

CS-GY 9223 D: Lecture 3 Supplemental

The Johnson-Lindenstrauss Lemma

NYU Tandon School of Engineering, Prof. Christopher Musco

Abstract architecture of a sketching algorithm:

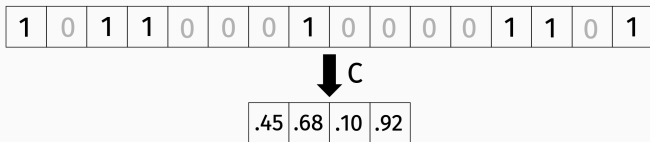
- Given a dataset $D = d_1, \dots, d_n$ with n pieces of data, we want to output $f(D)$ for some function f .
- **Sketch phase:** For each $i \in 1, \dots, n$, compute $s_i = C(d_i)$, where C is some compression function and $|s_i| \ll d_i$.
- **Process phase:** Using (lower dimensional) dataset s_1, \dots, s_n , compute an approximation to $f(D)$.



Better space complexity,
communication complexity,
runtime, all at once.

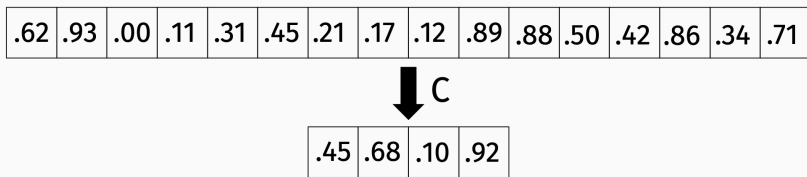
BINARY VECTOR COMPRESSION

We already saw a powerful application of sketching (the MinHash algorithm) to compressing binary vectors.



Let us estimate the Jaccard similarity between any two binary vectors \mathbf{q} and \mathbf{y} using the information in $C(\mathbf{q})$ and $C(\mathbf{y})$ alone.

TODAY: EUCLIDEAN DIMENSIONALITY REDUCTION



Euclidean norm / distance:

- Given $\mathbf{q} \in \mathbb{R}^d$, $\|\mathbf{q}\|_2 = \sqrt{\sum_{i=1}^d q(i)^2}$.
- Given $\mathbf{q}, \mathbf{y} \in \mathbb{R}^d$, distance defined as $\|\mathbf{q} - \mathbf{y}\|_2$.

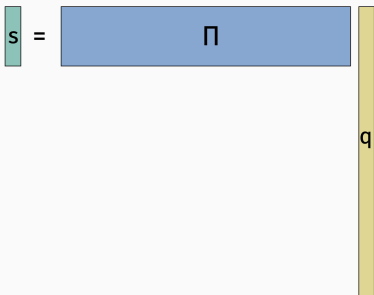
Can we find compact sketches that preserve Euclidean distance, just as we did for Jaccard similarity?

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 $\Pi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ where $k = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all i, j ,

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



Please remember: 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 linear map $\Pi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ where $k = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all i, j ,

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

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

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 linear map $\Pi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ where $k = O\left(\frac{\log n}{\epsilon^2}\right)$ such that for all i, j ,

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

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

Remarkably, Π can be chosen completely at random!

One possible construction: Random Gaussian.

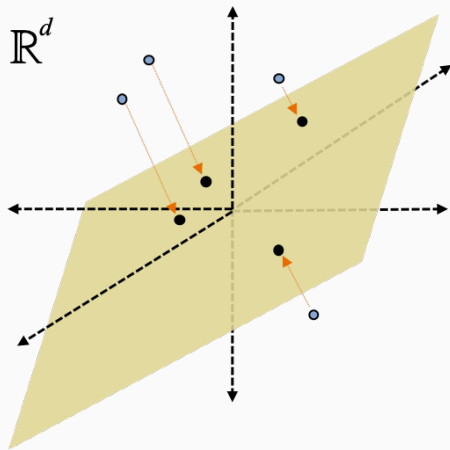
$$\Pi_{i,j} = \frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$$

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

[Indyk, Motwani 1998] [Arriaga, 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.

RANDOM PROJECTION



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

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 $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, then for any vector \mathbf{x} , with probability $(1 - \delta)$:

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

Given this lemma, how do we prove the traditional Johnson-Lindenstrauss lemma?

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 $1 - \delta$,

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

Finally, set $\delta = \frac{1}{n^2}$. 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:

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

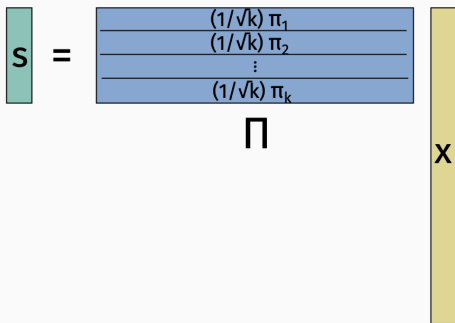
PROOF OF DISTRIBUTIONAL JL

Want to argue that, with probability $(1 - \delta)$,

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

Claim: $\mathbb{E}\|\Pi\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$.

Some notation:



So each π_i contains $\mathcal{N}(0, 1)$ entries.

$$\|\Pi \mathbf{x}\|_2^2 = \sum_i^k s(i)^2 = \sum_i^k \left(\frac{1}{\sqrt{k}} \langle \boldsymbol{\pi}_i, \mathbf{x} \rangle \right)^2 = \frac{1}{k} \sum_i^k (\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle)^2$$

$$\begin{aligned} \mathbb{E} [\|\Pi \mathbf{x}\|_2^2] &= \frac{1}{k} \sum_i^k \mathbb{E} [(\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle)^2] \\ &= \mathbb{E} [(\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle)^2] \end{aligned}$$

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

$$\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle = Z_1 \cdot \mathbf{x}(1) + Z_2 \cdot \mathbf{x}(2) + \dots + Z_d \cdot \mathbf{x}(d)$$

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

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

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

What type of random variable is $\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle$?

Fact (Stability of Gaussian random variables)

$$\mathcal{N}(\mu_1, \sigma_1^2) + \mathcal{N}(\mu_2, \sigma_2^2) = \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$\begin{aligned}\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle &= \mathcal{N}(0, \mathbf{x}(1)^2) + \mathcal{N}(0, \mathbf{x}(2)^2) + \dots + \mathcal{N}(0, \mathbf{x}(d)^2) \\ &= \mathcal{N}(0, \|\mathbf{x}\|_2^2).\end{aligned}$$

So $\mathbb{E}\|\boldsymbol{\Pi}\mathbf{x}\|_2^2 = \mathbb{E}\left[\left(\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle\right)^2\right] = \|\mathbf{x}\|_2^2$, as desired.

Want to argue that, with probability $(1 - \delta)$,

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

1. $\mathbb{E}\|\Pi\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$.
2. Need to use a concentration bound.

$$\|\Pi\mathbf{x}\|_2^2 = \frac{1}{k} \sum_{i=1}^k (\langle \pi_i, \mathbf{x} \rangle)^2 = \frac{1}{k} \sum_{i=1}^k \mathcal{N}(0, \|\mathbf{x}\|_2^2)$$

“Chi-squared random variable with k degrees of freedom.”

Lemma

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

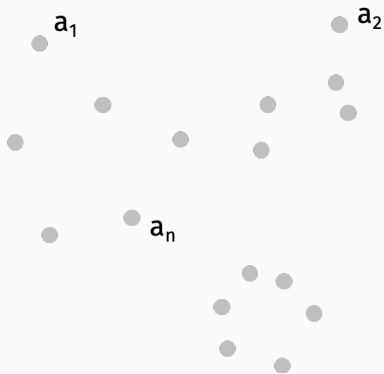
$$\Pr[|\mathbb{E}Z - Z| \geq \epsilon \mathbb{E}Z] \leq 2e^{-k\epsilon^2/8}$$

Goal: Prove $\|\mathbf{nx}\|_2^2$ concentrates within $1 \pm \epsilon$ of its expectation, which equals $\|\mathbf{x}\|_2^2$.

SAMPLE APPLICATION

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

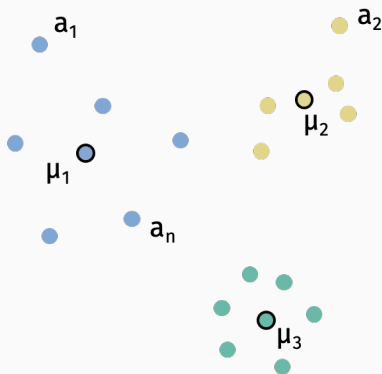
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{x}_i\|_2^2$$



SAMPLE APPLICATION

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

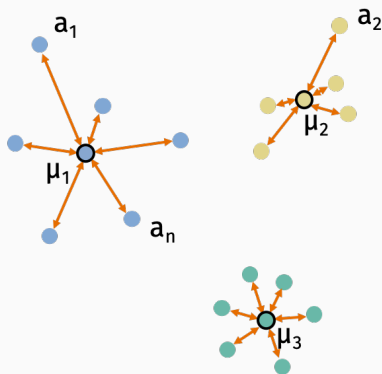
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{x}_i\|_2^2$$



SAMPLE APPLICATION

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

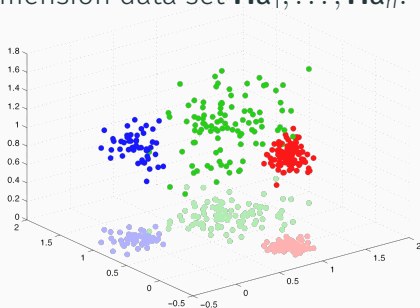
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



K-MEANS CLUSTERING

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

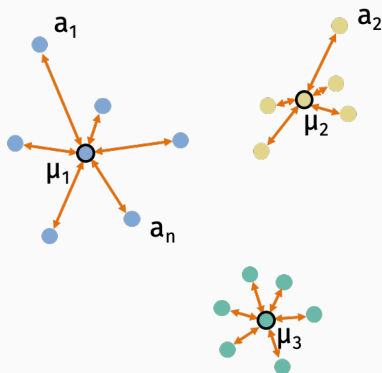
Approximation scheme: Find clusters $\tilde{C}_1, \dots, \tilde{C}_k$ for the $k = O\left(\frac{\log n}{\epsilon^2}\right)$ dimension data set $\mathbf{a}_1, \dots, \mathbf{a}_n$.



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

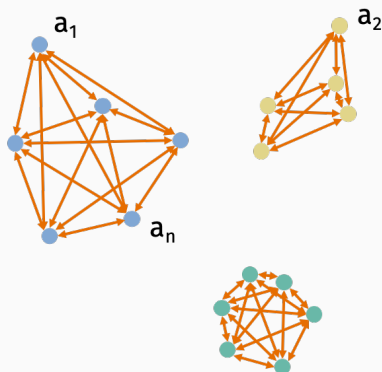
Equivalent formulation: Find clusters $C_1, \dots, C_k \subseteq \{1, \dots, n\}$ to minimize:

$$\text{Cost}(C_1, \dots, C_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u, v \in C_j} \|a_u - a_v\|_2^2.$$



Equivalent formulation: Find clusters $C_1, \dots, C_k \subseteq \{1, \dots, n\}$ to minimize:

$$\text{Cost}(C_1, \dots, C_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u, v \in C_j} \|a_u - a_v\|_2^2.$$



$$\text{Cost}(C_1, \dots, C_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\mathbf{a}_u - \mathbf{a}_v\|_2^2$$

$$\widetilde{\text{Cost}}(C_1, \dots, C_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\Pi \mathbf{a}_u - \Pi \mathbf{a}_v\|_2^2$$

Let $Cost^* = \min Cost(C_1, \dots, C_k)$ and
 $\widetilde{Cost}^* = \min \widetilde{Cost}(C_1, \dots, C_k)$.

Claim: $(1 - \epsilon)Cost^* \leq \widetilde{Cost}^* \leq (1 + \epsilon)Cost^*$.

Suppose we use an approximation algorithm to find clusters B_1, \dots, B_k such that:

$$\widetilde{\text{Cost}}(B_1, \dots, B_k) \leq (1 + \alpha) \widetilde{\text{Cost}}^*$$

Then:

$$\begin{aligned} \text{Cost}(B_1, \dots, B_k) &\leq \frac{1}{1 - \epsilon} \widetilde{\text{Cost}}(B_1, \dots, B_k) \\ &\leq (1 + \alpha)(1 + O(\epsilon)) \widetilde{\text{Cost}}^* \\ &\leq (1 + \alpha)(1 + O(\epsilon))(1 + \epsilon) \text{Cost}^* \\ &= 1 + O(\alpha + \epsilon) \text{Cost}^* \end{aligned}$$

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

Hard case: $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ are all mutually orthogonal unit vectors:

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = 2 \quad \text{for all } i, j.$$

From our result earlier, in $O(\log n/\epsilon^2)$ dimensions, there exists $2^{O(\epsilon^2 \cdot \log n/\epsilon^2)} \geq n$ unit vectors that are close to mutually orthogonal.

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