

CS-GY 6763: Lecture 4

High Dimensional Geometry, the Johnson-Lindenstrauss Lemma

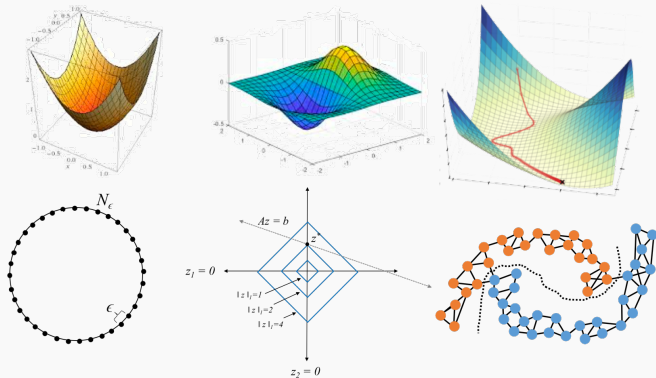
NYU Tandon School of Engineering, Prof. Christopher Musco

How do we deal with data (vectors) in high dimensions?

- Locality sensitive hashing for similarity search.
- Iterative methods for optimizing functions that depend on many variables.
- SVD + low-rank approximation to find and visualize low-dimensional structure.
- Convert large graphs to high dimensional vector data.

HIGH DIMENSIONAL IS NOT LIKE LOW DIMENSIONAL

Often visualize data and algorithms in 1,2, or 3 dimensions.



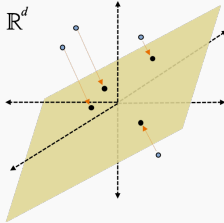
First part of lecture: Prove that high-dimensional space looks **very different** from low-dimensional space. These images are rarely very informative!

SKETCHING AND DIMENSIONALITY REDUCTION

Second part of lecture: Ignore our own advice.

Learn about **sketching, aka dimensionality reduction** techniques that seek to approximate high-dimensional vectors with much lower dimensional vectors.

- Johnson-Lindenstrauss lemma for ℓ_2 space.
- MinHash for binary vectors.

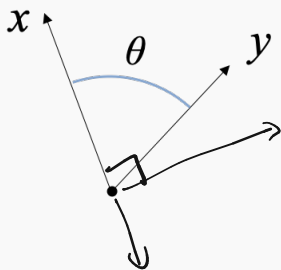
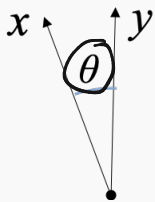


First part of lecture should help you understand the potential and limitations of these methods.

ORTHOGONAL VECTORS

Recall the inner product between two d dimensional vectors:

$$\langle \underline{x}, \underline{y} \rangle = \underline{x}^T \underline{y} = \underline{y}^T \underline{x} = \sum_{i=1}^d x[i]y[i]$$



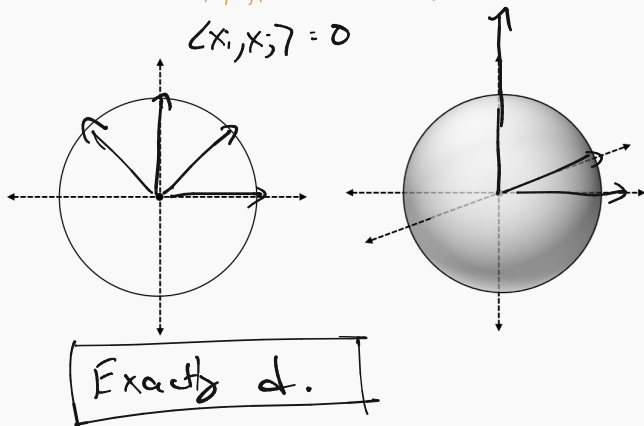
$$\langle x, y \rangle \approx 0$$

$$\theta \approx 90^\circ$$

$$\langle \underline{x}, \underline{y} \rangle = \cos(\theta) \cdot \|\underline{x}\|_2 \cdot \|\underline{y}\|_2$$

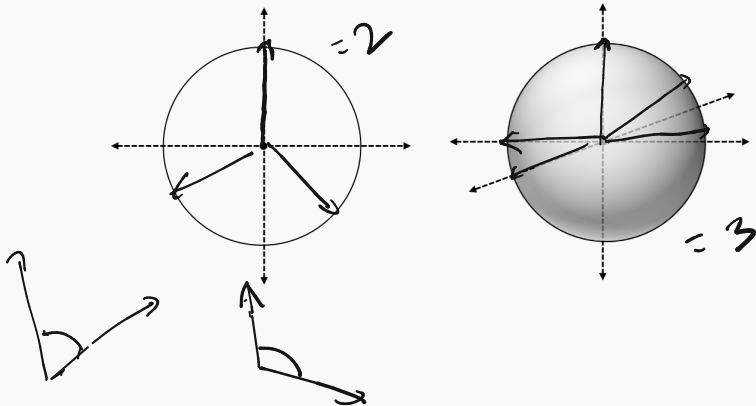
ORTHOGONAL VECTORS

What is the largest set of **mutually orthogonal** unit vectors x_1, \dots, x_t in d -dimensional space? I.e. with inner product $|x_i^T x_j| = 0$ for all i, j .



ORTHOGONAL VECTORS

What is the largest set **nearly orthogonal** unit vectors x_1, \dots, x_t in d dimensions. I.e., with inner product $|x_i \cdot x_j| \leq \epsilon$ for all i, j . Consider the case when ϵ is a constant. E.g. $\epsilon = 1/10$.



ORTHOGONAL VECTORS

What is the largest set **nearly orthogonal** unit vectors $\mathbf{x}_1, \dots, \mathbf{x}_t$ in d dimensions. I.e., with inner product $|\mathbf{x}_i^T \mathbf{x}_j| \leq \epsilon$ for all i, j . Consider the case when ϵ is a constant. E.g. $\epsilon = 1/10$.

1. d

2. $\Theta(d)$

3. $\Theta(d^2)$

4. $2^{\Theta(d)}$

ORTHOGONAL VECTORS

Claim: There is an exponential number of nearly orthogonal unit vectors in d dimensional space (i.e., $\sim 2^d$).

2^{cd} for $c < 1$

Proof strategy: Use the **Probabilistic Method!** For $t = 2^{\Theta(d)}$ define a random process which generates random vectors $(\mathbf{x}_1, \dots, \mathbf{x}_t)$ that are unlikely to have large inner product.

1. Claim that, with non-zero probability, $|\mathbf{x}_i^T \mathbf{x}_j| \leq \epsilon$ for all i, j .
2. Conclude that there must exist some set of t unit vectors with all pairwise inner-products bounded by ϵ .

PROBABILISTIC METHOD

Claim: There is an exponential number (i.e., $2^{\Theta(d)}$) of nearly orthogonal unit vectors in d dimensional space.

Proof: Let $\underline{x}_1, \dots, \underline{x}_t$ all have independent random entries, each set to $\pm \frac{1}{\sqrt{d}}$ with equal probability.

$$\cdot \|\underline{x}_i\|_2^2 = \sum_{j=1}^d [x_i[j]]^2 = \sum_{j=1}^d \frac{1}{d} = 1.$$

$$\cdot \mathbb{E}[x_i^T x_j] = \mathbb{E}\left(\sum_{k=1}^d x_i[k] x_j[k]\right) = \sum_{k=1}^d \mathbb{E}[x_i[k]] \mathbb{E}[x_j[k]] = 0.$$

$$\cdot \text{Var}[x_i^T x_j] = \sum_{k=1}^d \text{Var}(x_i[k] x_j[k]) = \frac{1}{d}$$

$$\begin{cases} \frac{1}{\sqrt{d}} & \text{w.p. } 1/2 \\ -\frac{1}{\sqrt{d}} & \text{w.p. } 1/2 \end{cases}$$

$$\frac{1}{d^2} \rightarrow C_k$$

INFORMAL PROOF

Let $Z = \underbrace{x_i^T x_j}_{\text{dot product}} = \sum_{i=1}^d C_i$ where each C_i is random $\pm \frac{1}{\sqrt{d}}$.

Z is a sum of many i.i.d. random variables, so looks approximately Gaussian. Roughly, we expect that:

$$\Pr[|Z - \mathbb{E}Z| \geq \alpha \cdot \sigma] \leq O(e^{-\alpha^2})$$

$$\Pr\{|Z| \gg \varepsilon\} \leq \underbrace{O(e^{-\varepsilon^2/d})}_{\downarrow \sqrt{d}} \quad \begin{array}{l} \varepsilon = \alpha \sigma = \frac{\alpha}{\sqrt{d}} \\ \alpha = \varepsilon \sqrt{d} \end{array}$$

t vectors, $\binom{t}{2}$ pairs, $\leq t^2$ pairs.

$$\Pr(\text{all inner products} < \varepsilon) \geq 1 - O(e^{-\varepsilon^2/d}) t^2$$

By a union bound, we can claim that the above holds for all pairs in a set of $t = \frac{1}{\sqrt{O(e^{-\varepsilon^2/d})}} = O(e^{\varepsilon^2/2d})$ vectors.

Use an exponential concentration inequality!

Theorem (Chernoff Bound) ↗

Let X_1, X_2, \dots, X_d be independent $\{0, 1\}$ -valued random variables and let $S = \sum_{i=1}^d X_i$. We have for any $\epsilon < 1$:

$$\Pr[|S - \mathbb{E}[S]| \geq \epsilon \mathbb{E}[S]] \leq 2e^{-\frac{\epsilon^2 \mathbb{E}[S]}{3}}.$$

Does not immediately apply because we have random variables that are $\pm 1/d$, not 0, 1.

Common trick: shift and scale to transform to the binary case.

FORMAL PROOF

$$\begin{aligned}
 \underline{x}_i^T \underline{x}_j = \underline{z} &= \sum_{i=1}^d C_i = \frac{2}{d} \sum_{i=1}^d \underbrace{\frac{d}{2}}_{E_i} \cdot C_i \\
 &= \frac{2}{d} \cdot \left(\sum_{i=1}^d (B_i - 1/2) \right) \\
 &= \frac{2}{d} \cdot \left(\underbrace{-\frac{d}{2}} + \underbrace{\sum_{i=1}^d B_i} \right)
 \end{aligned}$$

$\sum \frac{1}{d}$ w.p. $\frac{1}{2}$
 $\sum -\frac{1}{d}$ w.p. $\frac{1}{2}$

where each B_i is uniform in $\{0, 1\}$.

$$E_i = \begin{cases} \frac{1}{2} & \text{w.p. } \frac{1}{2} \\ -\frac{1}{2} & \text{w.p. } \frac{1}{2} \end{cases}$$

$$E_i + \frac{1}{2} = B_i$$

$$B_i = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ 0 & \text{w.p. } \frac{1}{2} \end{cases}$$

FORMAL PROOF

$$= \Pr\{Z > a\} + \Pr\{Z < -a\}$$

$$\underline{x_i^T x_j} = \underline{Z} = \frac{2}{d} \cdot \left(-\frac{d}{2} + \sum_{i=1}^d B_i \right)$$

where each B_i is uniform in $\{0, 1\}$.

$$\begin{aligned} \Pr[|Z| > \epsilon] &= \Pr\left[\sum_{i=1}^d B_i \geq \frac{d}{2} + \frac{\epsilon d}{2}\right] + \Pr\left[\sum_{i=1}^d B_i \leq \frac{d}{2} - \frac{\epsilon d}{2}\right] \\ &= \Pr\left[\sum_{i=1}^d B_i \geq (1 + \epsilon)\mathbb{E}[B_i]\right] + \Pr\left[\sum_{i=1}^d B_i \leq (1 - \epsilon)\mathbb{E}[B_i]\right] \end{aligned}$$

\mathbb{E}

$\sum_{i=1}^d B_i$

$\sum_{i=1}^d B_i$

Theorem (Chernoff Bound)

Let X_1, X_2, \dots, X_d be independent $\{0, 1\}$ -valued random variables and let $S = \sum_{i=1}^d X_i$. We have for any $\epsilon < 1$:

$$\Pr[|S - \mathbb{E}[S]| \geq \epsilon \mathbb{E}[S]] \leq 2e^{-\frac{\epsilon^2 \mathbb{E}[S]}{3}}.$$

$$S = \sum_{i=1}^d B_i;$$

$$\begin{aligned} \Pr[|S - \mathbb{E}[S]| \geq \epsilon \mathbb{E}[S]] &\leq 2e^{-\epsilon^2 \mathbb{E}[S]/3} \\ &= 2e^{-(\epsilon^2 d/2)/3} \\ &= 2e^{-\epsilon^2 d/6} \end{aligned}$$

Formally, using a Chernoff bound:

$$\Pr[|Z - \mathbb{E}Z| \geq \epsilon] \leq 2e^{-\epsilon^2 d/6}$$

For any i, j pair, $\Pr[|x_i^T x_j| < \epsilon] \geq 1 - 2e^{-\epsilon^2 d/6}$

By a union bound:

\downarrow
 χ_i

For all i, j pairs simultaneously, $\Pr[|x_i^T x_j| < \epsilon] \geq 1 - \binom{t}{2} \cdot 2e^{-\epsilon^2 d/6}$.

positive if $t < \frac{1}{\sqrt{2}} e^{\epsilon^2 d/12}$

$\geq 1 - t^2 2e^{-\epsilon^2 d/6}$

ORTHOGONAL VECTORS

Final result: In d -dimensional space, there are $2^{\Theta(\epsilon^2 d)}$ unit vectors with all pairwise inner products $\leq \epsilon$.

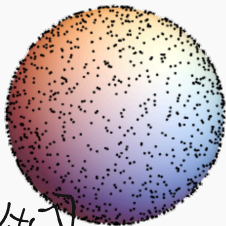
Corollary of proof: Random vectors tend to be far apart (and roughly equidistant) in high-dimensions.

$$\|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle$$

$$= 2 - 2\langle x, y \rangle$$

$$\approx 2 \pm 2\epsilon$$

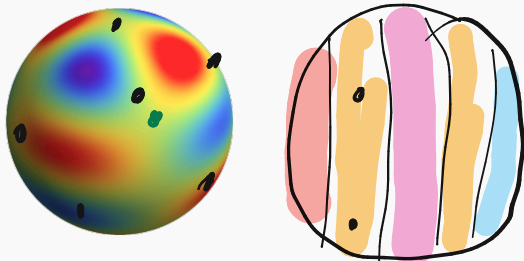
random unit vectors



$$\langle x - y, x - y \rangle = \langle x, x \rangle + \langle y, y \rangle - 2\langle x, y \rangle$$

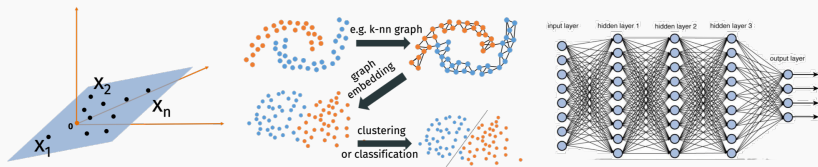
CURSE OF DIMENSIONALITY

Curse of dimensionality: Suppose we want to use e.g. k -nearest neighbors to learn a function or classify points in \mathbb{R}^d . If our data distribution is truly random, we typically need an exponential amount of data.



The existence of lower dimensional structure in our data is often the only reason we can hope to learn.

Low-dimensional structure.



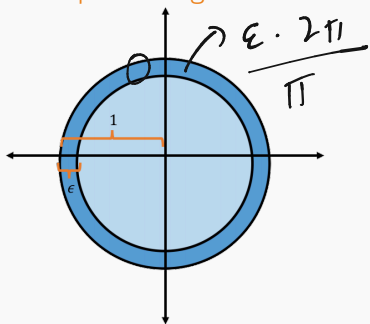
For example, data lies on low-dimensional subspace, or does so after transformation. Or function can be represented by a restricted class of functions, like neural net with specific architecture.

UNIT BALL IN HIGH DIMENSIONS

Let B_d be the unit ball in d dimensions:

$$B_d = \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}.$$

What percentage of volume of B_d falls with ϵ of its surface?



$$\frac{\text{vol}(B_d(1)) - \text{vol}(B_d(1-\epsilon))}{\text{vol}(B_d(1))}$$

$$= \frac{c \cdot 1^d - c \cdot (1-\epsilon)^d}{c \cdot 1^d}$$

$$= 1 - (1-\epsilon)^d$$

$$= 1 - \left((1-\epsilon)^{1/\epsilon} \right)^{d\epsilon}$$

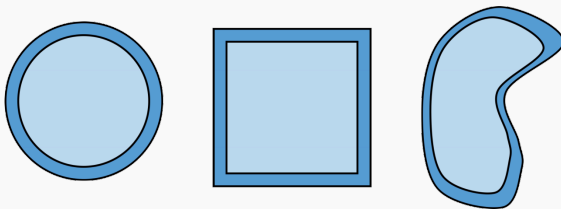
$$= 1 - (1/e)^{d\epsilon} = 1 - 2^{O(-d\epsilon)}$$

Volume of radius R ball is $\frac{\pi^{d/2}}{(d/2)!} \cdot R^d$.

ISOPERIMETRIC INEQUALITY

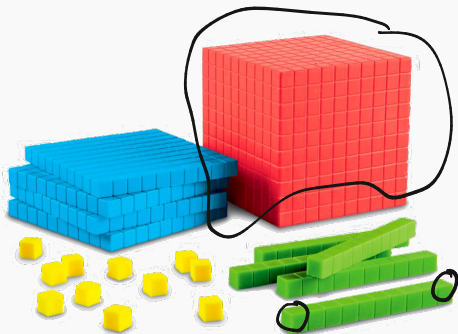
All but a vanishing small $2^{-\Theta(\epsilon d)}$ fraction of a unit ball's volume is within ϵ of its surface.

Isoperimetric Inequality: the ball has the minimum surface area/volume ratio of any shape.



- If we randomly sample points from any high-dimensional shape, nearly all will fall near its surface.
- 'All points are outliers.'

INTUITION



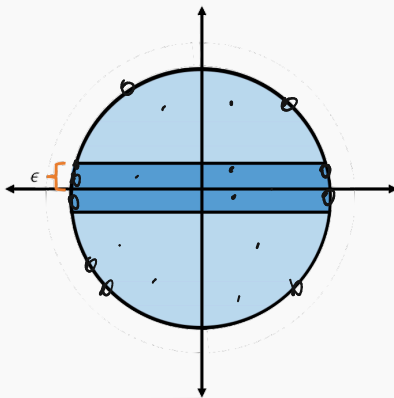
$$1D: \frac{\text{surface cubes}}{\text{total cubes}} = \frac{2}{10} = .2$$

$$2D: \frac{\text{surface cubes}}{\text{total cubes}} = \frac{10^2 - 2^2}{10^2} = \frac{100 - 64}{100} = .36$$

$$3D: \frac{\text{surface cubes}}{\text{total cubes}} = \frac{10^3 - 5^3}{10^3} = \frac{1000 - 512}{1000} = .488$$

SLICES OF THE UNIT BALL

What percentage of the volume of \mathcal{B}_d falls within ϵ of its equator?

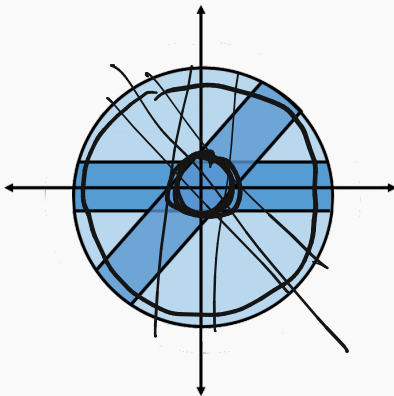


$$S = \{ \underline{x} \in \underline{\mathcal{B}}_d : \underline{|x_1|} \leq \epsilon \}$$

SLICES OF THE UNIT BALL

What percentage of the volume of \mathcal{B}_d falls within ϵ of its equator? Answer: all but a $2^{-\Theta(\epsilon^2 d)}$ fraction.

$$1 - 2^{-\Theta(\epsilon^2 d)}$$

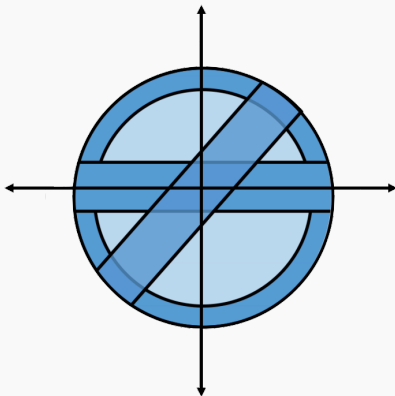


By symmetry, this is true for any equator:

$$S_{\mathbf{t}} = \{\mathbf{x} \in \mathcal{B}_d : \mathbf{x}^T \mathbf{t} \leq \epsilon\}.$$

BIZARRE SHAPE OF UNIT BALL

1. $(1 - 2^{-\Theta(\epsilon^d)})$ fraction of volume lies ϵ close to surface.
2. $(1 - 2^{-\Theta(\epsilon^2 d)})$ fraction of volume lies ϵ close to any equator.



High-dimensional ball looks nothing like 2D ball!

CONCENTRATION AT EQUATOR

Claim: All but a $2^{-\Theta(\epsilon^2 d)}$ fraction of the volume of the ball falls within ϵ of its equator.

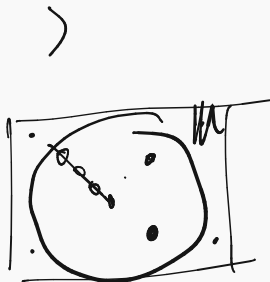
Equivalent: If we draw a point x randomly from the unit ball, $|x_1| \leq \epsilon$ with probability $\geq 1 - 2^{-\Theta(\epsilon^2 d)}$.

first entry of x

$(0, 0)$

$(-1, 1)$

$x(1)$



$\begin{bmatrix} .1 \\ .2 \\ -.1 \\ \vdots \end{bmatrix}$

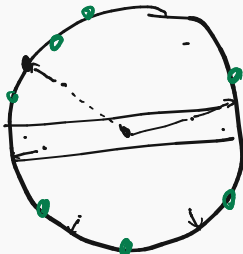
CONCENTRATION AT EQUATOR

Let $\underline{w} = \frac{\overset{\circlearrowleft}{x}}{\|x\|_2}$. Because $\|x\|_2 \leq 1$,

First entry of vector x .

$$\Pr[|x_1| \leq \epsilon] \geq \Pr[|w_1| \leq \epsilon].$$

Claim: $|w_1| \leq \epsilon$ with probability $\geq 1 - 2^{-\Theta(\epsilon^2 d)}$. This then proves our statement from the previous slide.



How can we generate w , which is a random vector taken from the unit sphere (the surface of the ball)?

IMPORTANT FACT IN HIGH DIMENSIONAL GEOMETRY

Rotational Invariance of Gaussian distribution: Let \mathbf{g} be a random Gaussian vector, with each entry drawn from $\mathcal{N}(0, 1)$. Then $\mathbf{w} = \mathbf{g}/\|\mathbf{g}\|_2$ is distributed uniformly on the unit sphere.

Why? Consider the probability density function of a high dimensional Gaussian:

$$\begin{aligned} P(\mathbf{g}) &= P(\mathbf{g}[1]) \cdot \dots \cdot P(\mathbf{g}[d]) = \prod_{i=1}^d c e^{-\mathbf{g}[i]^2/2} \\ &= c^d e^{-\sum_{i=1}^d \mathbf{g}[i]^2/2} \\ &= c^d e^{-\|\mathbf{g}\|_2^2/2} \end{aligned}$$

Handwritten notes: $c e^{-\mathbf{g}[i]^2/2}$ (with an arrow pointing from $P(\mathbf{g}[1])$ to this expression)

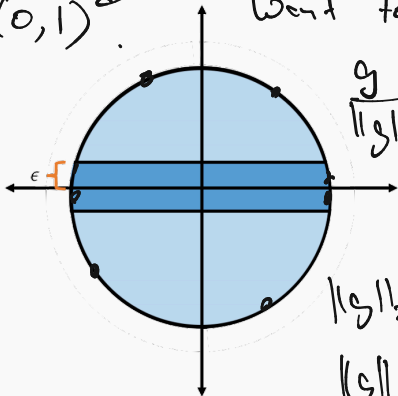
$\text{randN}(d, 1)$



PROOF STRATEGY:

Let $g \sim N(0, 1)^d$.

Want to show that



$$\frac{g_1}{\|g\|_2} \leq \epsilon$$

with prob.

$$1 - 2^{-O(\epsilon^2 d)}$$

$$\|g\|_2^2$$

$$\|g\|_2 \approx \sqrt{d}$$

1. Prove that with high probability, the first entry of g/\sqrt{d} is small.
2. Prove that g/\sqrt{d} is very very close to $g/\|g\|_2$, so this vector also has small first entry.

CONCENTRATION AT EQUATOR

Let \mathbf{g} be a random Gaussian vector and $\mathbf{w} = \mathbf{g}/\|\mathbf{g}\|_2$.

$$\cdot \mathbb{E}[\|\mathbf{g}\|_2^2] = \mathbb{E}\left[\sum_{i=1}^d g_i^2\right] = \sum_{i=1}^d \mathbb{E}[g_i^2] = \sum_{i=1}^d \text{Var}(g_i^2) = d$$

Excercise for home: Prove that $\Pr[\|\mathbf{g}\|_2^2 \leq \frac{1}{2}\mathbb{E}[\|\mathbf{g}\|_2^2]] \leq 2^{-\Theta(d)}$.

This should intuitively make sense. Can you tell me why?

$$\int x^k e^{-x^2/2} \dots \quad \text{Var}(g_i) = \mathbb{E}[g_i^2] - \mathbb{E}[g_i]^2 = 0$$

With high prob. $\|\mathbf{g}\|_2^2 \geq \frac{1}{2}d$
 $\rightarrow \|\mathbf{g}\|_2 \geq \frac{1}{\sqrt{2}}\sqrt{d}$

CONCENTRATION AT EQUATOR

For $1 - 2^{-\Theta(d)}$ fraction of vectors \mathbf{g} , $\|\mathbf{g}\|_2 \geq \sqrt{d/2}$. Condition on the event that we get a random vector in this set.
 If $*$ and $*$ holds then

Given this event:

$$\begin{aligned} \Pr[|w_1| \leq \epsilon] &= \Pr[|w_1| \cdot \sqrt{d/2} \leq \epsilon \cdot \sqrt{d/2}] \\ &\geq \Pr[|g_1| \leq \epsilon \cdot \sqrt{d/2}] \\ &\geq 1 - 2^{-\Theta((\epsilon \cdot \sqrt{d/2})^2)} \end{aligned}$$



By union bound, overall we have:

$$\Pr[|w_1| \leq \epsilon] \geq 1 - 2^{-\Theta(\epsilon^2 d)} - 2^{-\Theta(d)}$$

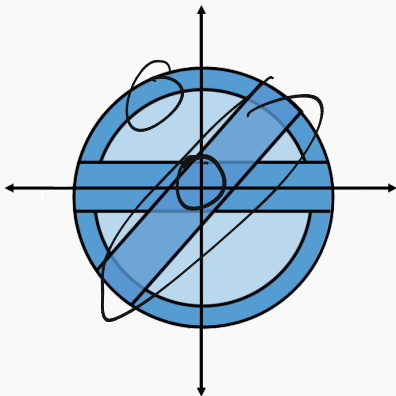
Recall: $w = \frac{\mathbf{g}}{\|\mathbf{g}\|_2}$. So after conditioning, we have $|w_1| \leq \frac{|g_1|}{\sqrt{d/2}}$.

$$= 1 - 2^{-\Theta(\epsilon^2 d)}$$

$$\rightarrow |w_1| \sqrt{d/2} \leq |g_1|$$

BIZARRE SHAPE OF UNIT BALL

1. $(1 - 2^{-\Theta(\epsilon^d)})$ fraction of volume lies ϵ close to surface.
2. $(1 - 2^{-\Theta(\epsilon^2 d)})$ fraction of volume lies ϵ close to any equator.



High-dimensional ball looks nothing like 2D ball!

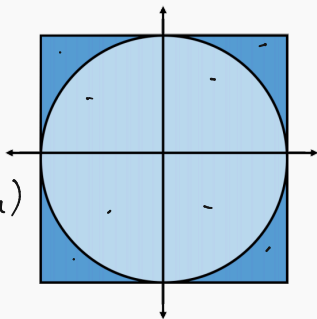
HIGH DIMENSIONAL CUBE

Let C_d be the d -dimensional cube:

$$C_d = \{x \in \mathbb{R}^d : |x(i)| \leq 1 \forall i\}.$$

$$O(d)^{O(d)}$$

$$(d/2)! \approx (d/4)^{(d/4)}$$



$$\frac{\sqrt{\pi}^d}{(d/2)!}$$

$$\cancel{B^d} 1$$

$$\underline{2^d}$$

$$\sqrt{\pi}^d$$

In two dimensions, the cube is pretty similar to the ball.

But volume of C_d is 2^d while volume of unit ball is $\frac{\sqrt{\pi}^d}{(d/2)!}$.

This is a huge gap! Cube has $O(d)^{O(d)}$ more volume.

Some other ways to see these shapes are very different:

- $\max_{\mathbf{x} \in \mathcal{B}_d} \|\mathbf{x}\|_2^2 = 1$

- $\max_{\mathbf{x} \in \mathcal{C}_d} \|\mathbf{x}\|_2^2 = \|(1 \ 1 \ \dots \ 1)\|_2^2 = d$

Some other ways to see these shapes are very different:

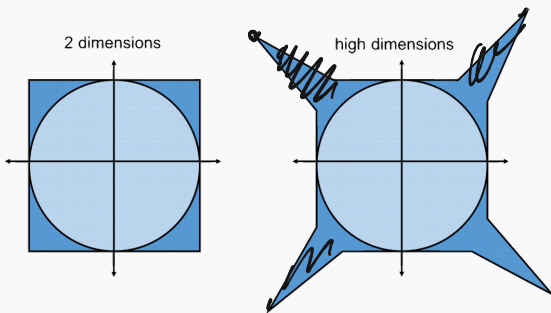
$$\begin{aligned} \bullet \mathbb{E}_{\mathbf{x} \sim \mathcal{B}_d} \|\mathbf{x}\|_2^2 &< 1 \\ \bullet \mathbb{E}_{\mathbf{x} \sim \mathcal{C}_d} \|\mathbf{x}\|_2^2 &= \mathbb{E} \sum_{i=1}^d \chi_i^2 \cdot \sum_{i=1}^d \mathbb{E}(\chi_i^2) = d/3. \end{aligned}$$

$$\chi_i \in [-1, 1]$$

$$= \frac{1}{2} \int_{-1}^1 x^2 dx = 1/3$$

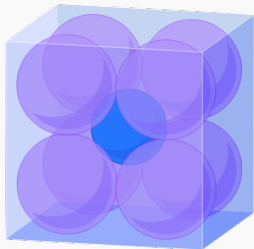
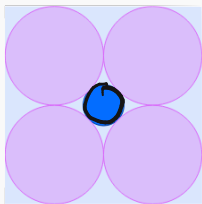
HIGH DIMENSIONAL CUBE

Almost all of the volume of the unit cube falls in its corners, and these corners lie far outside the unit ball.



$$\|g\|_r^2 \in (1 \pm 1/2) \|e\|_r^2 \text{ if } k_{\text{length}} \text{ is } O\left(\frac{\log(1/\delta)}{(1/2)^d}\right)$$

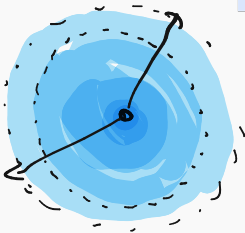
See [The Journey to Define Dimension](#) from Quanta Magazine for another fun example comparing cubes to balls!



If k_{length}
 $g = d,$
 $d = 4 \log_2(1/\delta)$

$$2^d = 2^4 \frac{1}{\delta}$$

$$\delta = \frac{1}{\alpha(2^d)}$$



Despite **all this** warning that low-dimensional space looks nothing like high-dimensional space, next we are going to learn about how to **compress high dimensional vectors to low dimensions**.

We will be very careful not to compress things too far. An extremely simple method known as Johnson-Lindenstrauss Random Projection pushes right up to the edge of how much compression is possible.

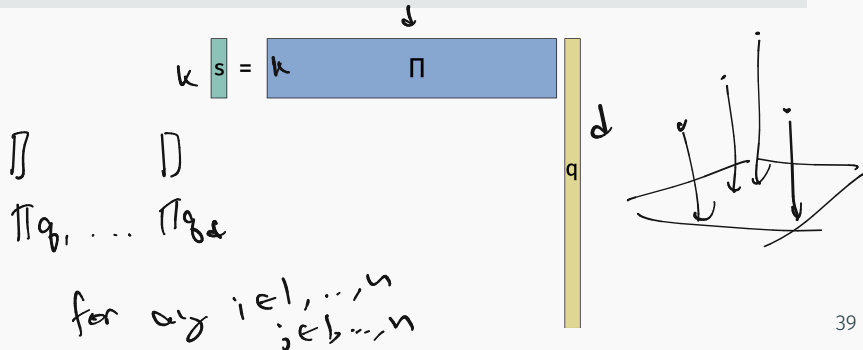
BREAK

EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Johnson-Lindenstrauss, 1984)

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

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



EUCLIDEAN DIMENSIONALITY REDUCTION

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 $(1 - \epsilon)^2 = 1 - O(\epsilon)$.

$$\downarrow \\ \epsilon \text{ or } 2\epsilon$$

$$\hookrightarrow \geq 1 - 2\epsilon$$

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 ,

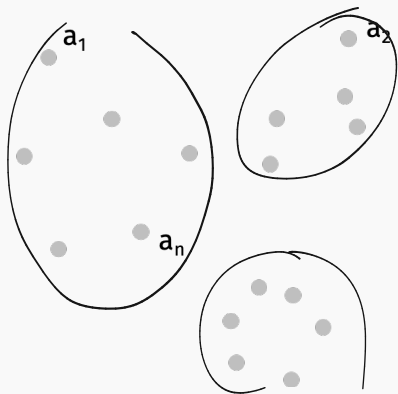
$$(1 - \epsilon) \|\Pi \mathbf{q}_i - \Pi \mathbf{q}_j\|_2^2 \leq \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \leq (1 + \epsilon) \|\Pi \mathbf{q}_i - \Pi \mathbf{q}_j\|_2^2.$$

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

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:

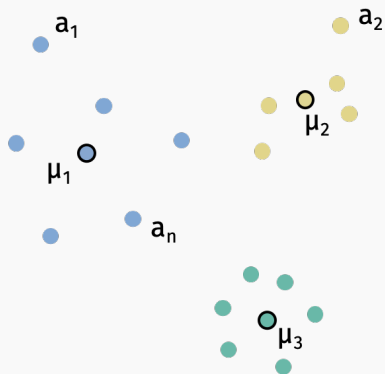
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



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:

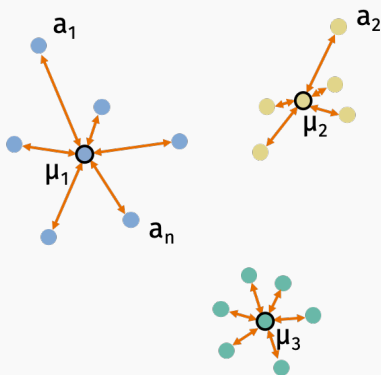
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



SAMPLE APPLICATION

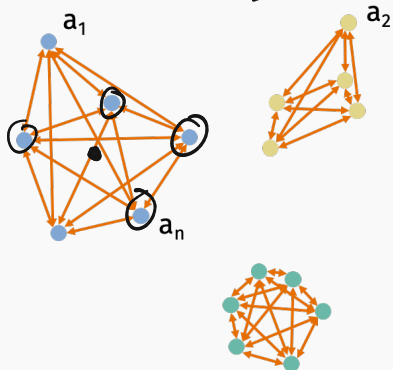
k-means clustering: Give data points $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^d$, find centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$ to minimize:


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



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

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

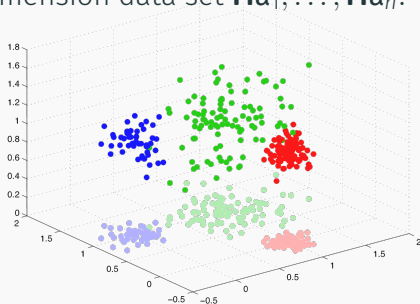


Exercise: Prove this to your self. 

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 $\mathbf{a}_1, \dots, \mathbf{a}_n$.

K-MEANS CLUSTERING

$$(1-\epsilon) \|a_u - a_v\|_h^2 \leq \|\Pi a_u - \Pi a_v\|_v^2 \leq (1+\epsilon) \|a_u - a_v\|_h^2$$

$$\underline{\text{Cost}}(C_1, \dots, C_k) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|a_u - a_v\|_2^2 \quad \left. \vphantom{\sum} \right) \text{focus!}$$

u, v

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

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

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

$$\widetilde{\text{Cost}}^* \leq (1 + \epsilon) \text{Cost}^*$$

$$\widetilde{\text{Cost}}^* = \widetilde{\text{Cost}}(\hat{C}_1^*, \dots, \hat{C}_k^*) \leq \widetilde{\text{Cost}}(C_1^*, \dots, C_k^*) \leq (1 + \epsilon) \text{Cost}(C_1^*, \dots, C_k^*)$$

K-MEANS CLUSTERING

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:

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

$$\text{Cost}^* = \min_{C_1, \dots, C_k} \text{Cost}(C_1, \dots, C_k) \text{ and}$$

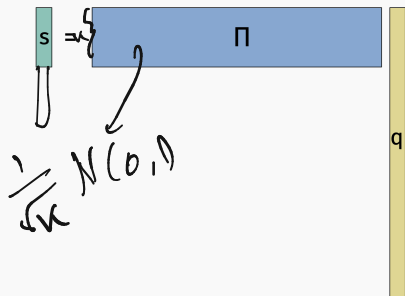
$$\widetilde{\text{Cost}}^* = \min_{C_1, \dots, C_k} \widetilde{\text{Cost}}(C_1, \dots, C_k)$$

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 ,

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



Remarkably, Π can be chosen completely at random!

One possible construction: Random Gaussian.

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

(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 $\mathbf{\Pi} \in \mathbb{R}^{k \times d}$ 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.

-2.1384	2.9888	-0.3338	0.0229	0.5201	-0.2938	-1.3320	-1.3617	-0.1952
-0.8396	0.8252	-0.0336	-0.2620	-0.0200	-0.8479	-2.3299	0.4550	-0.2176
1.3546	1.3790	-1.5771	-1.7502	-0.0348	-1.1201	-1.4491	-0.8487	-0.3031
-1.0722	-1.0582	0.0880	-0.2857	-0.7982	2.5268	0.3335	-0.3349	0.0230
0.9610	-0.4686	0.0820	-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	1.0904	-1.3337	-1.1564	-0.7145	-1.2571	-0.1393	-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	0.3502	-2.0026	-0.2248	-0.1765	-0.4762	0.6601	-0.2097
-1.2078	-2.0518	-0.2991	0.9642	-0.5890	0.7914	0.8620	-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

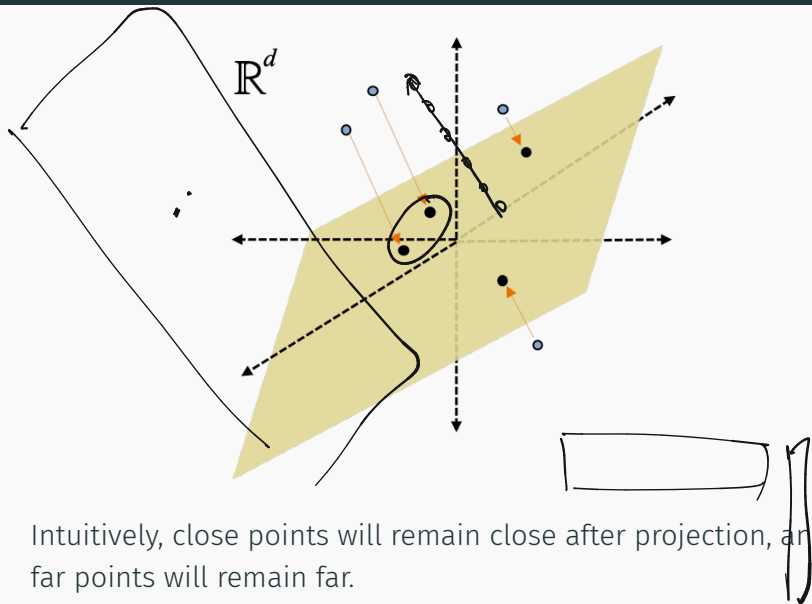
```
>> Pi = 2*randi(2,m,d)-3;
>> s = (1/sqrt(m))*Pi*q;
```

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

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.

$$\boxed{\text{random matrix}} = \boxed{\mathbf{I}}$$

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 $(1 - \delta)$:

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

$$\mathbf{x} = \mathbf{g}_i - \mathbf{g}_j$$

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

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:

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

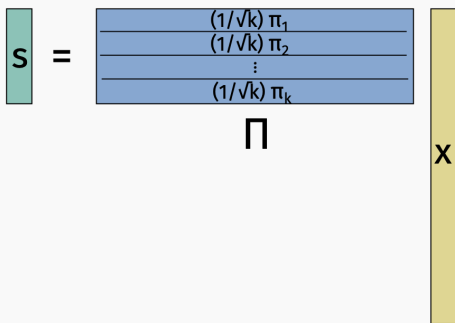
PROOF OF DISTRIBUTIONAL JL

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

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

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

Some notation:



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

Intermediate Claim:

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

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

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

We have 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: $\mathbb{E}\|\Pi\mathbf{x}\|_2^2 = \mathbb{E}\left[\sum (\langle \boldsymbol{\pi}_i, \mathbf{x} \rangle)^2\right]$

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

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.

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. Start on that next week!