CS-GY 6763: Lecture 4 Dimensionality reduction + Near neighbor search in high dimensions

NYU Tandon School of Engineering, Prof. Christopher Musco

EUCLIDEAN DIMENSIONALITY REDUCTION



Please remember: This is equivalent to:

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 $k = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all $\underline{i, j}$,

$$(1-\epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \le \|\mathbf{\Pi}\mathbf{q}_i - \mathbf{\Pi}\mathbf{q}_j\|_2^2 \le (1+\epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2.$$

because for small ϵ , $(1 + \epsilon)^2 = 1 + O(\epsilon)$ and $(1 - \epsilon)^2 = 1 - O(\epsilon)$. $\leq 1 + 3\epsilon$

And this is equivalent to:

Lemma (Johnson-Lindenstrauss, 1984)

For any set of n data points $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^d$ there exists a <u>linear map</u> $\Pi : \mathbb{R}^d \to \mathbb{R}^k$ where $k = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all <u>i</u>, j, $(1-\epsilon) \|\Pi \mathbf{q}_i - \Pi \mathbf{q}_j\|_2^2 \le \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \le (1+\epsilon) \|\Pi \mathbf{q}_i - \Pi \mathbf{q}_j\|_2^2$.

because for small ϵ , $\frac{1}{1+\epsilon} = 1 - O(\epsilon)$ and $\frac{1}{1-\epsilon} = 1 + O(\epsilon)$.

k-means clustering: Give data points $\mathbf{a}_1, \ldots, \mathbf{a}_n \in \mathbb{R}^d$, find centers $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k \in \mathbb{R}^d$ to minimize:

$$Cost(\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1,\ldots,k} \left\| \boldsymbol{\mu}_j - \boldsymbol{a}_i \right\|_2^2$$



k-means clustering: Give data points $\mathbf{a}_1, \ldots, \mathbf{a}_n \in \mathbb{R}^d$, find centers $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k \in \mathbb{R}^d$ to minimize:

$$Cost(\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_k) = \sum_{j=1}^n \min_{j=1,\ldots,k} \|\boldsymbol{\mu}_j - \mathbf{a}_j\|_2^2$$



k-means clustering: Give data points $\mathbf{a}_1, \ldots, \mathbf{a}_n \in \mathbb{R}^d$, find centers $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k \in \mathbb{R}^d$ to minimize:

$$Cost(\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1,\ldots,k} (|\boldsymbol{\mu}_j - \underline{\mathbf{a}}_i||_2^2)$$



NP hard to solve exactly, but there are many good approximation algorithms. All depend at least linearly on the dimension *d*.

Approximation scheme: Find clusters $\tilde{C}_1, \ldots, \tilde{C}_k$ for the $k = O\left(\frac{\log n}{\epsilon^2}\right)$ dimension data set $\underline{\Pi} a_1, \ldots, \underline{\Pi} a_n$.



Argue these clusters are near optimal for $\mathbf{a}_1, \ldots, \mathbf{a}_n$.

K-MEANS CLUSTERING



Exercise: Prove this to your self.

K-MEANS CLUSTERING

$$\underline{Cost}(C_{1},...,C_{k}) = \sum_{j=1}^{k} \frac{1}{2|C_{j}|} \sum_{u,v \in C_{j}} \frac{\|\mathbf{a}_{u} - \mathbf{a}_{v}\|_{2}^{2}}{\sum_{i=1}^{k} \frac{1}{2|C_{j}|} \sum_{u,v \in C_{j}} \frac{\|\Pi \mathbf{a}_{u} - \Pi \mathbf{a}_{v}\|_{2}^{2}}$$

Claim: For any clusters C_1, \ldots, C_k :

$$(1-\epsilon)Cost(C_1,\ldots,C_k) \leq \widetilde{Cost}(C_1,\ldots,C_k) \leq (1+\epsilon)Cost(C_1,\ldots,C_k)$$

Suppose we use an approximation algorithm to find clusters $B_1, \ldots, \underline{B_k}$ such that: α)Cosť

Cost

Then:

$$\underbrace{Cost(B_1, \dots, B_k)}_{\leq (1 + O(\epsilon))(1 + \alpha)\widetilde{Cost}^*} \leq (1 + O(\epsilon))(1 + \alpha)\widetilde{Cost}^*}_{\leq (1 + O(\epsilon))(1 + \alpha)(1 + \epsilon)Cost^*}$$
$$= (1 + O(\alpha + \epsilon))Cost^*$$
$$= (1 + O(\alpha + \epsilon))Cost^*$$
$$= (1 + O(\alpha + \epsilon))Cost^*$$
$$\underbrace{Cost^* = \min_{C_1, \dots, C_k}Cost(C_1, \dots, C_k)and}_{Cost^* = \min_{C_1, \dots, C_k}Cost(C_1, \dots, C_k)and}$$

EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Johnson-Lindenstrauss, 1984)

For any set of n data points $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^d$ there exists a linear map $\prod : \mathbb{R}^d \to \mathbb{R}^k$ where $k = O\left(\frac{\log n^{\lambda}}{\epsilon^2}\right)$ such that for all $\underline{i}, \underline{j}, \qquad 1 - \delta$

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



Remarkably, Π can be chosen <u>completely at random</u>!

One possible construction: Random Gaussian.

$$\int \mathbf{\Pi}_{i,j} = \frac{1}{\sqrt{k}} \mathcal{N}(0,1)$$

The map **Π** is **oblivious to the data set**. This stands in contrast to e.g. PCA, amoung other differences.

[Indyk, Motwani 1998] [Arriage, Vempala 1999] [Achlioptas 2001] [Dasgupta, Gupta 2003].

Many other possible choices suffice – you can use random $\{+1, -1\}$ variables, sparse random matrices, pseudorandom Π . Each with different advantages.

Let $\Pi \in \mathbb{R}^{k \times d}$ be chosen so that each entry equals $\frac{1}{\sqrt{h}}\mathcal{N}(0,1)$ or each entry equals $\frac{1}{\sqrt{b}} \pm 1$ with equal probability. TIT'S T -0.3538 -0.2938 2,9888 -1.3617-0.2620 1.3546 -1.7502 -0.0348 -1.1201 -1.4491-0.8487 -0.7982 h 0.1240 -0.1332 0.3502 -2.8026 -0.2248 -0.2097 -2.8518 -0.2991 >> Pi = randn(m,d); >> s = (1/sqrt(m))*Pi*q; **)** Pi = 2*randi(2,m,d)-3; s = (1/sqrt(m))*Pi*q:

A random orthogonal matrix also works. I.e. with $O^T = I_{k \times k}$. For this reason, the JL operation is often called a "random projection", even though it technically isn't a <u>projection</u> when entries are i.i.d.

RANDOM PROJECTION



Intuitively, close points will remain close after projection, and far points will remain far.

EUCLIDEAN DIMENSIONALITY REDUCTION

x = 8; - 8; B1 1. 2-Intermediate result: Lemma (Distributional JL Lemma) Let $\mathbf{\Pi} \in \mathbb{R}^{k \times d}$ be chosen so that each entry equals $(\frac{1}{\sqrt{k}}\mathcal{N}(0,1))$ where $\mathcal{N}(0,1)$ denotes a standard Gaussian random variable. If we choose $\underline{k} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any vector **x**, with probability $(1 - \delta)$: $(1-\epsilon)\|\mathbf{x}\|_{2}^{2} \leq \|\mathbf{\Pi}\mathbf{x}\|_{2}^{2} \leq (1+\epsilon)\|\mathbf{x}\|_{2}^{2}$ $(1-\varepsilon) \| q_i - q_j \|_{2^2} \leq \| \pi(q_i - q_j) \|_{2^2} \leq (1+\varepsilon) \| q_i - q_j \|_{2^2}$ ≤ 1TK - T 9; 12 € Given this lemma, how do we prove the traditional Johnson-Lindenstrauss lemma? $\left| o_{j} \left(\left| c n^{2} \right) \right| = O\left(\left| o_{j} \left(r \right) \right\rangle \right)$ SE IONZ 9/10

We have a set of vectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Fix $i, j \in 1, \dots, n$. Let $\mathbf{x} = \mathbf{q}_i - \mathbf{q}_j$. By linearity, $\mathbf{\Pi} \mathbf{x} = \mathbf{\Pi} (\mathbf{q}_i - \mathbf{q}_j) = \mathbf{\Pi} \mathbf{q}_i - \mathbf{\Pi} \mathbf{q}_j$. By the Distributional JL Lemma, with probability $1 - \delta$, $(1 - \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2 \le \|\mathbf{\Pi} \mathbf{q}_i - \mathbf{\Pi} \mathbf{q}_j\|_2 \le (1 + \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2$. By the Distributional JL Lemma, with probability $1 - \delta$, we have that for any \mathbf{q}_i , \mathbf{q}_j ,

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

Finally, set $\delta = \frac{1}{10n^2}$. Since there are $< n^2$ total *i*, *j* pairs, by a union bound we have that with probability 9/10, the above will hold for all *i*, *j*, as long as we compress to:

$$k = O\left(\frac{\log(1/(1/n^2))}{\epsilon^2}\right) = O\left(\frac{\log n}{\epsilon^2}\right) \text{ dimensions.} \quad \Box$$

Want to argue that, with probability $(1 - \delta)$, $(1 - \epsilon) ||\mathbf{x}||_2^2 \le |\mathbf{\Pi}\mathbf{x}||_2^2 \le (1 + \epsilon) ||\mathbf{x}||_2^2$ Claim: $\mathbb{E} [||\mathbf{\Pi}\mathbf{x}||_2^2] = ||\mathbf{x}||_2^2$. Some notation:



19



Goal: Prove
$$\mathbb{E} [\|\Pi x\|_2^2] = \|x\|_2^2$$
.

where each Z_1, \ldots, Z_d is a standard normal $\mathcal{N}(0, 1)$ random variable.

We have that $Z_i(i)$ is a normal $\mathcal{N}(0, \mathbf{x}(i)^2)$ random variable.

Goal: Prove $\mathbb{E}\left[\|\mathbf{\Pi}\mathbf{x}\|_2^2\right] = \|\mathbf{x}\|_2^2$. Have: $\mathbb{E}\left[\|\mathbf{\Pi}\mathbf{x}\|_2^2\right] = \mathbb{E}\left[\left(\langle \boldsymbol{\pi}_1, \mathbf{x} \rangle\right)^2\right]$

What type of random variable is $\langle \pi_i, x \rangle$? Fact (Stability of Gaussian random variables) $\mathcal{N}(\mu_1, \sigma_1^2) + \mathcal{N}(\mu_2, \sigma_2^2) = \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$ $(\pi_{1}, \mathbf{x}) = \mathcal{N}(0, \mathbf{x}(1)^{2}) + \mathcal{N}(0, \mathbf{x}(2)^{2}) + \dots + \mathcal{N}(0, \mathbf{x}(d)^{2})$ $= \mathcal{N}(0, \|\mathbf{x}\|_{2}^{2}).$ $\chi(\mathbf{u})^{2} + \dots + \chi(\mathbf{d})^{2} = \|\mathbf{x}\|_{2}^{2}$ So $\mathbb{E}\left[\|\Pi \mathbf{x}\|_{2}^{2}\right] = \mathbb{E}\left[\left(\langle \pi_{1}, \mathbf{x} \rangle\right)^{2}\right] = \|\mathbf{x}\|_{2}^{2}$, as desired.

Goal: Prove $\mathbb{E}\left[\|\mathbf{\Pi}\mathbf{x}\|_2^2\right] = \|\mathbf{x}\|_2^2$. Have: $\mathbb{E}\left[\|\mathbf{\Pi}\mathbf{x}\|_2^2\right] = \mathbb{E}\left[\left(\langle \pi_1, \mathbf{x} \rangle\right)^2\right]$

Want to argue that, with probability
$$(1 - \delta)$$
,
 $(1 - \epsilon) ||\mathbf{x}||_2^2 \le ||\mathbf{\Pi}\mathbf{x}||_2^2 \le (1 + \epsilon) ||\mathbf{x}||_2^2$
1. $\mathbb{E}[||\mathbf{\Pi}\mathbf{x}||_2^2] = ||\mathbf{x}||_2^2$
2. Need to use a concentration bound.
 $||\mathbf{\Pi}\mathbf{x}||_2^2 = \frac{1}{k} \sum_{i=1}^k (\langle \pi_i, \mathbf{x} \rangle)^2 = \frac{1}{k} \sum_{i=1}^k \mathcal{N}(0, ||\mathbf{x}||_2^2)^2$
"Chi-squared random variable with k degrees of freedom."

 \frown

Lemma

Let Z be a Chi-squared random variable with k degrees of freedom.

$$\Pr[|\mathbb{E}Z - Z| \ge \epsilon \mathbb{E}Z] \le 2e^{-k\epsilon^2/8} \le \delta$$

$$\Pr(\left||\mathbb{T}\chi||_{*}^{*} - ||\chi||_{*}^{*}| \neg \varepsilon ||\chi||_{*}^{*}| \neg \varepsilon ||\chi||_{*}^{2}| \le 2e^{-k\epsilon^2/8} \le \delta$$

$$e^{-k\epsilon^{1/8}} = \delta/2 - k\epsilon^{1/8} = \log(\delta/2)$$

$$k\epsilon^{1/8} = \log(\delta/2) - k\epsilon^{1/8} = \log(\delta/2)$$

$$= \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)})$$

$$= \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)})$$

$$= \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)})$$

$$= \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/8)})$$

$$= \delta(1^{6}\delta^{(1/8)}) = \delta(1^{6}\delta^{(1/$$

If high dimensional geometry is so different from low-dimensional geometry, why is <u>dimensionality reduction</u> <u>possible?</u> Doesn't Johnson-Lindenstrauss tell us that high-dimensional geometry can be approximated in low dimensions?

CONNECTION TO DIMENSIONALITY REDUCTION

Hard case
$$x_1, \dots, x_n \in \mathbb{R}^d$$
 are all mutually orthogonal unit
vectors:
$$\|x_i - x_j\|_2^2 = 2$$
From our result earlier, in $O(\log n/\epsilon^2)$ dimensions, there exists $2^{O(\epsilon^2 \cdot \log n/\epsilon^2)} \ge n$ unit vectors that are close to mutually orthogonal

orthogonal.

 $O(\log n/\epsilon^2)$ = just enough dimensions.

The Johnson-Lindenstrauss Lemma let us sketch vectors and preserve their ℓ_2 Euclidean distance.

We also have dimensionality reduction techniques that preserve alternative measures of similarity!

JACCARD SIMILARITY

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).

$$\mathbf{y} = \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{y}} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{y}} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{y}} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{y}} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{y}} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \hat$$

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



Given **q**, find any nearby "fingerprint" **y** in a database – i.e. any **y** with dist(**y**, **q**) small.

Challenges:

> Kad

- Database is possibly huge: O(nd) bits.
- Expensive to compute dist(y, q): O(d) time.

Olyk)

"Bag-of-words" model:



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



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.

JACCARD SIMILARITY FOR SEISMIC DATA



Feature extract pipeline for earthquake data.

(see paper by Rong et al. posted on course website)

Goal: Design a compact sketch $C : \{0, 1\} \rightarrow \mathbb{R}^k$:



Homomorphic Compression: Want to use C(q), C(y) to approximately compute the Jaccard similarity J(q, y).

MINHASH

MinHash (Broder, '97):


MINHASH

• Choose *k* random hash functions

$$h_1, \ldots, h_k : \{1, \ldots, n\} \to [0, 1].$$

• For
$$i \in 1, \ldots, k$$
,
• Let $c_i = \min_{i \in [n] = 1} h_i($

$$C(\mathbf{q}) = [c_1, \ldots, c_k].$$



MINHASH ANALYSIS



Proof:

1. For $c_i(\mathbf{q}) = c_i(\mathbf{y})$, we need that $\arg\min_{i \in \mathbf{q}} h(i) = \arg\min_{i \in \mathbf{y}} h(i)$.

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

Unbiased estimate for Jaccard similarity:
 $\mathbb{E}\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} (\mathbf{F} \mathbf{1}(\ldots)) - \frac{1}{k} \sum_{i=1}^{k} \mathbf{j}(\mathbf{q}, \mathbf{y}) =$

The more repetitions, the lower the variance.

Let $J = J(\mathbf{q}, \mathbf{y})$ denote the true Jaccard similarity. $\mathcal{L} \circ \rho \mathcal{J}$ Estimator: $\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} \underline{\mathbb{1}}[c_i(\mathbf{q}) = c_i(\mathbf{y})].$ $\operatorname{Var}[\tilde{J}] = \frac{l}{kr} \sum_{i=1}^{k} \operatorname{Vor}\left(\operatorname{I}\left(\operatorname{c}_{i}(q) = \operatorname{c}_{i}(\gamma)\right)\right) = \frac{l}{kr} \sum_{i=1}^{k} \overline{J} - \overline{J}^{r}$ so that with probability > $1 - \delta$, $|J - \tilde{J}| \leq \epsilon$? 1 N= 1 $Pr[[\tilde{J}-J] > \varepsilon] \in \delta$ Pr [] f - T [Z 46] ≤ 1 02 $\Pr\left(\left|\hat{J}-J\right| \geq \frac{1}{6} \cdot \frac{1}{16}\right] \leq d$ 41 **Chebyshev inequality:** As long as $k = O\left(\frac{1}{\epsilon^2 \delta}\right)$, then with prob. $1 - \delta$,

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

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?

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(nd) to O(nk).
- 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 $d = o(\log n)$.

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]
 - Probably the most practical. This is most similar to the custom method e.g. Shazam uses.

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

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 $g : [0,1] \to \{1, \dots, m\}$ 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 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 $J(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to \mathbf{y} .

Note: The meaning of "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.

Suppose the nearest database point q has J(y, 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?

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*(**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?

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 random MinHash.
- + Let $g: [0,1]^r \to \{1,\ldots,m\}$ be a uniform random hash function.

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



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:



61

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

S-CURVE TUNING

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

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



r = 5, t = 40

S-CURVE TUNING

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

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



65

S-CURVE TUNING

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 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:

 $.95 \cdot 10 + .95 \cdot 10,000 + .005 \cdot 9,989,990 \approx 60,000 \ll 10,000,000.$ 67

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

Near Neighbor Search Problem

Concrete worst case result:

Theorem (Indyk, Motwani, 1998)

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: $O(n^{1+1/C})$.

 $\|\boldsymbol{q}-\boldsymbol{y}\|_0=$ "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 $\|\tilde{q}-y\|_0 \leq C \cdot \|q-y\|_0$ in:

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

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 $\|\mathbf{x} - \mathbf{y}\|_2^2$:

• Suppose for simplicity that $\|\mathbf{x}\|_2^2 = \|\mathbf{y}\|_2^2 = 1$.

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 $Pr[h(\mathbf{x}) == h(\mathbf{y})]$?

SIMHASH ANALYSIS

Theorem: If $cos(\theta(\mathbf{x}, \mathbf{y})) = v$, then $Pr[h(\mathbf{x}) == h(\mathbf{y})] = 1 - \frac{\theta}{\pi} + O(\frac{1}{m}) = 1 - \frac{cos^{-1}(v)}{\pi} + O(\frac{1}{m})$



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:

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



76

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} .