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?

- High-dimensional similarity search.
- Iterative methods for optimizing functions in high-dimensions.
- SVD + low-rank approximation to find and visualize low-dimensional structure.
- Convert large graphs to high-dimensional vector data to uncover interesting things.

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! 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 (next class).



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

ORTHOGONAL VECTORS

Recall the inner product between two dimensional vectors:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^{\mathsf{T}} \mathbf{y} = \mathbf{y}^{\mathsf{T}} \mathbf{x} = \sum_{j=1}^{a} x_{j} \underline{y_{j}}$$



What is the largest set of mutually orthogonal unit vectors $\mathbf{x}_1, \ldots, \mathbf{x}_t$ in *d*-dimensional <u>space</u>? I.e. with inner product $|\mathbf{x}_i^T \mathbf{x}_i| = 0$ for all *i*, *j*.

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



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



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

Formally: In *d*-dimensional space, there are $2^{\Theta(\epsilon^2 d)}$ unit vectors with all pairwise inner products $\leq \epsilon$. $2^{\Theta(l/(d)^{\nu} \cdot d)} = 2^{O(r)}$

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

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

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

Proof: Let x_1 ... x_t all have independent random entries, each set to $\pm \frac{1}{\sqrt{d}}$ with equal probability. $\chi_{1} = (1/\sqrt{d} - 1/\sqrt{d} + 1/\sqrt{d})$ $\|\mathbf{x}_i\|_2^2 = \sum_{i=1}^d x_i (j_i)^2 = \sum_{j=1}^d \frac{1}{d} = 1$ $\mathbb{E}[\mathbf{X}_{i}^{\mathsf{T}}\mathbf{X}_{j}] = \mathbb{E}\left(\sum_{k=1}^{d} (\mathbf{X}_{i}(\mathbf{k}) \mathbf{X}_{j}(\mathbf{k}))\right) = \sum_{k=1}^{d} \mathbb{E}\left(\mathbf{X}_{i}(\mathbf{k}) \mathbf{X}_{j}(\mathbf{k})\right) = \sum_{k=1}^{d} \mathbb{E}\left(\mathbf{X}_{i}(\mathbf{k}) \mathbf{X}_{j}(\mathbf{k})\right) = \sum_{k=1}^{d} \mathbb{E}\left(\mathbf{X}_{i}(\mathbf{k}) \mathbf{X}_{j}(\mathbf{k})\right) = \mathbb{E}\left(\mathbf{X}_{i}(\mathbf{k}) \mathbf{X}_{i}(\mathbf{k})\right) = \mathbb{$ $\cdot \operatorname{Var}[X_{j}^{T}X_{j}] = \operatorname{Vcc}\left(\underbrace{\overset{d}{\underset{k=1}{\overset{}}}}_{k=1}^{T}\chi_{i}\left(k\right)\chi_{j}\left(k\right)\right) = \underbrace{\overset{d}{\underset{k=1}{\overset{}}}\operatorname{Vcc}\left(\chi_{i}\left(k\right)\chi_{j}\left(k\right)\right)$ $\chi_{i}(h)\chi_{i}(h) = \begin{cases} V_{\perp} & op & V_{\nu} \\ \sigma & \sigma & \sigma \end{cases}$ = 5 1/2 <

INFORMAL PROOF

Use an exponential concentration inequality!

Theorem (Chernoff Bound) Let $X_1, X_2, ..., X_d$ be independent (0, 1)-valued random variables and let $S = \sum_{i=1}^{d} X_i$. We have for any $\epsilon < 1$: $\Pr[|S - \mathbb{E}[S]| \ge \epsilon \mathbb{E}[S]] \le 2e^{\frac{-\epsilon^2 \mathbb{E}[S]}{3}}$.

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

$$\begin{aligned} = \frac{2}{d} \bigvee_{i=1}^{d} (i - \frac{d}{2}) &= \frac{2}{d} \bigvee_{i=1}^{d} (i - \frac{d}{2}) &= \frac{2}{d} \bigvee_{i=1}^{d} (i - \frac{d}{2}) &= \frac{2}{d} (i - \frac{d}{2}) &= \frac{2}$$

FORMAL PROOF

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

$$\Pr[\underline{|Z|} > \underline{\epsilon}] = \left(\Pr\left[\sum_{i=1}^{d} B_i \ge \frac{d}{2} + \frac{\epsilon d}{2}\right] + \Pr\left[\sum_{i=1}^{d} B_i \le \frac{d}{2} - \frac{\epsilon d}{2}\right]\right)$$

$$= \Pr\left[\sum_{i=1}^{d} B_i \ge (1+\epsilon)\mathbb{E}\left[\sum_{i=1}^{d} B_i\right]\right]$$

$$+ \Pr\left[\sum_{i=1}^{d} B_i \le (1-\epsilon)\mathbb{E}\left[\sum_{i=1}^{d} B_i\right]\right]$$

Theorem (Chernoff Bound)

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

$$(\Pr[|S - \mathbb{E}[S]| \ge \epsilon \mathbb{E}[S]] \le 2e^{\frac{-\epsilon^2 \mathbb{E}[S]}{3}}.)$$
Apply with $X_1, \dots, X_d = \underline{B}_1, \dots, \underline{B}_d$:

$$\Pr[|S - \mathbb{E}[S]| \ge \epsilon \mathbb{E}[S]] \le 2 \cdot e^{-c^2 \cdot (d/2)/3}$$

$$= 2 \cdot e^{-c^2 \cdot d/6}$$

$$|7| \le \epsilon \quad \text{w.p.} \quad 2e^{-\epsilon^2 \cdot d/6}$$

Conclusion from Chernoff bound:

For any *i*, *j* pair,
$$\Pr[|\mathbf{x}_i^T \mathbf{x}_j| < \epsilon] \ge 1 - 2e^{-\epsilon^2 d/6}$$
.

By a union bound:

For <u>all</u> *i*, *j* pairs simultaneously, $\Pr[|\mathbf{x}_i^T \mathbf{x}_j| < \epsilon] \ge 1 - {t \choose 2} \cdot 2e^{-\epsilon^2 d/6}$.

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: <u>Random vectors</u> tend to be far apart (and roughly equidistant) in high-dimensions.

$$\frac{(|x_{1} - x_{3}||_{v}^{v} - ||x_{3}||_{v}^{v} - 24x_{3}, x_{3})}{1} = \frac{2}{2} - 24x_{3}, x_{3}} = \frac{2}{2} - 24x_{3} - x_{3}, x_{3}}{\frac{1}{\sqrt{1}}}$$

$$= \frac{2}{\sqrt{1}} - 24x_{3} - x_{3}, x_$$

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 is our data is often the only reason we can hope to learn.

CURSE OF DIMENSIONALITY

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. Let \mathcal{B}_d be the unit ball in d dimensions:

$$\mathcal{B}_d = \{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \le 1 \}.$$

What percentage of volume of \mathcal{B}_d falls with ϵ of its surface?



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. $\mathfrak{O}(\mathfrak{f}_{\mathcal{F}})$



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

• 'All points are outliers.'

INTUITION



SLICES OF THE UNIT BALL



 $S = \{ \mathbf{x} \in \mathcal{B}_d : |\mathbf{x}_1| \le \epsilon \}$

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.



By symmetry, this is true for any equator: $S_{t} = \{ \mathbf{x} \in \mathcal{B}_{d} : \mathbf{x}^{\mathsf{T}} \mathbf{t} \leq \epsilon \}.$

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!

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 andomly from the unit ball, $|x_1| \le \epsilon$ with probability $\ge 1 - 2^{-\Theta(\epsilon^2 d)}$.



CONCENTRATION AT EQUATOR

Let
$$\underline{\mathbf{w}} = \underbrace{\mathbf{w}}_{\|\underline{\mathbf{x}}\|_{2}}$$
. Because $\|\underline{\mathbf{x}}\|_{2} \leq 1$, $\mathbf{w} \colon \underline{\mathbf{x}} \cdot \|\underline{\mathbf{x}}\|_{2}$,
 $\Pr[|\underline{x}_{1}| \leq \epsilon] \geq \Pr[|\underline{w}_{1}| \leq \epsilon]$.
Claim: $|w_{1}| \leq \epsilon$ with probability $\geq 1 - 2^{-\Theta(\epsilon^{2}d)}$. This then proves
our statement from the previous slide.
 $p(\mathbf{x})$
How can we generate w, which is a random vector taken from

the unit sphere (the surface of the ball)?

$$\omega = \left(N(0, 1), \dots, N(n, 1) \right) \quad \omega = \frac{\omega}{|1|^{\omega}|^{1}} 2^{7}$$

Rotational Invariance of Gaussian distribution: Let 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? Copsider the probability density function of a high - 2 Žg(i) -e-1/2 [1g](i) dimensional Gaussian: $p(g[d]) = \prod^{a} ce^{-g[i]^{2}/2}$ $c^{d} e^{-\|\mathbf{g}\|_{2}^{2}}$ · (\(0))) $-7 p(x) = C \cdot e^{-x^2/2}$

PROOF STRATEGY

Draw $\mathbf{g} \sim \mathcal{N}(\mathbf{0.I})$. Show that first entry of $\mathbf{w} = \mathbf{g}/\|\mathbf{g}\|_2 \le \epsilon$ with very high probability.



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

CONCENTRATION AT EQUATOR

Let **g** be a random Gaussian vector and $\mathbf{w} = \mathbf{g}(\|\mathbf{g}\|_{2})$ $\cdot \mathbb{E}[\|\mathbf{g}\|_{2}^{2}] = \mathcal{E}\left(\sum_{i=1}^{d} g(i)^{2}\right) + \sum_{i=1}^{d} \mathcal{E}\left(g(i)^{2}\right) = \mathcal{E}\left(g(i)^{2}\right)$ **Excersize for home:** Prove that $\Pr\left[\|\mathbf{g}\|_2^2 \le \frac{1}{2}\mathbb{E}[\|\mathbf{g}\|_2^2]\right] \le 2^{-\Theta(d)}$. This should intuitively make sense. Can you tell me why?

$$|w_{1}| \leq \alpha$$
 $|P_{1}(||y||_{2} \leq \frac{1}{2}d] \leq 2^{-\Theta(d)}$

CONCENTRATION AT EQUATOR

For
$$1-2^{-\Theta(d)}$$
 fraction of vectors \mathbf{g} , $\|\mathbf{g}\|_{2} \ge \sqrt{d/2}$. Condition on
the event that we get a random vector in this set. \mathbf{g}
Recall that $\mathbf{w} = \frac{\mathbf{g}}{\|\mathbf{g}\|_{2}}$. Given this event:
 $\Pr[\|w_{1}| \le \epsilon] = \Pr[\|g_{1}/\|\mathbf{g}\|_{2}| \le \epsilon]$
 $\ge \Pr[\|g_{1}|/\sqrt{d/2} \le \epsilon]$
 $= \Pr[\|g_{1}| \le \epsilon \cdot \sqrt{d/2}]$
 $\ge 1-2^{-\Theta((\epsilon \cdot \sqrt{d/2})^{2})} = 1 - 2^{O((\alpha^{2} + 1))}$
By union bound, overall we have:
 $\Pr[|w_{1}| \le \epsilon] \ge (1-2^{-\Theta(\epsilon^{2}d)}-2^{-\Theta(d)})$

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!

Let C_d be the *d*-dimensional cube:



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 $\sqrt{\pi^d}$. This is a huge gap! Cube has $O(d)^{O(d)}$ more volume. Some other ways to see these shapes are very different:

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

$$\cdot \max_{\mathbf{x} \in \mathcal{C}_d} \|\mathbf{x}\|_2^2 = \mathcal{J}$$

Some other ways to see these shapes are very different:

$$\begin{array}{c} \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{B}_{d}} \|\mathbf{x}\|_{2}^{2} & \mathbf{x} \\ \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{C}_{d}} \|\mathbf{x}\|_{2}^{2} = \mathbb{E}\left(\frac{1}{2} \mathbf{x}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) \\ = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) \\ = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) \\ = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{x}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{r}; \mathbf{r}, \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{x}; \mathbf{r}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{r}; \mathbf{r}; \mathbf{r}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{r}; \mathbf{r}; \mathbf{r}; \mathbf{r}\right) = \frac{1}{2} \mathbb{E}\left(\mathbf{r}; \mathbf{r}; \mathbf{r}\right) = \frac{1}$$

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



RECENT ARTICLE

See **The Journey to Define Dimension** from Quanta Magazine for another fun example comparing cubes to balls!



Place 2^d unit balls in box with side length 4. Look at sphere they enclose. It has radius $\sqrt{d-1}$. So for d > 9, it sticks out of the box... $\sqrt{d} - \frac{1}{\sqrt{3}}$ 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 <u>too</u> far. An extremely simple method known as Johnson-Lindenstrauss Random Projection pushes right up to the edge of how much compression is possible.

EUCLIDEAN DIMENSIONALITY REDUCTION

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 <u>for all</u> <u>i,j</u>,

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



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

Make pretty much any computation involving vectors faster and more space efficient.

- Faster vector search (used in image search, AI-based web search, Retrieval Augmented Generation (RAG), etc.).
- Faster machine learning (next class we will see an application to speeding up clustering).
- Faster numerical linear algebra.

Only useful if we can explicity construct a JL map **Π** and apply efficiently to vectors.

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

One possible construction: Random Gaussian.

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

The map **Π** is **oblivious to the data set**. This stands in contrast to other vector compression methods you might know like PCA.

[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 $\frac{1}{\sqrt{k}}\mathcal{N}(0,1)$ or each entry equals $\frac{1}{\sqrt{k}} \pm 1$ with equal probability.

-2.1384	2,9888	-0.3538	8.8229	8.5281	-0.2938	-1.3328	-1.3617	-0.1952
-8.8396	0.8252	-0,8236	-8,2620	-0.0200	-0,8479	-2,3299	0.4550	-0,2176
1.3546	1.3798	-1.5771	-1.7502	-0.0348	-1.1201	-1.4491	-0.8487	-0.3031
-1.0722	-1.0582	0.5080	-8.2857	-0.7982	2.5260	0.3335	-0.3349	0.0230
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	-0.1303	-1.1176	1.5270
-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	8.9642	-0.5898	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	1	-1	1
-1	1	-1	1	-1	1	1	-1	-1	1	-1	1	-1	-1	1	-1	1
-1 -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 **Q** also works. I.e. with $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_{k \times k}$. For this reason, the JL operation is often called a "random projection", even though it technically is not a projection when $\mathbf{\Pi}$'s entries are i.i.d. Can anyone see why Π is similar to a projection matrix? I.e., a matrix satisfying $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_{k \times k}$.

RANDOM PROJECTION



Intuition: Multiplying by a random matrix mimics the process of projecting onto a random *k* dimensional subspace in *d* dimensions.

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 <u>any vector **x**</u>, 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?

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 $1 - \delta$,

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

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 <u>for all</u> *i*, *j*, as long as we compress to:

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

PROOF OF DISTRIBUTIONAL JL

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:



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

PROOF OF DISTRIBUTIONAL JL

Intermediate Claim: Let π be a length d vector with $\mathcal{N}(0, 1)$ entries.

$$\mathbb{E}\left[\|\mathbf{\Pi}\mathbf{x}\|_{2}^{2}
ight] = \mathbb{E}\left[\left(\langle \boldsymbol{\pi}, \mathbf{x}
angle
ight)^{2}
ight].$$

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

$$\langle \boldsymbol{\pi}, \mathbf{X} \rangle = Z_1 \cdot x[1] + Z_2 \cdot x[2] + \ldots + Z_d \cdot x[d]$$

where each Z_1, \ldots, Z_d is a standard normal $\mathcal{N}(0, 1)$. We have that $Z_i \cdot x[i]$ is a normal $\mathcal{N}(0, x[i]^2)$ random variable.

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

What type of random variable is $\langle \pi, x \rangle$?

Fact (Stability of Gaussian random variables)

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

$$\langle \pi, \mathbf{x} \rangle = \mathcal{N}(0, x[1]^2) + \mathcal{N}(0, x[2]^2) + \ldots + \mathcal{N}(0, x[d]^2)$$

= $\mathcal{N}(0, \|\mathbf{x}\|_2^2).$

So
$$\mathbb{E} \| \mathbf{\Pi} \mathbf{x} \|_2^2 = \mathbb{E} \left[\left(\langle \boldsymbol{\pi}, \mathbf{x} \rangle \right)^2 \right] = \mathbb{E} \left[\mathcal{N}(0, \|\mathbf{x}\|_2^2) \right] = \|\mathbf{x}\|_2^2$$
, as desired.

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.

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

"Chi-squared random variable with k degrees of freedom."

Lemma

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

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

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 <u>dimensionality reduction</u> <u>possible?</u>

Doesn't Johnson-Lindenstrauss tell us that high-dimensional geometry can be approximated in low dimensions?

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

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

When we reduce to *k* dimensions with JL, we still expect these vectors to be nearly orthogonal. Why?

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

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

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

Alternative view: Without additional structure, we expect that learning/inference in *d* dimensions requires $2^{O(d)}$ data points. If we really had a data set that large, then the JL bound would be vacous, since $\log(n) = O(d)$.