, where $h(\mathbf{x}) = f(\operatorname{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))$ and *f* is uniformly random hash function.



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 <u>space-partitioning method</u>.



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?

