

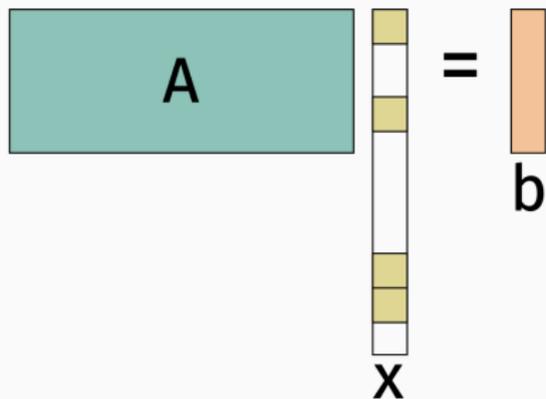
CS-GY 9223 I: Lecture 14

High Dimensional Geometry (+ finish up compressed sensing)

NYU Tandon School of Engineering, Prof. Christopher Musco

SPARSITY RECOVERY/COMPRESSED SENSING

Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n$, $\mathbf{b} \in \mathbb{R}^m$ with $\mathbf{b} = \mathbf{A}\mathbf{x}$ Recover \mathbf{x} , under the assumption that it is k -sparse.



What properties of \mathbf{A} let us solve this problem efficiently.

- **Measurement efficiency:** Small m .
- **Computational efficiency:** polynomial in n, k .

BASIC RESULT FROM LAST CLASS

If \mathbf{A} is matrix satisfying the $(O(k), O(1))$ - **Restricted Isometry Property** then \mathbf{x} can be uniquely recovered from $\mathbf{b} = \mathbf{Ax}$ in polynomial time using the basis pursuit linear program:

$$x = \arg \min_z \|\mathbf{z}\|_1 \quad \text{subject to} \quad \mathbf{Az} = \mathbf{b}.$$

Measurement complexity $m \approx O(k \log n)$ for:

- Random matrices (JL matrices)
- Subsampled Fourier matrices

In other words, this is how many rows these matrices need to satisfy $(O(k), O(1))$ -RIP.

APPLICATIONS

Random matrices used in algorithms. E.g. to solve the heavy-hitters problem:

The diagram shows a matrix equation $Ax = b$. The matrix A is a 4x10 grid of values (+1, -1) with columns labeled 1, 2, 3, ..., n. The vector x is a 10x1 column of yellow boxes. The vector b is a 4x1 orange box.

+1	-1	-1	+1	+1	+1	-1	+1	-1	-1
-1	+1	+1	+1	-1	+1	-1	-1	+1	-1
+1	-1	-1	+1	-1	-1	-1	+1	+1	+1
-1	-1	+1	+1	+1	-1	+1	-1	-1	+1
1	2	3	...	n					

x

b

Subsampled Fourier matrices: $b = Ax$ is the evaluation of the Fourier transform Fx at a set of random frequencies $f_1, \dots, f_m \in \{0, \dots, n-1\}$. If we had entirety of Fx could recover x using inverse Fourier transform.

Input: Stream of number: 1, 2, 1, 5, 1, 1, 6,

Output: Majority element (if one exists)

- $item = 0; count = 0;$
- For each e in our stream:
 - If $count == 0, item = e; count = 1;$
 - Else if $item == e; count = count + 1;$
 - Else $count = count - 1;$
- Return $item$

Misra-Gries algorithm

Slick generalization of this algorithm for finding any element e which appears more than $1/k$ of the time. Uses $O(k)$ space.

A lot of interest in developing even faster algorithms that avoid using the “heavy hammer” of linear programming and run in even faster than $O(n^{3.5})$ time.

- **Iterative Hard Thresholding:** Looks a lot like projected gradient descent. Solve $\min_z \|\mathbf{Az} - \mathbf{b}\|$ while continually projecting z back to the set of k -sparse vectors. Runs in time $\sim O(nk \log n)$ for Gaussian measurement matrices and $O(n \log n)$ for subsampled Fourier matrices.
- Other “first order” type methods: Orthogonal Matching Pursuit, CoSaMP, Subspace Pursuit, etc.

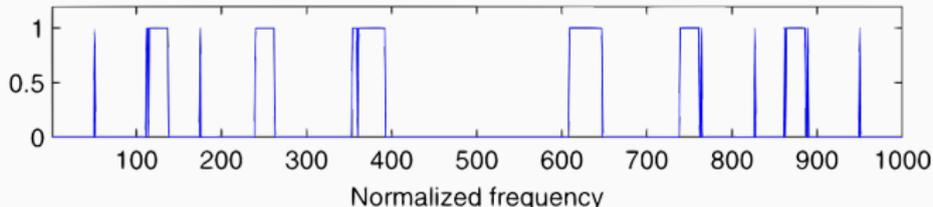
When \mathbf{A} is a subsampled Fourier matrix, there are now methods that run in $O(k \log^c n)$ time [Hassanieh, Indyk, Kapralov, Katabi, Price, Shi, etc. 2012+].

SPARSE FOURIER TRANSFORM

Corollary: When \mathbf{x} is k -sparse, we can compute the inverse Fourier transform $\mathbf{F}^*\mathbf{F}\mathbf{x}$ of $\mathbf{F}\mathbf{x}$ in $O(k \log^c n)$ time!

- Randomly subsample $\mathbf{F}\mathbf{x}$.
- Feed that input into our sparse recovery algorithm to extract \mathbf{x} .

Fourier and inverse Fourier transforms in sublinear time when the output is sparse.



Applications in: Wireless communications, GPS, protein imaging, radio astronomy, etc. etc.

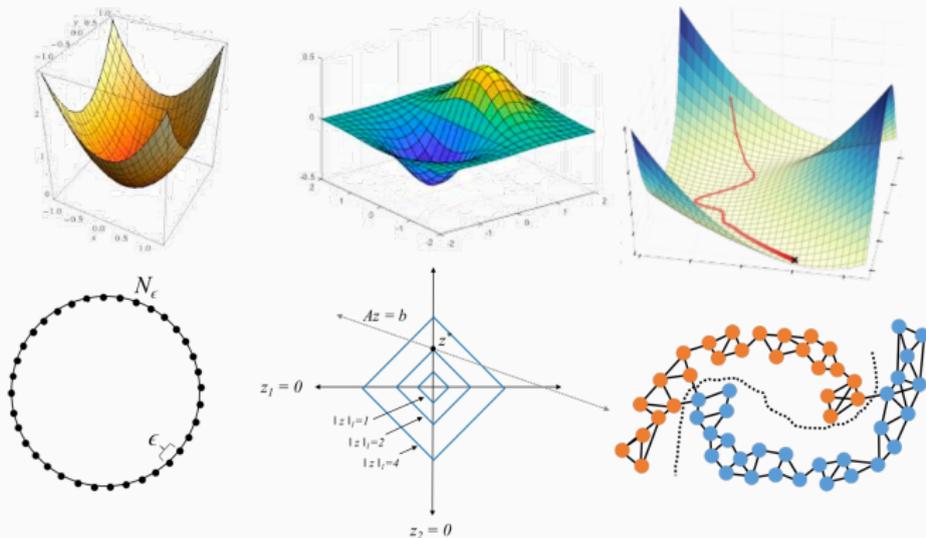
HIGH DIMENSIONAL GEOMETRY

How do we deal with data in high dimensions?

- Randomized sketching + dimensionality reduction.
- 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.

VISUALIZING HIGH DIMENSIONAL DATA

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



This is not always a good thing to do: high-dimensional space looks **very different** from low-dimensional space.

What is the largest set of mutually orthogonal unit vectors in d -dimensional space?

What is the largest set of unit vectors in d -dimensional space with inner product $|\mathbf{x}^T \mathbf{y}| \leq \epsilon$?

1. d

2. $\Theta(d)$

3. $\Theta(d^2)$

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

Claim: There is an exponential number of nearly orthogonal unit vectors in d dimensional space.

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

- $\|\mathbf{x}_i\|_2 =$

- $\mathbb{E}[\mathbf{x}_i^T \mathbf{x}_j] =$

For any i, j pair, $\Pr[|\mathbf{x}_i^T \mathbf{x}_j| < \epsilon] \geq 1 - 2e^{-\epsilon^2 d/3}$.

By a union bound:

For all i, j pairs simultaneously, $\Pr[|\mathbf{x}_i^T \mathbf{x}_j| < \epsilon] \geq 1 - t^2 \cdot 2e^{-\epsilon^2 d/3}$.

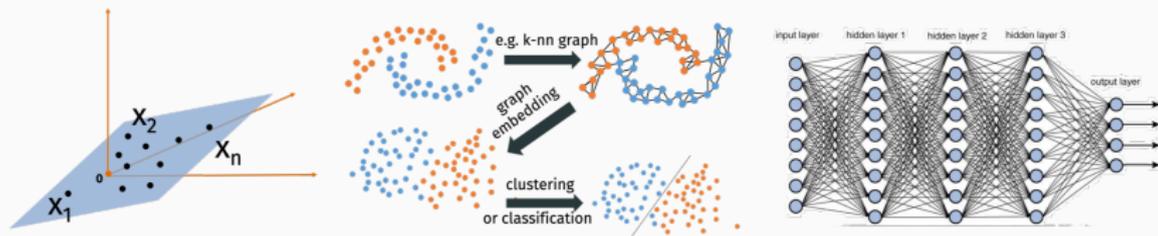
Final result: In d -dimensional space, there are $2^{\theta(\epsilon^2 d)}$ unit vectors with all pairwise inner products $\leq \epsilon$.

Corollary: Random vectors are all approximately the same distance from each other.

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.

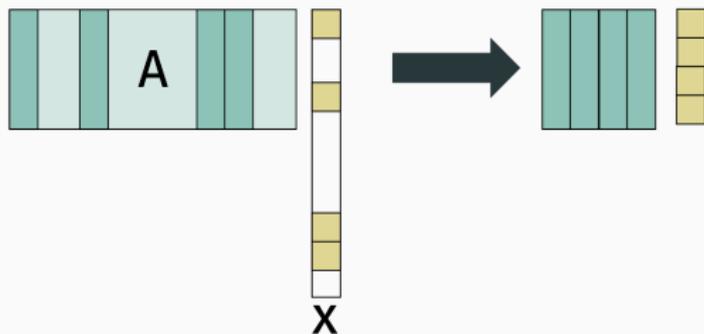


VALUE OF MANY ORTHOGONAL VECTORS

Definition ((q, ϵ)-Restricted Isometry Property)

A matrix \mathbf{A} satisfies (q, ϵ)-RIP if, for all \mathbf{x} with $\|\mathbf{x}\|_0 \leq q$,

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



Every subset of k columns $\mathbf{U} \in \mathbb{R}^{m \times k}$ is approximate isometry.

$$\mathbf{U}^T \mathbf{U} \approx \mathbf{I}.$$

If $\mathbf{U}^T \mathbf{U} \approx \mathbf{I}$, it better be that any two columns $\mathbf{u}_i, \mathbf{u}_j$ are approximately orthogonal.

Deduce: All pairs of columns in \mathbf{A} are approximately orthogonal.

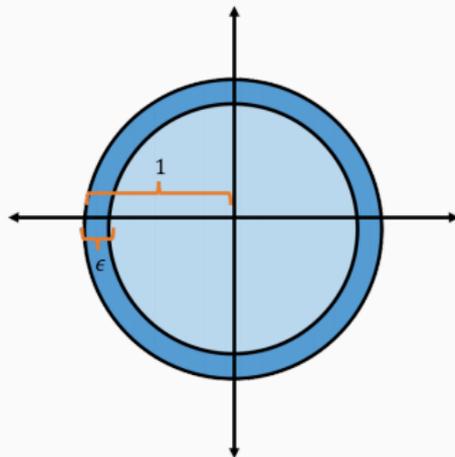
Think of k as a constant: $k = O(1)$. We have d nearly orthogonal vectors living in $O(k \log d) = O(\log d)$ dimensional space.

UNIT BALL IN HIGH DIMENSIONS

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

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

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

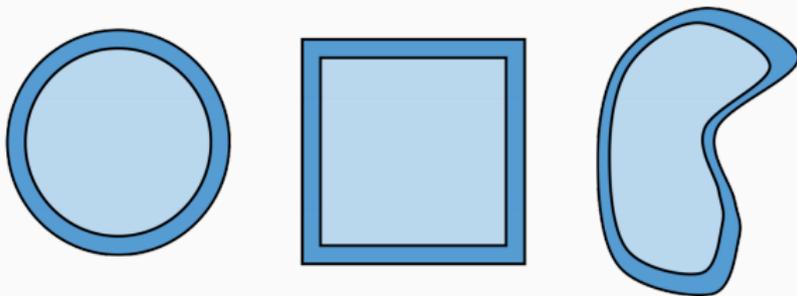


Volume of radius R ball is $\frac{\pi^{d/2}}{(d/2)!} \cdot R^d$.

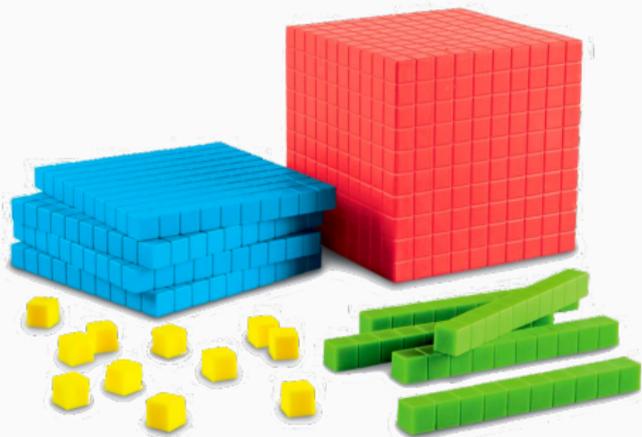
ISOPERIMETRIC INEQUALITY

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

Isoperimetric Inequality:

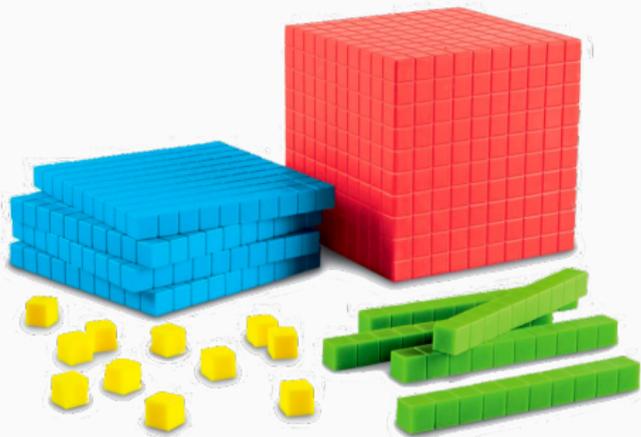


INTUITION



$$\frac{\text{surface cubes}}{\text{total cubes}} =$$

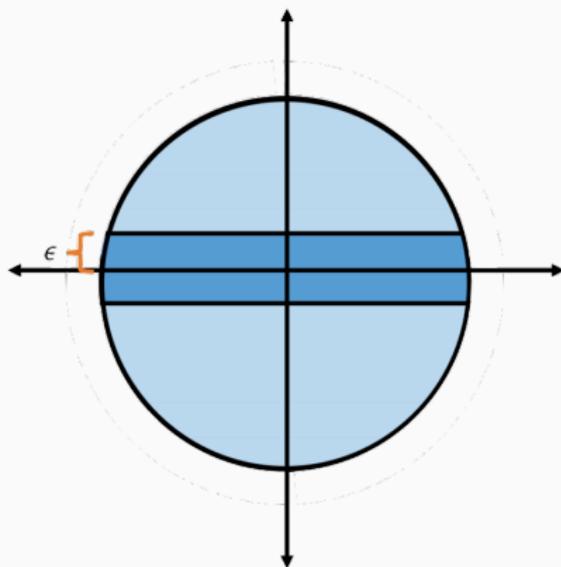
INTUITION



$$\frac{\text{surface cubes}}{\text{total cubes}} =$$

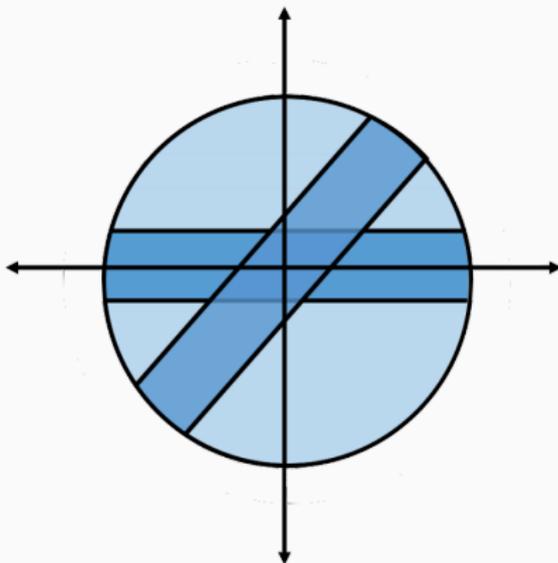
- 1 dimension: $2/10 = .2$
- 2 dimension: $38/100 = .38$
- 3 dimension: $484/1000 = .484$

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



$$S = \{\mathbf{x} \in \mathcal{B}_d : |x(1)| \leq \epsilon\}$$

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.

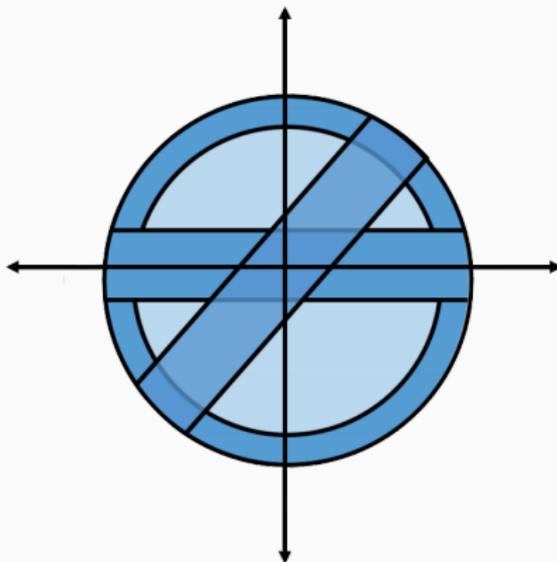


By symmetry, this is true for any equator:

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

BIZARRE SHAPE OF UNIT BALL

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



High-dimensional ball looks nothing like 2D ball!

CONCENTRATION AT EQUATOR

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

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

CONCENTRATION AT EQUATOR

Let $\mathbf{w} = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$.

$$\Pr[|\mathbf{x}(1)| \leq \epsilon] \geq \Pr[|\mathbf{w}(1)| \leq \epsilon].$$

How can we generate \mathbf{w} , which is a random vector taken by scaling a random $\mathbf{x} \in \mathcal{B}_d$?

CONCENTRATION AT EQUATOR

Let \mathbf{g} be a random Gaussian vector – each entry is $\mathcal{N}(0, 1)$. Set $\mathbf{w} = \mathbf{g}/\|\mathbf{g}\|_2$.

- $\mathbb{E}[\|\mathbf{g}\|_2^2] =$

- $\Pr [\|\mathbf{g}\| \leq \frac{1}{2}\mathbb{E}[\|\mathbf{g}\|_2^2]] \leq$

CONCENTRATION AT EQUATOR

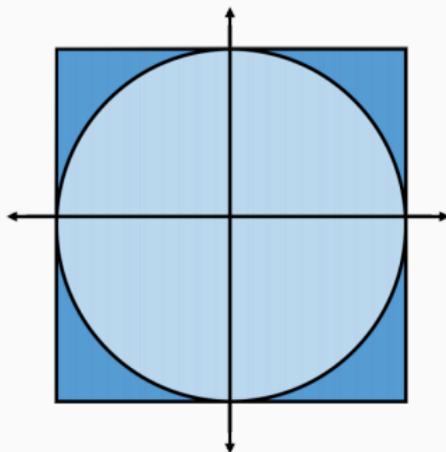
For $1 - 2^{\Theta(d)}$ fraction of vectors \mathbf{g} , $\|\mathbf{g}\|_2 \geq \sqrt{d/2}$. Condition on even that we get a random vector in this set.

$$\begin{aligned}\Pr[|\mathbf{w}(1)| \leq \epsilon] &= \Pr\left[|\mathbf{w}(1)| \cdot \sqrt{d/2} \leq \epsilon \cdot \sqrt{d/2}\right] \\ &\geq \Pr\left[|\mathbf{g}(1)| \leq \epsilon \cdot \sqrt{d/2}\right] \\ &\geq 1 - 2^{\theta\left(-(\epsilon \cdot \sqrt{d/2})^2\right)}\end{aligned}$$

HIGH DIMENSIONAL CUBE

Let \mathcal{C}_d be the d -dimensional cube:

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



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

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

This is a huge gap!

Some other ways to see these shapes are very different:

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

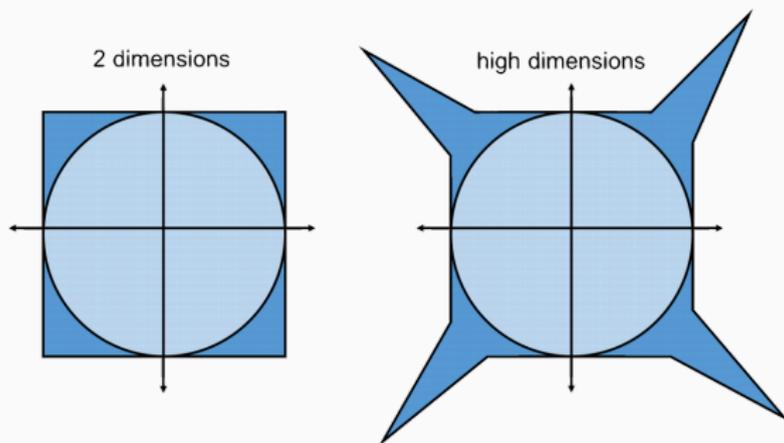
- $\max_{\mathbf{x} \in \mathcal{C}_d} \|\mathbf{x}\|_2^2 =$

Some other ways to see these shapes are very different:

- $\mathbb{E}_{\mathbf{x} \sim \mathcal{B}_d} \|\mathbf{x}\|_2^2$
- $\mathbb{E}_{\mathbf{x} \sim \mathcal{C}_d} \|\mathbf{x}\|_2^2 =$

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.



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.

THANK YOU FOR A GREAT SEMESTER!

Exam format:

- Can have 1 double sided sheet of notes/equations.
- Will be designed for 1.5 hours, but won't cut off until 2 hours.

Exam topics:

- Very end of convex optimization (preconditioning, coordinate descent, + gradient descent for non-convex functions).
- Singular value decomposition and low-rank approximation.
- Spectral graph theory (stochastic block model, matrix perturbation, etc).
- Randomized linear algebra (subspace embeddings, approximate regression, ϵ -nets
- Sparse recovery (restricted isometry property, basis pursuit)

Please write a course review!

- https://m.albert.nyu.edu/app/student/nyuCrseEval/crseEval/1198/24247/Y_LEC/10