

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:

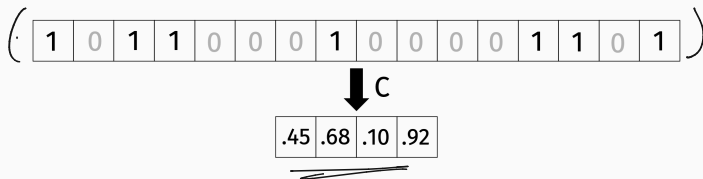
- Given a dataset $\underline{D} = \underline{d}_1, \dots, \underline{d}_n$ with n pieces of data, we want to output $f(D)$ for some function f .
- **Sketch phase:** For each $i \in 1, \dots, n$, compute $s_i = C(d_i)$, where C is some compression function and $|s_i| \ll d_i$.
- **Process phase:** Using (lower dimensional) dataset s_1, \dots, s_n , compute an approximation to $f(D)$.



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

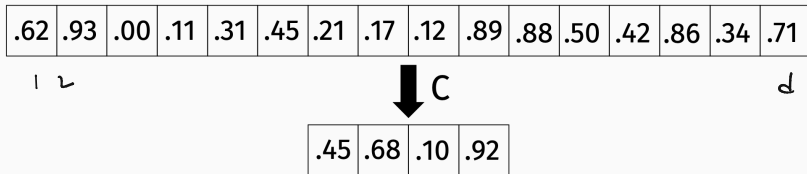
BINARY VECTOR COMPRESSION

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 $\mathbf{q} \in \mathbb{R}^d$, $\|\mathbf{q}\|_2 = \sqrt{\sum_{i=1}^d q(i)^2}$. ↗ *ith entry of \mathbf{q}*
- 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)

$$k \ll d$$

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

$$\left((1 - \epsilon) \|\underline{q}_i - \underline{q}_j\|_2 \right) \leq \left(\|\underline{\Pi} \underline{q}_i - \underline{\Pi} \underline{q}_j\|_2 \right) \leq \left((1 + \epsilon) \|\underline{q}_i - \underline{q}_j\|_2 \right)$$

$$\begin{array}{c}
 \begin{array}{c} k \\ \downarrow \\ k = O\left(\frac{\log n}{\epsilon^2}\right) \end{array}
 \left\{ \begin{array}{c} \text{green box} \\ \underline{s} \end{array} \right\} = \begin{array}{c} \begin{array}{cc} & d \\ \text{blue box} & \Pi \end{array} \\ \left. \begin{array}{c} \text{yellow box} \\ \underline{q} \end{array} \right\} d
 \end{array}$$

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 ,

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

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

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.

RANDOMIZED JL CONSTRUCTIONS

Let $\Pi \in \mathbb{R}^{k \times d}$ be chosen so that each entry equals $\frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$.

... or each entry equals $\frac{1}{\sqrt{k}} \pm 1$ with equal probability.

Π

k

-2.1384	2.9888	-0.3538	0.0229	0.5201	-0.2938	-1.3328	-1.3617	-0.1952
-0.8396	0.8252	-0.8236	-0.2620	-0.0200	-0.8479	-2.3299	0.4550	-0.2176
1.3546	1.3798	-1.5771	-1.7582	-0.0348	-1.1201	-1.4491	-0.8487	-0.3831
-1.0722	-1.0582	0.5088	-0.2857	-0.7982	2.5268	0.3335	-0.3349	0.0238
0.9610	-0.4686	0.2820	-0.8314	1.0187	1.6555	0.3914	0.5528	0.0513
0.1240	-0.2725	0.0335	-0.9792	-0.1332	0.3075	0.4517	1.0391	0.8261
1.4367	1.0984	-1.3337	-1.1564	-0.7145	-1.2571	-0.1383	-1.1176	1.5278
-1.9689	-0.2779	1.1275	-0.5336	1.3514	-0.8655	0.1837	1.2607	0.4669
-0.1977	0.7015	0.3582	-2.0026	-0.2248	-0.1765	-0.4762	0.6501	-0.2097
-1.2078	-2.0518	-0.2991	0.9642	-0.5898	0.7914	0.8628	-0.0679	0.6252

```
>> Pi = randn(m,d);  
>> s = (1/sqrt(m))*Pi*q;
```

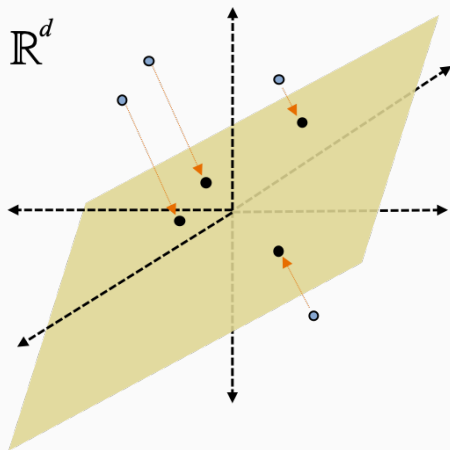
1	1	-1	-1	-1	-1	-1	-1	1	-1	-1	1	-1	-1	1	1	-1
1	1	1	-1	1	-1	-1	-1	1	1	1	-1	1	-1	-1	-1	-1
1	1	-1	-1	-1	1	-1	-1	1	1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	1	1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	1	1
1	-1	1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	-1	1
1	-1	-1	1	-1	1	1	-1	-1	-1	1	-1	-1	-1	1	1	1
1	1	-1	1	1	-1	1	-1	1	-1	1	-1	1	1	1	1	-1
-1	-1	-1	-1	-1	-1	1	1	1	-1	1	-1	1	-1	-1	1	1
-1	-1	1	1	1	1	-1	-1	1	-1	1	1	1	-1	1	-1	1
-1	1	-1	1	-1	1	1	-1	-1	1	-1	1	-1	-1	1	-1	1

```
>> Pi = 2*randi(2,m,d)-3;  
>> s = (1/sqrt(m))*Pi*q;
```

A random orthogonal matrix also works. I.e. with $\Pi\Pi^T = I_{k \times k}$.

For this reason, the JL operation is often called a “random projection”, even though it technically isn’t a projection when entries are i.i.d.

RANDOM PROJECTION



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

Intermediate result:

Lemma (Distributional JL Lemma)

Let $\Pi \in \mathbb{R}^{k \times d}$ 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 $(1 - \delta)$:

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

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

JL FROM DISTRIBUTIONAL JL

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

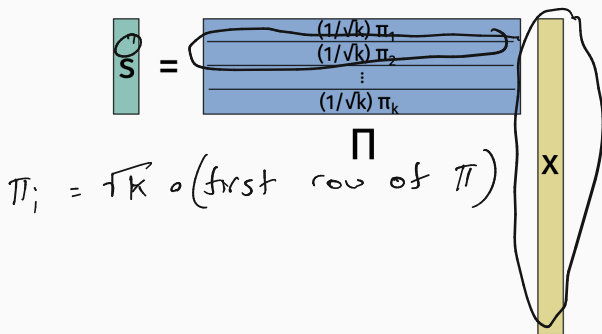
Want to argue that, with probability $(1 - \delta)$,

$$n = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$$

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

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

Some notation:



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

PROOF OF DISTRIBUTIONAL JL

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

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

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

PROOF OF DISTRIBUTIONAL JL

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

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

STABLE RANDOM VARIABLES

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

Fact (Stability of Gaussian random variables)

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

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

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

\downarrow
 $= \text{var} [\langle \pi_i, \mathbf{x} \rangle]$

PROOF OF DISTRIBUTIONAL JL

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

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

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."

CONCENTRATION OF CHI-SQUARED RANDOM VARIABLES

Lemma

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

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

$$Z = \|\Pi x\|_2^2$$

$$2e^{-k\epsilon^2/8} = \delta \quad \begin{aligned} -k\epsilon^2 &= O(\log \delta) \\ k\epsilon^2 &= O(\log 1/\delta) \end{aligned}$$

$$\mathbb{E}[Z] = \|x\|_2^2$$

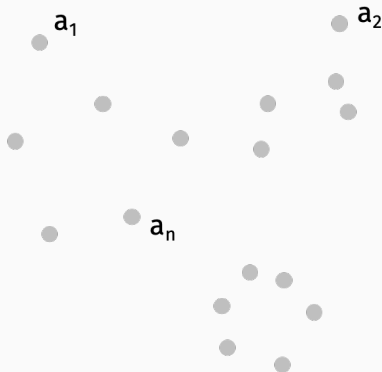
$$\Pr \left[\underbrace{|\|x\|_2^2 - \|\Pi x\|_2^2|}_{\leq \delta} \geq \epsilon \|x\|_2^2 \right] \leq 2e^{-k\epsilon^2/8} \quad k = O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$$

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

SAMPLE APPLICATION

k-means clustering: Give data points $\underline{a_1}, \dots, \underline{a_n} \in \mathbb{R}^d$, find centers $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ to minimize:

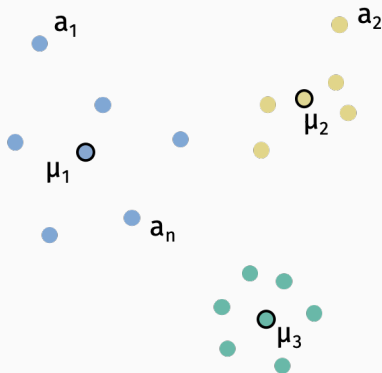
$$\text{Cost}(\underline{\mu_1}, \dots, \underline{\mu_k}) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\underline{\mu_j} - \underline{a_i}\|_2^2$$



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:

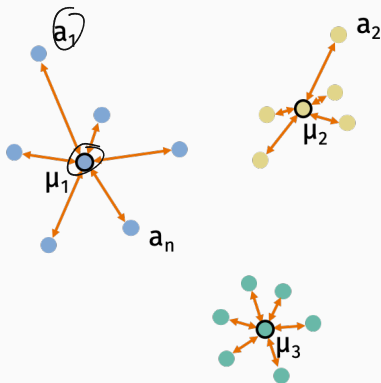
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



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:

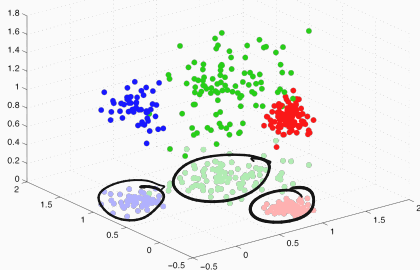
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



K-MEANS CLUSTERING

NP hard to solve exactly, but there are many good approximation algorithms. All depend at least linearly on the dimension d . $a_1, \dots, a_n \in \mathbb{R}^d$

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 $\underline{a}_1, \dots, \underline{a}_n \in \mathbb{R}^k$



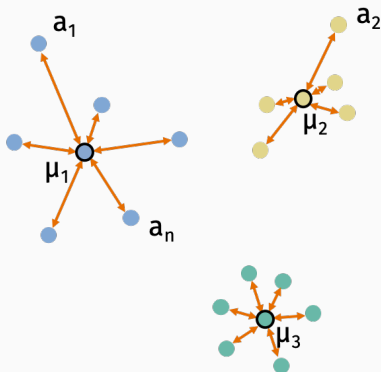
Argue these clusters are near optimal for a_1, \dots, a_n .

K-MEANS CLUSTERING

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

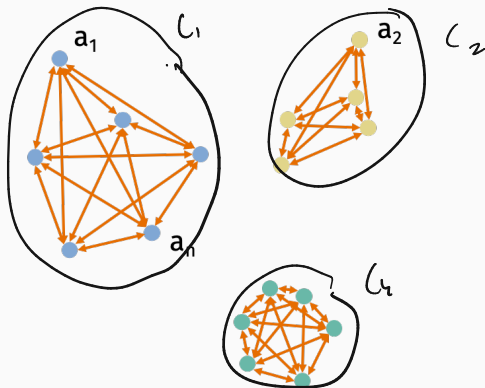
$$Cost(\underline{C}_1, \dots, \underline{C}_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|a_u - a_v\|_2^2.$$

$C_1 \cup C_2 \cup \dots \cup C_k = \{1, \dots, n\}$



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

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



K-MEANS CLUSTERING

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

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

For any clusters C_1, \dots, C_k

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

$$= \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|\Pi a_u - \Pi a_v\|_2^2 \leq \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v} (1+\epsilon) \|a_u - a_v\|_2^2$$

K-MEANS CLUSTERING

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

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

B_1, \dots, B_n be optimal clusters for
original data a_1, \dots, a_n

$$\begin{aligned}\widetilde{\text{Cost}}^* &\leq \widetilde{\text{Cost}}(B_1, \dots, B_n) \leq (1 + \epsilon) \text{Cost}(B_1, \dots, B_n) \\ &= (1 + \epsilon) \text{Cost}^*\end{aligned}$$

K-MEANS CLUSTERING

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

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

Then:

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

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?

CONNECTION TO DIMENSIONALITY REDUCTION

$$\|x_i\|_2 = \|x_j\|_2 = 1$$

Hard case: $x_1, \dots, x_n \in \mathbb{R}^d$ are all mutually orthogonal unit vectors:

$$\|x_i - x_j\|_2^2 = 2 \pm \epsilon \quad \text{for all } i, j.$$



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

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