CS-GY 6763: : Lecture 4 Near neighbor search in high dimensions + locality sensitive hashing

NYU Tandon School of Engineering, Prof. Christopher Musco

Given two length *d* vectors \underline{y} and \underline{q} , construct compact representations (sketches) $\underline{\tilde{y}}$ and $\underline{\tilde{q}}$ such that $\underline{\text{dist}(y, q)}$ can be estimated accurately from $\underline{\tilde{y}}$ and $\underline{\tilde{q}}$.

Each of $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{q}}$ should require $k \ll d$ space.

EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Johnson-Lindenstrauss, 1984)

For any two data points $\mathbf{y}, \mathbf{q} \in \mathbb{R}^d$ there exists a <u>linear map</u> $\Pi : \mathbb{R}^d \to \mathbb{R}^k$ where $\underline{k} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ such that with probability $1 - \delta, \quad (1 - \epsilon) \quad \|\mathbf{x}\|_{\mathbf{L}} \leq \|\mathbf{f}_{\mathbf{X}}\|_{\mathbf{L}} \leq (1 + \epsilon) \quad \|\mathbf{x}\|_{\mathbf{L}}$ $(1 - \epsilon) \|\mathbf{q} - \mathbf{y}\|_2 \leq \|\mathbf{\Pi}\mathbf{q} - \mathbf{\Pi}\mathbf{y}\|_2 \leq (1 + \epsilon) \|\mathbf{q} - \mathbf{y}\|_2.$



Lemma (Johnson-Lindenstrauss, 1984)

For any set of n data points $\mathbf{q}_1, \ldots, \mathbf{q}_n \in \mathbb{R}^d$ there exists a <u>linear map</u> $\Pi : \mathbb{R}^d \to \mathbb{R}^k$ where $k = O\left(\frac{\log(n/\delta)}{\epsilon^2}\right)$ such that with probability $(1 - \delta)$, for all i, j,

$$(1-\epsilon)\|\mathbf{q}_i-\mathbf{q}_j\|_2 \leq \|\mathbf{\Pi}\mathbf{q}_i-\mathbf{\Pi}\mathbf{q}_j\|_2 \leq (1+\epsilon)\|\mathbf{q}_i-\mathbf{q}_j\|_2.$$

Extends to approximating all pairwise distances in a set of *n* vectors via a **union bound**.

JACCARD SIMILARITY

B

Another distance measure (actually a similarity measure) between <u>binary</u> vectors in $\{0, 1\}^d$:

Definition (Jaccard Similarity)

$$J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|} = \frac{\text{\# of non-zero entries in common}}{\text{total \# of non-zero entries}}$$

Natural similarity measure for binary vectors. $0 \le J(q, y) \le 1$.

Can be applied to any data which has a natural binary representation (more than you might think).

$$y = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$q = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

How does **Shazam** match a song clip against a library of 8 million songs (32 TB of data) in a fraction of a second?



Spectrogram extracted from audio clip.

Processed spectrogram: used to construct audio "fingerprint" $\mathbf{q} \in \{0,1\}^d$.

Each clip is represented by a high dimensional binary vector **q**.

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

JACCARD SIMILARITY FOR DOCUMENT COMPARISON



How many words do a pair of documents have in common?

JACCARD SIMILARITY FOR DOCUMENT COMPARISON



How many bigrams do a pair of documents have in common?

- Finding duplicate or new duplicate documents or webpages.
- Change detection for high-speed web caches.
- Finding near-duplicate emails or customer reviews which could indicate spam.

MINHASH



MINHASH

- Choose k random hash functions $(h_1) \dots, h_k : \{1, \dots, n\} \rightarrow [0, 1].$ • For $i \in 1, \dots, k$, • Let $c_i = \min_{j, q_i=1} h_i(j).$
- $C(\mathbf{q}) = [c_1, \ldots, c_k].$



MINHASH ANALYSIS



Proof:

1. For $c_i(\mathbf{q}) = c_i(\mathbf{y})$, we need that $\arg\min_{i \in \mathbf{q}} h(i) = \arg\min_{i \in \mathbf{y}} h(i)$. if $\mathfrak{g} \Rightarrow \mathfrak{f} \mathfrak{s} \mathfrak{f} \mathfrak{s} \mathfrak{f}$

MINHASH ANALYSIS

Claim: $Pr[c_i(q) = c_i(y)] = J(q, y).$



2. Every non-zero index in $\mathbf{q} \cup \mathbf{y}$ is equally likely to produce the lowest hash value. $c_i(\mathbf{q}) = c_i(\mathbf{y})$ only if this index is 1 in <u>both</u> \mathbf{q} and \mathbf{y} . There are $\mathbf{q} \cap \mathbf{y}$ such indices. So:

$$\Pr[c_i(\mathbf{q}) = c_i(\mathbf{y})] = \frac{\mathbf{q} \cap \mathbf{y}}{\mathbf{q} \cup \mathbf{y}} = J(\mathbf{q}, \mathbf{y})$$

Let J = J(q, y) denote the Jaccard similarity between q and y.

Return: $\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})].$

Unbiased estimate for Jaccard similarity:

$$\mathbb{E}\tilde{J} = \frac{l}{k} \sum_{i=1}^{k} \mathbb{E}\left(\int \left(c_{i}(c_{b}) = c_{i}(c_{b}) \right) \right) = \frac{l}{k} \sum_{i=1}^{k} J$$

$$C(q) \boxed{.12} 24 .76 .35 C(y) \boxed{.12} .98 .76 .11} = J$$

The more repetitions, the lower the variance.

MINHASH ANALYSIS

$$\Gamma = \sum_{i=1}^{k} \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})].$$

Estimator: $\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})].$

 $\Gamma = \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})].$

$$\operatorname{Var}[\tilde{J}] = \frac{1}{k^{2}} \sum_{i=1}^{k} \operatorname{Var}\left(\operatorname{I}\left(\operatorname{C}_{i}\left(\operatorname{P}\right) = \operatorname{C}_{i}\left(\operatorname{P}\right)\right)\right) = \frac{1}{k^{2}} \sum_{i=1}^{k} \operatorname{J} - \operatorname{J}^{2}$$

$$= \frac{1}{k} \sum_{i=1}^{k} \operatorname{Var}\left(\operatorname{I}\left(\operatorname{C}_{i}\left(\operatorname{P}\right) = \operatorname{C}_{i}\left(\operatorname{P}\right)\right)\right) = \frac{1}{k^{2}} \sum_{i=1}^{k} \operatorname{J} - \operatorname{J}^{2}$$

Plug into Chebyshev inequality. How large does k need to be so that with probability $> 1 - \delta$:

$$Pr\left(\left|\tilde{J} - \tilde{E}\tilde{J}\right| \geqslant z \cdot \frac{1}{K}\right] \leq \frac{1}{2^{2}}$$

$$E \qquad z \cdot \frac{1}{K} = E \qquad \frac{1}{6} \cdot \frac{1}{K} = E \qquad \frac{1}{6^{2}} \cdot \frac{1}{K}$$

15

MINHASH ANALYSIS

Chebyshev inequality: As long as $k = O\left(\frac{1}{\epsilon^2 \delta}\right)$, then with prob. $1 - \delta$,

$$J(q, y) - \epsilon \leq \tilde{J}(C(q), C(y)) \leq J(q, y) + \epsilon.$$

And \tilde{J} only takes O(k) time to compute! Independent of original fingerprint dimension d.

Can be improved to $\log(1/\delta)$ dependence. Can anyone tell me how? $\frac{1}{2}$



SIMILARITY SKETCHING



BREAK

NEAR NEIGHBOR SEARCH

Common goal: Find all vectors in database $\mathbf{q}_1, \ldots, \mathbf{q}_n \in \mathbb{R}^d$ that are close to some input query vector $\mathbf{y} \in \mathbb{R}^d$. I.e. find all of \mathbf{y} 's "nearest neighbors" in the database.

- The Shazam problem.
- Audio + video search.
- Finding duplicate or near duplicate documents.
- Detecting seismic events.

How does similarity sketching help in these applications?

- Improves runtime of "linear scan" from $O(n_{d})$ to $O(n_{k})$.
- Improves space complexity from O(nd) to O(nk). This can be super important – e.g. if it means the linear scan only accesses vectors in fast memory.

New goal: Sublinear o(n) time to find near neighbors.

This problem can already be solved for a small number of dimensions using space partitioning approaches (e.g. kd-tree).



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

HIGH DIMENSIONAL NEAR NEIGHBOR SEARCH



Only been attacked much more recently:

Ocality-sensitive hashing [Indyk, Motwani, 1998]
 Spectral hashing [Weiss, Torralba, and Fergus, 2008]
 Vector quantization [Jégou, Douze, Schmid, 2009]

Key Insight: Trade worse space-complexity for better time-complexity.

0(nd) 0(nd)

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



$$M = O(N)$$

LSH for Jaccard similarity:

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

If $J(\mathbf{q}, \mathbf{y}) = \mathbf{v}$, $\Pr[h(\mathbf{q}) == h(\mathbf{y})] = \mathbf{V} + (|-\mathbf{v}|) \cdot \frac{1}{m} \quad \mathbf{x} \quad \mathbf{v}$ $1. \quad c(\mathbf{x}) = c(\mathbf{y}) \quad hoppens \quad with \quad prob. \quad \mathbf{v}$ $\mathbf{y} \cdot c(\mathbf{x}) \neq c(\mathbf{y}) \quad hoppens \quad \mathbf{v} \cdot \mathbf{p}. \quad |-\mathbf{v}|$ Basic approach for near neighbor search in a database. $g_{1/\dots -}, g_{n}$ 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 $\mathbf{q} \in T(\underline{h}(\underline{y}))$ 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 $J(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to \mathbf{y} .

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

Suppose the nearest database point **q** has $J(\mathbf{y}, \mathbf{q}) = .4$. What's the probability we 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 $T_j(h_j(\mathbf{q}_i))$.

t= 10 94%

Query:

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

$$L \cdot (T_1(h(5)) \cup T_1(h_1(5)) \dots)$$

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

What's the probability we find q?

$$1 - 6^{+}$$
, 99 for $f = 10$.

89%

Suppose there is some other database point 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**.

+= 10

In the new scheme with t = 10 tables?

Change our locality sensitive hash function.Tunable LSH for Jaccard similarity: \mathcal{N} $\begin{pmatrix} g \\ g \end{pmatrix}$: \mathcal{J} \mathcal{L} \mathcal{N}

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

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



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
$$J(\mathbf{q}, \mathbf{y}) = \mathbf{v}$$
, then $\Pr[h(\mathbf{q}) == h(\mathbf{y})] = \bigvee_{\mathbf{y}} + \underbrace{(I = \mathcal{O}^{r})}_{\mathbf{y}} \times \bigvee_{\mathbf{y}} \nabla$

2. This doesn't heppen 33

TUNABLE LSH



Full LSH cheme has two parameters to tune:



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

False Negatives	False Positives
Deccare	Increase

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



t= 3 78°%

Choose tables *t* large enough so false negative rate to 1%.

Parameter: r = 1.

Chance we find **q** with $J(\mathbf{y}, \mathbf{q}) = \underline{.8}$: $\begin{array}{c} \underbrace{+ : 3}{2} \\ \underbrace{- .2^{+}}{7} & \underline{.99^{\circ}}_{0} \end{array}$

Chance we need to check z with J(y, z) = .4:

Choose tables *t* large enough so false negative rate to 1%.



Chance we need to check z with J(y, z) = .4:

Choose tables t large enough so false negative rate to 1%.

Parameter: r = 5.

Chance we find **q** with
$$J(\mathbf{y}, \mathbf{q}) = .8$$
:
 $\begin{pmatrix} - ((-, \mathcal{G})^{\dagger} \gg .9 \mathcal{G}) \end{pmatrix}$

Chance we need to check z with J(y, z) = .4:

$$(-(1-y^{-})^{+}=.12$$

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

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



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

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



r = 40, t = 5

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

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 <u>*t* = 40</u>,

- 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 $\lesssim 1-(1-.7^{25})^{40}=.\underbrace{.005}$

Upper bound on total number of items checked:

 $(.95 \cdot 10) + (.95 \cdot 10, 000) + (.005 \cdot 9, 989, 990) \approx 60,000 \ll 10,000,000.$ 44

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



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

Theorem (Indyk, Motwani, 1998)

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

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

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.

Euclidean distance
$$\|\mathbf{x} - \mathbf{y}\|_{2}^{2}$$
: = $(x - y)^{\dagger} (x - y) = x^{\intercal} - 2x^{\intercal} y$
• Suppose for simplicity that $\|\mathbf{x}\|_{2}^{2} = \|\mathbf{y}\|_{2}^{2} = 1$.
• $\|\mathbf{x}\|_{2}^{2} - 2x^{\intercal} y$
+ $\|\mathbf{y}\|_{2}^{2} = 1$.

SIMHASH

(Locality sensitive hash for cosine similarity:

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

•
$$h: \mathbb{R}^d \to \{1, \ldots, m\}$$
 is definied $h(\mathbf{x}) = f(\operatorname{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)).$

If $cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $Pr[h(\mathbf{x}) == h(\mathbf{y})]$? = 1

SIMHASH ANALYSIS



SimHash can be tuned, 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
 $h(\mathbf{x}) = f([\operatorname{sign}(\langle \mathbf{g}_1, \mathbf{x} \rangle), \dots, \operatorname{sign}(\langle \mathbf{g}_r, \mathbf{x} \rangle)]).$

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

SIMHASH ANALYSIS

To prove:

$$\Pr[h(\mathbf{x}) == h(\mathbf{y})] = 1 - \frac{\theta}{\pi}, \text{ where } h(\mathbf{x}) = f(\operatorname{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)).$$

$$\begin{array}{c} \mathbf{x} \\ \mathbf{y} \\ \mathbf{y}$$

SIMHASH ANALYSIS



 $Pr[h(\mathbf{x}) == h(\mathbf{y})] \approx$ probability \mathbf{x} and \mathbf{y} are on the same side of hyperplane orthogonal to \mathbf{g} .