

Note: These lecture notes are closely based on lecture notes by Sanjeev Arora [1] and Matt Weinberg [2].

1 Curse and Blessing of Dimensionality

High dimensional vectors are common in data mining and machine learning (e.g. items purchased by a Amazon customer, gene expression data). The phrase “curse of dimensionality” refers to the fact that algorithms are frequently harder to design in high-dimensional space - we’ve seen this with the convex hull algorithm. But, there is sometimes a flip side called “blessing of dimensionality”, wherein high-dimensional spaces can sometimes make life easier to analyze. For example, we can pack vectors more tightly in high-dimensional space, it is easier to route around obstacles there, and many random samples are more likely to be tightly clustered around a mean (e.g. via Chernoff bounds).

The fact is that high dimensional spaces behave differently than our intuition suggests (living as we are in 3-dimensional space). Following are some examples, but first some notation.

For a vector $x \in \mathbb{R}^d$, its ℓ_2 -norm is $|x|_2 = (\sum_i x_i^2)^{1/2}$ and ℓ_1 -norm is $|x|_1 = (\sum_i |x_i|)$. For any two vectors x, y , their Euclidean distance is $|x - y|_2$ and their Manhattan distance is $|x - y|_1$.

Some generalizations of geometric objects to higher dimensions:

- The n -cube in \mathbb{R}^d : $\{(x_1, \dots, x_d : 0 \leq x_i \leq 1)\}$. In \mathbb{R}^4 , if you are looking at one of the faces, say where $x_1 = 1$, then you are looking at a cube in \mathbb{R}^3 . The volume of the n -cube is 1.
- The unit n -ball in \mathbb{R}^d : $B_d = \{(x_1, \dots, x_d : \sum_i x_i^2 \leq 1)\}$. In \mathbb{R}^4 , if you slice through it with a hyperplane, say $x_1 = 1/2$, then this slice is a ball in \mathbb{R}^3 with radius of $\sqrt{1 - 1/2^2} = \sqrt{3}/2$. Every parallel slice also gives a ball. The volume of B_d is $\frac{\pi^{d/2}}{(d/2)!}$ (assuming d even). This is $\frac{1}{d^{\Theta(d)}}$

1.1 High Dimensionality Weirdness (and Intuition)

1.2 Near Orthogonal Vectors

How many “almost orthogonal” unit vectors can we have such that all pairwise angles lie between say 89 and 91 degrees? In \mathbb{R}^2 , the answer is 2. In \mathbb{R}^3 , it is 3. In \mathbb{R}^d , it is e^{cd} for some constant $c > 0$. Intuitively, to see this note that to get the angle close to 90, we just need to get the dot product of all vector pairs “close” to 0. When there are many entries in the vector, this is much easier to do. (more on this later).

1.2.1 Unit Ball

What is the ratio of the unit ball to its circumscribing cube (cube of side length 2)? In \mathbb{R}^2 , it is $\pi/4$ or about .78. In \mathbb{R}^3 it is $\pi/6$ or about .52. In d dimensions, it is $\frac{1}{d^{\Theta(d)}}/2^d = d^{-cd}$ for some constant $c > 0$.

2 Some Probability

Some tools from probability will be surprisingly useful for us to both get intuition about high dimensional geometry and also to do our projections to lower dimensional spaces. To start recall that a random variable (rv), X is informally a variable whose value depends on the outcome of some

random phenomena. Typically, random variables have a finite number of possible values in the real numbers, and we let X also refer to the set of possible outcomes. In this case, the expectation of a random variable, $E(X)$, is defined as $E(X) = \sum_{x \in X} x Pr(X = x)$.

First we prove linearity of expectation. Note that in the following lemma and proof, the random variables do *not* need to be independent. This makes the result extremely powerful.

Lemma 1. (*Linearity of Expectation*) Given a set of random variables X_1, \dots, X_n , $E(\sum_{i=1}^n X_i) = \sum_{i=1}^n E(X_i)$.

Proof: We first prove this for two random variables X and Y .

$$\begin{aligned} E(X + Y) &= \sum_{x \in X} \sum_{y \in Y} (x + y) Pr(X = x, Y = y) \\ &= \sum_{x \in X} \sum_{y \in Y} x \cdot Pr(X = x, Y = y) + \sum_{y \in Y} \sum_{x \in X} y \cdot Pr(X = x, Y = y) \\ &= \sum_{x \in X} x \cdot Pr(X = x) + \sum_{y \in Y} y \cdot Pr(Y = y) \\ &= E(X) + E(Y) \end{aligned}$$

The general result for n random variables now follows by induction. □

Lemma 2. (*Markov's Inequality*) Let X be a random variable that only takes on nonnegative values (i.e. $X \geq 0$ always). Then for any $\lambda > 0$,

$$Pr(X \geq \lambda) \leq \frac{E(X)}{\lambda}.$$

Proof: Assume not. Then for some value $\lambda > 0$, $Pr(X \geq \lambda) > \frac{E(X)}{\lambda}$. If this is true, then the expected value of X can be bounded as:

$$\begin{aligned} E(X) &\geq \sum_{i \geq \lambda} i Pr(X = i) \\ &\geq \sum_{i \geq \lambda} \lambda Pr(X = i) \\ &= \lambda Pr(X \geq \lambda) \\ &> \lambda \frac{E(X)}{\lambda} \\ &= E(X) \end{aligned}$$

But this sequence of inequalities implies that $E(X) > E(X)$, which is clearly a contradiction. □

2.1 Chernoff Bounds

The following important bound only works for independent random variables. We prove it for 0/1-valued random variables, which only take on the values 0 or 1, and we prove an upper bound. The lemma generalizes easily to also bound the probability of deviation below the mean.

Lemma 3. (Chernoff bounds) Let X_1, \dots, X_n be independent 0/1-valued random variables and let $p_i = E(X_i)$, where $0 \leq p_i < 1$ for all i . Then the sum $X = \sum_i X_i$, which has mean $\mu = E(X) = \sum_i p_i$ satisfies

$$\Pr(X \geq (1 + \delta)\mu) \leq (c_\delta)^\mu,$$

where $c_\delta = \frac{e^\delta}{(1+\delta)^{1+\delta}}$.

Proof: Consider an arbitrary positive constant t , to be set later, and consider the random variable e^{tX} . (If $X = 2$, say, this rv is e^{2t} .) A nice property of this random variable is the following:

$$\begin{aligned} E(e^{tX}) &= E(e^{t\sum_i X_i}) \\ &= E\left(\prod_i e^{tX_i}\right) \\ &= \prod_i E(e^{tX_i}) \end{aligned}$$

The last inequality holds since the X_i random variables are independent, and hence so are the e^{tX_i} random variables; and since $E(XY) = E(X)E(Y)$ if X and Y are independent. Note that

$$E(e^{tX_i}) = (1 - p_i) + p_i e^t.$$

Thus, we have:

$$\begin{aligned} \prod_i E(e^{tX_i}) &= \prod_i [1 + p_i(e^t - 1)] \\ &\leq \prod_i e^{p_i(e^t - 1)} \\ &\leq e^{\mu(e^t - 1)} \end{aligned}$$

In the above, the second step holds by the inequality $1 + x \leq e^x$ (via Taylor expansion of e . Recall that $e^x = 1 + x + x^2/2! + x^3/3! + \dots$). Now, we apply Markov's inequality to the random e^{tX} to get:

$$\begin{aligned} \Pr(X \geq (1 + \delta)\mu) &= \Pr(e^{tX} \geq e^{t(1+\delta)\mu}) \\ &\leq \frac{e^{\mu(e^t - 1)}}{e^{t(1+\delta)\mu}} \\ &\leq e^{\mu((e^t - 1) - t(1+\delta))} \end{aligned}$$

Recall that Markov's inequality says that for any positive random variable Y , and any $\lambda > 0$,

$$\Pr(Y \geq \lambda) \leq E(Y)/\lambda.$$

We let $Y = e^{tX}$, and note that $E(Y) \leq e^{\mu(e^t - 1)}$; and we let $\lambda = e^{t(1+\delta)\mu}$.

This holds for any positive t , and is minimized when $t = \ln(1 + \delta)$ (to see this, differentiate to get the minimum). This gives the lemma statement. \square

Using a symmetric argument, we can bound the probability of deviation below the mean. Combining the results and using some approximations gives the following extremely useful lemma.

Lemma 4. Let X_1, \dots, X_n be independent Poisson trials such that $P(X_i = 1) = p_i$. Let $X = \sum_i X_i$ and $\mu = E(X)$. Then for $0 \leq \delta \leq 1$,

$$Pr(|X - \mu| \leq \delta\mu) \leq 2e^{-\mu\delta^2/3}$$

A related inequality, which is even more closely like the central limit theorem, is the Bernstein inequality below. (From <https://www.cs.princeton.edu/~smattw/Teaching/Fa19Lectures/lec3/lec3.pdf> and <http://cseweb.ucsd.edu/~klevchen/techniques/chernoff.pdf> original proof due to Van Vu at UCSD, local copy available at <https://www.cs.unm.edu/~saia/classes/506-s20/lec/bernstein.pdf>)

Lemma 5. (Bernstein Inequality) Let X_1, \dots, X_n be independent random variables with $E(X_i) = 0$ and $|X_i| \leq 1$ for all i . Let $X = \sum_i X_i$, $\sigma_i^2 = E(X_i^2) - (E(X_i))^2$ and $\sigma^2 = \sum_i \sigma_i^2$. Then for all $0 \leq k \leq \sigma/2$ we have:

$$Pr(|X| \geq k\sigma) \leq 2e^{-k^2/4}$$

2.2 Example Using Chernoff Bounds

Assume we flip a fair coin n times and let X be the number of heads. Note that $E(X) = n/2$. Then by Chernoff bounds, we have that:

$$Pr(|X - n/2| \leq \delta n/2) \leq 2e^{-n\delta^2/6}$$

Q: What is the smallest value of δ that still ensures that we have polynomially small probability?

A: To ensure this, need $2e^{-n\delta^2/6} \leq n^{-1}$, which means that $-n\delta^2/6 \leq -\ln n$.

How about $\delta = 1$: we get $-n/6 \leq -\ln n$ which works

How about $\delta = 1/\sqrt{n}$: we get $-n(1/n)/6 = \Theta(1)$

How about $\delta = \sqrt{(\ln n)/n}$: we get $-n(\ln n)/n/6 = \Theta(-\ln n)$. That works!

2.3 Union Bounds

The following tool is frequently useful in conjunction with Chernoff bounds.

Lemma 6. (Union Bounds) Consider n events ξ_1, \dots, ξ_n . Then we have that

$$Pr(\cup_i \xi_i) \leq \sum_{i=1}^n Pr(\xi_i)$$

Proof: We'll show this for two events, the lemma statement then holds by an inductive argument. Let ξ_1 and ξ_2 be any two events. Then we have that

$$\begin{aligned} Pr(\xi_1 \cup \xi_2) &= Pr(\xi_1) + Pr(\xi_2) - Pr(\xi_1 \cap \xi_2) \\ &\leq Pr(\xi_1) + Pr(\xi_2) \end{aligned}$$

□

3 Number of Almost Orthogonal Vectors

One of the *benefits* of high-dimensional spaces are that they are very “roomy”. For example, we now show that there are $\Theta(e^d)$ vectors in \mathbb{R}^d that are “almost” orthogonal. Recall that the angle, θ , between two vectors can be found via the identity $\cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|}$, where $\|\cdot\|$ is the 2-norm.

Lemma 7. *Let a be a unit vector in \mathbb{R}^n . Let $x = (x_1, \dots, x_n)$ be a unit vector in \mathbb{R}^n created by choosing each x_i independently and uniformly in $\{\frac{-1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\}$. Let $X = a \cdot x = \sum_i a_i x_i$. Then for all $t > 0$,*

$$Pr(|X| > t) < 2e^{-nt^2/4}.$$

Proof: Note that $E(X) = E(\sum_i a_i x_i) = 0$. This is true since $E(a_i x_i) = 1/2(a_i \frac{-1}{\sqrt{n}}) + 1/2(a_i \frac{1}{\sqrt{n}}) = 0$. Also, since $\sigma^2 = E(X^2) - (E(X))^2 = E(X^2)$.

$$\begin{aligned} \sigma^2 &= E\left(\left(\sum_{i=1}^n a_i x_i\right)^2\right) \\ &= E\left(\sum_{1 \leq i \leq j \leq n} a_i a_j x_i x_j\right) \\ &= \sum_{1 \leq i \leq j \leq n} a_i a_j E(x_i x_j) \\ &= \sum_{1 \leq i \leq n} a_i^2 E(x_i^2) + \sum_{1 \leq i \neq j \leq n} a_i a_j E(x_i x_j) \\ &= \sum_{1 \leq i \leq n} a_i^2 (1/n) \\ &= 1/n \end{aligned}$$

For the second to last step, note that if $i \neq j$, $E(x_i x_j) = (1/2) * (1/n) + (1/2)(-1/n) = 0$, and if $i = j$, $E(x_i^2) = 1/n$. Thus, using Bernstein’s inequality, we see that

$$Pr(|X| > t) < 2e^{-(t/\sigma)^2/4} \leq 2e^{-nt^2/4}$$

□

From the above, the dot product of any unit vector $x \in \mathbb{R}^n$ with a “randomly chosen” vector is “small” with high probability. Since the cosine of two unit vectors x and y equals $x \cdot y$, we have the following:

Lemma 8. *Let $\epsilon > 0$ be a fixed constant. Consider a set S of $e^{\epsilon^2 n/10}$ vectors in \mathbb{R}^n , where each entry is independently and uniformly chosen in $\{\frac{-1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\}$. For any pair of vectors $x, y \in S$, let $\theta_{x,y}$ be the angle between x and y . Then for all $x, y \in S$,*

$$Pr(|\cos \theta_{x,y}| > \epsilon) \leq e^{-\epsilon^2 n/21}$$

Proof: Consider some fixed pair of vectors $x, y \in S$. Let $\xi_{x,y}$ be the event that $x \cdot y > \epsilon$. Note that $Pr(|\cos \theta_{x,y}| > \epsilon) = Pr(|x \cdot y| > \epsilon)$ Thus, by Lemma 7,

$$Pr(|\cos \theta_{x,y}| > \epsilon) < 2e^{-\epsilon^2 n/4}$$

Now let ξ be the event that *any* pair of vertices violates the bound. In particular, $\xi = \cup_{x,y \in S} \xi_{x,y}$. Then by a Union bound, we have:

$$\begin{aligned} Pr(\xi) &\leq \sum_{x,y \in S, x \neq y} Pr(\xi_{x,y}) \\ &\leq |S|^2 2e^{-\epsilon^2 n/4} \\ &\leq 2e^{\epsilon^2 n/5} e^{-\epsilon^2 n/4} \\ &\leq 2e^{-\epsilon^2 n/20} \\ &\leq e^{-\epsilon^2 n/21} \end{aligned}$$

where the last step holds for n sufficiently large. \square

4 Dimension Reduction

In this problem, we're given n points $v_1, \dots, v_n \in \mathbb{R}^d$ and a fixed $\epsilon > 0$. We want to find a function $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$, where $m \ll d$ such that for all i and j :

$$|v_i - v_j| \leq |f(v_i) - f(v_j)| \leq (1 + \epsilon)|v_i - v_j|$$

In other words, the distances between points are (approximately) preserved.

Note that many naive ideas fail to achieve this such as: (1) taking a random sample of m coordinates out of d ; and (2) partition coordinates into m subsets and add up the values in each subset.

Idea 1 fails for the case where we have vector $x = (0, 0, \dots, 1)$ and $y = (1, 0, 0, \dots, 0)$. Note that $|x - y| = 1$, but any random sample of coordinates is unlikely to find the 1 entry in either of these vectors. Idea 2 fails for the case that $x = (0, 1, 0, 1, \dots)$ and $y = (1, 0, 1, 0, \dots)$. Note that $|x - y|$ is large but these sums would be very close.

4.1 Johnson-Lindenstrauss Projection

Let G be a m by d matrix where each entry is a normal random variable, i.e. $G_{i,j} \sim \mathcal{N}(0, 1)$. Let $\Pi = \frac{1}{\sqrt{m}}G$ and let

$$f(x) = \Pi x.$$

So each entry in $f(v)$ equals $v \cdot g$ for some vector g filled with scaled Normal random variables (note that Gaussian and Normal are synonymous). Other (simpler) approaches also work (See Section 4.5 below).

4.2 Analysis

4.3 Reduction to Norm Preservation

Our main lemma is below. Note that, by taking square roots, this theorem implies that

$$(1 - \epsilon)|x| \leq |\Pi x| \leq (1 + \epsilon)|x|.$$

(For example, by Theorem 1 we have that:

$$\sqrt{(1 - \epsilon)}|x| \leq |\Pi x|$$

This implies that

$$(1 - \epsilon)|x| \leq \sqrt{(1 - \epsilon)}|x| \leq |\Pi x|$$

Distance Preservation: Then to prove distance preservation, we note that by the linearity of $f = \Pi$,

$$|f(x) - f(y)| = |\Pi(x) - \Pi(y)| = |\Pi(x - y)|$$

So with probability $1 - \delta$, we preserve the distance of one pair by Theorem 1. Then we'll do a union bound over all pairs, which will increase the error probability by $\binom{n}{2}$.

4.4 Main Theorem

Theorem 1. ((ϵ, δ) -JL property) If $m = 9 \log(1/\delta)/\epsilon^2$ then, with probability $1 - \delta$, for any vector x ,

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

Proof: Let $w = \Pi x$. Then we have:

$$|w|^2 = |\Pi x|^2 = \left| \frac{1}{\sqrt{m}} Gx \right|^2 = \frac{1}{m} \sum_{i=1}^m w_i^2.$$

Consider the i -th entry of w , which we can write as:

$$w_i = \sum_{j=1}^d x_j g_j$$

where each $g_j \sim \mathcal{N}(0, 1)$. So $E(w_i) = \sum_{j=1}^d x_j E(g_j) = 0$. Recall that $\text{var}(X) = E(X^2) - E^2(X)$. Thus $\text{Var}(w_i) = E(w_i^2)$. It follows that

$$\text{Var}(w_i) = E(w_i^2) = \sum_{j=1}^d \text{Var}(x_j g_j) = \sum_{j=1}^d x_j^2 \text{Var}(g_j) = \sum_{j=1}^d x_j^2 = |x|^2.$$

The above follows since for independent random variables X and Y , $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$. Thus,

$$E(w) = E\left(\frac{1}{m} \sum_{i=1}^m w_i^2\right) = \frac{1}{m} \sum_{i=1}^m E(w_i^2) = \frac{1}{m} \sum_{i=1}^m |x|^2 = |x|^2$$

Now we make use of the following fact about normal random variables:

Fact 1: If X and Y are independent and $X \sim \mathcal{N}(0, a^2)$ and $Y \sim \mathcal{N}(0, b^2)$, then $X + Y \sim \mathcal{N}(0, a^2 + b^2)$. The property that the sum of Normal distributions remains normal is known as *stability*.

By this fact, $w_i \sim \mathcal{N}(0, |x|^2)$. It follows that w_i^2 is a χ^2 (chi-squared) random variable, and that $|w|^2 = \frac{1}{m} \sum_{i=1}^m w_i^2$ is a chi-squared random variable with m degrees of freedom. These random variables are very well studied and they concentrate around their mean essentially as well as a Normal random variable. In particular, if $X = \frac{1}{m} \sum_{i=1}^m w_i^2$, then we have the following tail bound¹ for any positive ϵ :

$$P(|X - E(X)| \geq \epsilon v) \leq 2e^{-m\epsilon^2/8}$$

So if we set $m = \Theta(\log(1/\delta)/\epsilon^2)$, we get that $|\Pi x|$ satisfies

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

¹See, e.g., <https://www.stat.berkeley.edu/~mjbain/stat210b/Chap2.TailBounds-Jan22-2015.pdf>

with probability at least $1 - \delta$.

In particular, let $m = 9 \log(1/\delta)/\epsilon^2$, then we know that

$$\begin{aligned} P(|X - E(X)| \geq \epsilon v) &\leq 2e^{-(9 \log(1/\delta)/\epsilon^2)(\epsilon^2/8)} \\ &= 2e^{-((9/8) \log(1/\delta))} \\ &= 2(\delta)^{9/8} \\ &\leq \delta \end{aligned}$$

where the last step holds for δ sufficiently small. In particular, we want $2\delta^{9/8} \leq \delta$. Dividing both sides by δ , we see that this holds when $2\delta^{1/8} \leq 1$ or $\delta^{1/8} \leq 1/2$ or $\delta \leq (1/2)^8$ or $\delta < 1/256$. \square

Now we can prove the main theorem.

Theorem 2. Assume we are given n points $v_1, \dots, v_n \in \mathbb{R}^d$ and a fixed $\epsilon > 0$. Let $m = O(\log n/\epsilon^2)$ and set $f = \Pi$, where Π is a m by d matrix of independent $\mathcal{N}(0, 1)$ random variables. Then, with probability $1 - 1/n$, for any i and j , $1 \leq i < j \leq n$:

$$(1 - \epsilon)|v_i - v_j| \leq |f(v_i) - f(v_j)| \leq (1 + \epsilon)|v_i - v_j|$$

Proof: Set $\delta = 1/n^3$ and $m = 27 \log n/\epsilon^2$. For any fixed pair of points v_i and v_j , let $\xi_{i,j}$ be the (bad) event that the following does not hold:

$$|v_i - v_j| \leq |f(v_i) - f(v_j)| \leq (1 + \epsilon)|v_i - v_j|$$

Then by Theorem 1, $Pr(\xi_{i,j}) \leq 1/n^3$. Let ξ be the (bad) event that $\xi_{i,j}$ occurs for any v_i and v_j .

Then, by a Union bound, we know that

$$\begin{aligned} Pr(\xi) &\leq \sum_{i,j} Pr(\xi_{i,j}) \\ &= \binom{n}{2} \frac{1}{n^3} \\ &\leq 1/n. \end{aligned}$$

\square

Interestingly, this bound is tight. Noga Alon has shown (in 2018) that there are point sets that can't be embedded in less than $O(\log n/\epsilon^2)$ dimensions if we want to preserve pairwise distances.

4.5 Another Simpler Johnson-Lindenstrauss Projection

Here is a simpler Johnson-Lindenstrauss Algorithm for projection that also works.

1. $x_1, \dots, x_m \leftarrow$ vectors in \mathbb{R}^m chosen as follows. Each coordinate is chosen independently and randomly from $\left\{ \sqrt{\frac{1+\epsilon}{m}}, -\sqrt{\frac{1+\epsilon}{m}} \right\}$
2. $u_i[j] \leftarrow x_i \cdot u_i$ for all $i : 1 \leq i \leq n$ and $j : 1 \leq j \leq m$

In other words, $u_i = (z_i \cdot x_1, \dots, z_i \cdot x_m)$ for $i = 1, \dots, m$. Note that we can think of this as a linear transformation $u = Az$ where A is a matrix with random and independent entries in $\left\{ \sqrt{\frac{1+\epsilon}{m}}, -\sqrt{\frac{1+\epsilon}{m}} \right\}$.

4.6 Analysis

We now do a “sketch” of the analysis. The following lemma shows that things work out well in expectation.

Lemma 9. For any $1 \leq i < j \leq n$, $E(|u_i - u_j|^2) = (1 + \epsilon)|z_i - z_j|^2$

Proof: According to the projection, we have the following for any $1 \leq i < j \leq n$:

$$|u_i - u_j|^2 = \sum_{k=1}^m \left(\sum_{\ell=1}^n (z_i[\ell] - z_j[\ell])x_k[\ell] \right)^2$$

Fix i and j . Let $z = z_i - z_j$ and let $u = u_i - u_j$. Then for any $1 \leq k \leq m$, we have

$$\begin{aligned} E(|u \cdot x_k|^2) &= E \left(\left(\sum_{\ell=1}^n (z[\ell]x_k[\ell]) \right)^2 \right) \\ &= \sum_{\ell} \sum_{\ell'} E(z[\ell]x_k[\ell]z[\ell']x_k[\ell']) \\ &= \sum_{\ell=1}^n E((z[\ell]x_k[\ell])^2) \\ &= \frac{1 + \epsilon}{m} |z|^2 \end{aligned}$$

Hence, by linearity of expectation $E(|u|^2) = (1 + \epsilon)|z|^2$. □

The rest of the analysis follows similar to that in Theorem 2. First, one establishes a (harder) tail-bound around this expectation and then does a union bound over all pairs of points. In this way, we can get the same result as Theorem 2.

5 Applications of JL Projection

- Approximate all-pairs distances in $O(n \log n + nd)$ vs $O(n^2d)$ time
- Approximate distance-based clustering
- Approximate support vector machine (SVM) classification
- Approximate Linear Regression

Note: For some of these Machine Learning type applications, we need it to be the case that distances are approximately preserved across *all* (infinite) vectors in the vector space. Thus, a simple union bound won't work and instead we need to make use of a technique called ϵ -nets. We discuss this technique below.

6 Linear Regression and ϵ -Nets

The following is the classic least-squares regression problem.

Given: n data vectors $a_1, \dots, a_n \in \mathbb{R}^d$, and n response values $y_1, \dots, y_n \in \mathbb{R}$. Let A be a $n \times d$ matrix with rows a_1, \dots, a_n ; let y be a length n vector with entries y_1, \dots, y_n .

Goal: Find $x \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n (a_i \cdot x - y_i)^2 = |Ax - y|^2$$

Usually, this problem requires $O(nd^2)$ time to solve (for example, by using singular value decomposition). We now show how to speed it up by reducing n using Johnson-Lidenstrauss.²

Let Π be chosen from the family of matrices from Theorem 2. To obtain an approximate solution, we solve the “sketched” problem where we find $x \in \mathbb{R}^d$ to minimize:

$$|\Pi Ax - \Pi y|^2.$$

This can be solved in $O(md^2)$ time (once ΠA and Πy are computed - we haven’t discussed this but there are JL transforms which are also fast, since they are sparse). We want to prove that a solution to this smaller problem is a good approximation to the big problem. Note that the following lemma is a direct consequence of Theorem 1, applied to the vector $Ax - y$:

Lemma 10. *Let $m = O(\log(1/\delta)/\epsilon^2)$. Then, for any particular vector x , with probability $1 - \delta$,*

$$(1 - \epsilon)|Ax - y|^2 \leq |\Pi Ax - \Pi y|^2 \leq (1 + \epsilon)|Ax - y|^2$$

Now if we could show this was true *for all* x , we’d be done. In particular, let x^* be the optimal solution for the original problem, and let \tilde{x}^* be the solution for the sketched problem. Then we’ve have, with probability $1 - \delta$.

$$|A\tilde{x}^* - y|^2 \leq \frac{1}{1 - \epsilon} |\Pi A\tilde{x}^* - \Pi y|^2 \leq \frac{1}{1 - \epsilon} |\Pi Ax^* - \Pi y|^2 \leq \frac{1 + \epsilon}{1 - \epsilon} |Ax^* - y|^2.$$

In the above the first and last inequalities hold via Lemma 10, and the middle inequality holds by noting that \tilde{x}^* minimizes $|\Pi Ax - \Pi y|$ over all vectors x .

If $\epsilon \leq .25$, then $\frac{1+\epsilon}{1-\epsilon} \leq 1 + 3\epsilon$, so we can get an approximation to the original regression problem. Q: Why do we need a bound for all x above??? The main problem is that \tilde{x}^* depends on the projection π , and so it’s not fixed ahead of time. How do we extend Lemma 10 to all x ? We can’t use union bounds since there are an infinite number of possible vectors x .

7 Beyond Union Bounds

Recall that we have $A \in \mathbb{R}^{n \times d}$ and want to approximately find x to minimize $|Ax - y|^2$, by instead solving the sketched problem $|\Pi Ax - \Pi y|^2$. We want to argue that for all $x \in \mathbb{R}^d$,

$$(1 - \epsilon)|Ax - y|^2 \leq |\Pi Ax - \Pi y|^2 \leq (1 + \epsilon)|Ax - y|^2 \tag{1}$$

But proving this requires establishing a JL-bound for an infinity of possible vectors, which clearly can’t be shown via union bounds. Instead, we use a different approach.

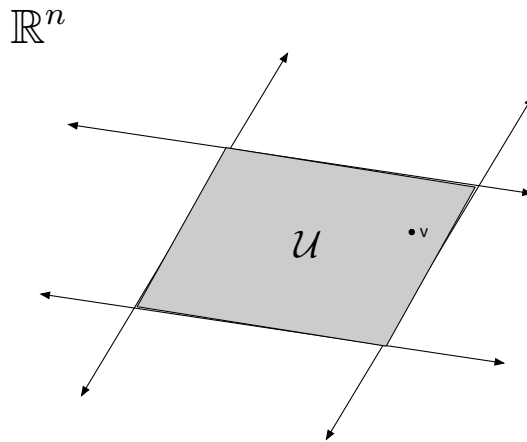


Figure 1. JL approximately preserves distances over any subspace \mathcal{U} of dimension d contained in \mathbb{R}^n

7.1 Subspace Embeddings

We will prove a more general statement that implies equation 1, and is useful in other applications.

Theorem 3. Let $\mathcal{U} \subset \mathbb{R}^n$ be a d -dimensional linear subspace in \mathbb{R}^n . If $\Pi \in \mathbb{R}^{m \times n}$ is chosen from any distribution \mathcal{D} satisfying Theorem 1, then with probability $1 - \delta$,

$$(1 - \epsilon)|v| \leq |\Pi v| \leq (1 + \epsilon)|v| \quad (2)$$

for all $v \in \mathcal{U}$, as long as $m = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$

(Note that it's possible to prove a slightly tighter bound of $m = O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$ that we won't discuss here.)

How does this theorem imply equation 1? We can apply it to the $d + 1$ dimensional subspace spanned by the d columns of A and the vector y . Every vector formed by inputting some vector x into the linear equation $Ax - y$ lies in this $d + 1$ dimensional subspace. So for the regression problem, we require dimension $m = O\left(\frac{(d+1) \log(1/\epsilon)}{\epsilon^2}\right)$. In particular, we can approximately solve linear regression over $n \gg d$ examples for the same amount of work as $O(d)$ examples, for fixed ϵ .

7.2 An Example

Let $n = 3$ and \mathcal{U} could be the 2 dimensional subspace spanned by $(1, -1, 1)$ and $(1, 1, -1)$. JL will basically find a low-dimensional sub-space that is not much higher than the dimensionality of \mathcal{U} .

7.3 Reduction to a Sphere

We first note that Theorem 3 holds so long as equation 2 holds for all points on the unit sphere in \mathcal{U} . This is a consequence of linearity of the Euclidean norm. In particular, denote the sphere $\mathcal{S}_{\mathcal{U}}$ as

$$\mathcal{S}_{\mathcal{U}} = \{v \mid v \in \mathcal{U} \text{ and } |v| = 1\}.$$

²Note that we are reducing n (number of vectors) and not d (dimension). Since we only care about the matrix A , you could think of n as the dimension and d as the number of vectors.

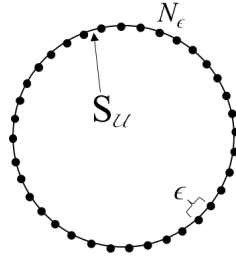


Figure 2. An ϵ -net N_ϵ for a sphere in a 2-dimensional subspace of \mathcal{U}

Now any point $v \in \mathcal{U}$ can be written as cx for some scalar c and some point $x \in \mathcal{S}_\mathcal{U}$. If $(1 - \epsilon)|x| \leq |\Pi x| \leq (1 + \epsilon)|x|$, then $c(1 - \epsilon)|x| \leq c|\Pi x| \leq c(1 + \epsilon)|x|$ and so $(1 - \epsilon)|cx| \leq |\Pi cx| \leq (1 + \epsilon)|cx|$.

Note that the last inequality holds since $|cx| = \sqrt{\sum_i (cx)_i^2} = c\sqrt{\sum_i x_i^2} = c|x|$ since x was on the unit sphere.

7.4 Constructing a Net

We prove Theorem 3 by showing that there is a large but finite set of points $N_\epsilon \subset \mathcal{S}_\mathcal{U}$ such that if equation 2 holds for all $v \in N_\epsilon$. The set N_ϵ is called an ϵ -net. In particular, we will show the following.

Lemma 11. For any positive $\epsilon < 1$, there exists a set $N_\epsilon \subset \mathcal{S}_\mathcal{U}$ with $|N_\epsilon| = \left(\frac{4}{\epsilon}\right)^d$ such that $\forall v \in \mathcal{S}_\mathcal{U}$,

$$\min_{x \in N_\epsilon} |v - x| \leq \epsilon.$$

Proof: We use the following greedy procedure to construct N_ϵ (note that this construction is just for proof of existence, our algorithms do not need to implement this). Initially $N_\epsilon \leftarrow \{\}$. Then:

- While there is a point $v \in \mathcal{S}_\mathcal{U}$ with distance greater than ϵ from any point in N_ϵ , add v to N_ϵ .

After running this procedure, we have $|N_\epsilon|$ points such that $\min_{x \in N_\epsilon} |v - x| \leq \epsilon$ for all $v \in \mathcal{S}_\mathcal{U}$. So we just need to bound $|N_\epsilon|$.

To do so, we first lower bound the volume taken up by balls around points in $N_\epsilon = \{x_1, x_2, \dots, x_{|N_\epsilon|}\}$. In particular, note that for all $i \neq j$, $|x_i - x_j| \geq \epsilon$. If not, then either x_i or x_j would not have been added to N_ϵ by our greedy algorithm. So if we place balls of radius $\epsilon/2$ around each x_i :

$$B(x_1, \epsilon/2) \dots B(x_{|N_\epsilon|}, \epsilon/2)$$

then for all $i \neq j$, $B(x_i, \epsilon/2)$ does not intersect $B(x_j, \epsilon/2)$.

So how do we now set up an inequality to bound $|N_\epsilon|$??? The volume of a d dimensional ball of radius r is cr^d for some fixed constant c . Thus, the amount of space taken up by all the balls surrounding points in N_ϵ is $c|N_\epsilon|(\epsilon/2)^d$.

Next note that the amount of space that these balls can exist in is at most the volume of a d dimensional sphere with radius $1 + \epsilon/2$. This volume is $c(1 + \epsilon/2)^d$.

Thus we have that

$$|N_\epsilon|c(\epsilon/2)^d \leq c(1 + \epsilon/2)^d$$

Solving for $|N_\epsilon|$, we have that

$$\begin{aligned} |N_\epsilon| &\leq \frac{c(1 + \epsilon/2)^d}{c(\epsilon/2)^d} \\ &\leq \frac{(1 + \epsilon/2)^d}{(\epsilon/2)^d} \\ &\leq \left(\frac{4}{\epsilon}\right)^d \end{aligned}$$

□

7.5 Proving Theorem 3

We can now prove Theorem 3 by extending to all vectors in the subspace.

Proof: Choose $m = O\left(\frac{\log(|N_\epsilon|/\delta)}{\epsilon^2}\right) = O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right)$ so that Equation 2 holds for all $x \in N_\epsilon$ (via Theorem 1 and a union bound).

Now consider any $v \in \mathcal{S}_U$. We claim that for some $x_0, x_1, x_2, \dots \in N_\epsilon$ that we can write v as:

$$v = x_0 + c_1 x_1 + c_2 x_2 + \dots$$

for constants c_1, c_2, \dots where $|c_i| \leq \epsilon^i$. To see this, note that there is some point x_0 within distance ϵ of v . Next we need to represent $v - x_0$, which has norm at most ϵ . So instead, we can write represent the point $\frac{v-x_0}{|v-x_0|}$, which has norm 1 and multiply the resulting coefficients by ϵ . Again there is some point x_1 within distance ϵ of *this* point. Continuing this process ad infinitum gives the claim.

Now, we can consider $|\Pi v|$ and make use of the triangle inequality in order to complete the proof.

$$\begin{aligned} |\Pi v| &= |\Pi(x_0 + c_1 x_1 + c_2 x_2 + \dots)| \\ &= |\Pi x_0 + \Pi c_1 x_1 + \Pi c_2 x_2 + \dots| \\ &\leq |\Pi x_0| + |\Pi c_1 x_1| + |\Pi c_2 x_2| + \dots \\ &\leq (1 + \epsilon)|x_0| + (1 + \epsilon)c_1|x_1| + (1 + \epsilon)c_2|x_2| + \dots \\ &\leq (1 + \epsilon)(|x_0| + c_1|x_1| + c_2|x_2| + \dots) \\ &\leq (1 + \epsilon)(1 + \epsilon + \epsilon^2 + \dots) \\ &\leq (1 + O(\epsilon)) \end{aligned}$$

In the above, the third step follows by the triangle inequality. The fourth step follows by the fact that each $x_i \in N_\epsilon$ (and so by a Union bound their norms are all approximately preserved). The last line follows since $\epsilon + \epsilon^2 + \epsilon^3 + \dots$ is a geometric summation, which has value that is a constant times its largest term.

The other direction of the proof is symmetric. It is included below for completeness.

$$\begin{aligned} |\Pi v| &= |\Pi(x_0 + c_1 x_1 + c_2 x_2 + \dots)| \\ &\geq |\Pi x_0| - \epsilon|\Pi x_1| - \epsilon^2|\Pi x_2| - \dots \\ &\geq (1 - \epsilon) - \epsilon(1 + \epsilon) - \epsilon^2(1 + \epsilon) + \dots \\ &\geq 1 - O(\epsilon) \end{aligned}$$



7.6 Other Applications of JL

Speed up Winnow by projecting “training data”??? Yes

Speed up Boosting by projecting “training data”??? Yes

Speed up Winnow by projecting attributes??? Not necessarily

Approximate solutions to System of Linear equations? Sometimes

Finding an ϵ -approximate convex hull?? Sometimes

References

- [1] Sanjeev Arora. Advanced Algorithm Design Class, Princeton University, 2013. <https://www.cs.princeton.edu/courses/archive/fall15/cos521/>.
- [2] Matt Weinberg. Dimensionality Reduction and the Johnson-Lindenstrauss Lemma, 2019. <https://www.cs.princeton.edu/~smattw/Teaching/Fa19Lectures/lec9/lec9.pdf>.