

Tuesday, Feb 17

• Exam 3/5

Plan

- Revisiting Load Balancing
- High dimensional geometry

## Exponential Concentration Inequalities

$X_1, \dots, X_n$  indep

$$S = \sum_{i=1}^n X_i \quad \mu = \mathbb{E}[S_i]$$

↳ Make assumptions on  $X_i$ :  
stronger assumption = stronger bound

↳ Proved using clever applications  
of Markov's

## Chernoff Bound

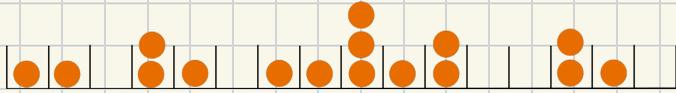
Binary  $X_i \in \{0, 1\}$ . For  $0 < \epsilon < 1$ ,

$$\Pr(|S - \mu| \geq \epsilon \mu) \leq 2 \exp\left(-\frac{\epsilon^2 \mu}{3}\right)$$

For  $\epsilon > 0$ ,

$$\Pr(S \geq (1 + \epsilon)\mu) \leq \exp\left(-\frac{\epsilon^2 \mu}{2 + \epsilon}\right)$$

## Revisiting Load Balancing



$$S_i = \sum_{j=1}^n \mathbb{1}[j \text{ to } i]$$

$$S = \max_i S_i$$

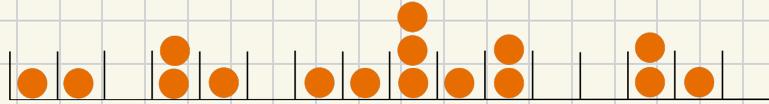
$$\Pr(S \geq c) \leq \frac{1}{10}$$

$$\Leftrightarrow \Pr(S_i \geq c) \leq 1/10n$$

with Chebyshev's,  $c = O(\sqrt{n})$

Q: Can we do better?

## Power of Two Choices



use two hash functions,  
choose least occupied

Then  $\Pr(S \geq c) \leq 1/10$  for  $c = \log n$

$$c = \log \log n$$

$$c = \log \log \log n$$

How about power of three choices?

## High dimensional geometry

Unifying theme...

- random projections
- locality sensitive hashing
- low rank approximation
- graph representations

## Setup

$$x, y \in \mathbb{R}^d$$

$$\langle x, y \rangle = x^T y = y^T x = \sum_{i=1}^d x_i y_i$$

$$\|x\|_2^2 = \langle x, x \rangle = \sum_{i=1}^d x_i^2$$

$$\langle x, y \rangle = \|x\|_2 \|y\|_2 \cos \theta$$

## High-dimensional geometry is weird

Q: In  $d$  dimensions, how many orthogonal vectors are there?

i.e., orthogonal if  $\langle x, y \rangle = 0$

Q: In  $d$  dimensions, how many nearly orthogonal vectors are there?

i.e., nearly orthogonal if  $|\langle x, y \rangle| < \epsilon$

## Probabilistic Method

Let  $t = 2^{ce^2d}$ , for constant  $c$ .

We'll prove  $\exists x_1, \dots, x_t$  nearly orthogonal.

Strategy: Define random process

and show, with non-zero probability,

$|\langle x_i, x_j \rangle| < \epsilon$  for all  $i \neq j$

$$x_i = \begin{bmatrix} +1/\sqrt{d} \\ -1/\sqrt{d} \\ -1/\sqrt{d} \\ \vdots \end{bmatrix}$$

$$x_i[j] = \begin{cases} +1/\sqrt{d} & \text{wp } 1/2 \\ -1/\sqrt{d} & \text{wp } 1/2 \end{cases}$$

$$\|x_i\|_2^2 = \sum_{k=1}^d x_i[k]^2 = \frac{d}{d} = 1$$

$$\begin{aligned} \mathbb{E}[\langle x_i, x_j \rangle] &= \sum_{k=1}^d \mathbb{E}[x_i[k] x_j[k]] \\ &\stackrel{\text{indep}}{=} \sum_{k=1}^d \mathbb{E}[x_i[k]] \mathbb{E}[x_j[k]] \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Var}(\langle x_i, x_j \rangle) &\stackrel{\text{indep}}{=} \sum_{k=1}^d \text{Var}(x_i[k] x_j[k]) \\ &= \sum_{k=1}^d \mathbb{E}[(x_i[k] x_j[k] - 0)^2] \\ &= \sum_{k=1}^d \frac{1}{d^2} = \frac{1}{d} \end{aligned}$$

Thursday, February 19

- High dimensional geometry

## Probabilistic Method

Let  $t = 2^{ce^2 d}$ , for constant  $c$ .

We'll prove  $\exists x_1, \dots, x_t$  nearly orthogonal.

Strategy: Define random process

and show, with non-zero probability,

$|\langle x_i, x_j \rangle| < \epsilon$  for all  $i \neq j$

$$x_i = \begin{bmatrix} +1/\sqrt{d} \\ -1/\sqrt{d} \\ -1/\sqrt{d} \\ \vdots \end{bmatrix}$$

$$x_i[j] = \begin{cases} 1/\sqrt{d} & \text{w.p. } 1/2 \\ -1/\sqrt{d} & \text{w.p. } 1/2 \end{cases}$$

$$\|x_i\|_2^2 = \sum_{k=1}^d x_i[k]^2 = \frac{d}{d} = 1$$

$$\begin{aligned} \mathbb{E}[\langle x_i, x_j \rangle] &= \sum_{k=1}^d \mathbb{E}[x_i[k] x_j[k]] \\ &\stackrel{\text{indep}}{=} \sum_{k=1}^d \mathbb{E}[x_i[k]] \mathbb{E}[x_j[k]] \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Var}(\langle x_i, x_j \rangle) &\stackrel{\text{indep}}{=} \sum_{k=1}^d \text{Var}(x_i[k] x_j[k]) \\ &= \sum_{k=1}^d \mathbb{E}[(x_i[k] x_j[k] - 0)^2] \\ &= \sum_{k=1}^d \frac{1}{d^2} = \frac{1}{d} \end{aligned}$$

Fix  $i, j$ . Let  $Z = \langle x_i, x_j \rangle = \sum_{k=1}^d C_k$

where  $C_k = \begin{cases} 1/d & \text{w.p. } 1/2 \\ -1/d & \text{else} \end{cases}$

Since  $Z$  iid sum, expect  $Z \sim \mathcal{N}$ .

If  $Z$  were Gaussian,

$$\Pr(|Z| \geq \alpha \frac{1}{\sqrt{d}}) = \Pr(|Z - \mathbb{E}Z| \geq \alpha \sigma) \leq O(e^{-\alpha^2})$$

and done by setting  $\alpha = \epsilon \sqrt{d}$ .

But  $Z$  is not Gaussian, so more work  $\ddot{\smile}$

Chernoff? Write  $Z$  as sum of binary.

$$B_k = \begin{cases} 1 & \text{w.p. } 1/2 \\ 0 & \text{else} \end{cases}$$

$$C_k = \frac{2}{d} (B_k - 1/2)$$

$$\begin{aligned} Z &= \sum_{k=1}^d C_k = \frac{2}{d} \sum_{k=1}^d \frac{d}{2} C_k \\ &= \frac{2}{d} \sum_{k=1}^d (-1/2 + B_k) \\ &= \frac{2}{d} \left( -\frac{d}{2} + \sum_{k=1}^d B_k \right) \end{aligned}$$

$$Z > \epsilon \iff$$

$$\sum_{k=1}^d B_k > \frac{d}{2} + \frac{d}{2} \epsilon$$

$$Z < -\epsilon \iff$$

$$\sum_{k=1}^d B_k < \frac{d}{2} - \frac{d}{2} \epsilon$$

$$B = \sum_{k=1}^d B_k, \quad \mathbb{E}[B] = \frac{d}{2}$$

$$\Pr(|B| \geq \epsilon) = \Pr(|B - \mathbb{E}B| \geq \epsilon \mathbb{E}B)$$

$$\leq 2 \exp\left(-\frac{\epsilon^2 \mathbb{E}B}{3}\right)$$

$$= 2 \exp\left(-\frac{\epsilon^2 d}{6}\right)$$

$$\Pr(\exists i \neq j : |\langle x_i, x_j \rangle| \geq \epsilon) \leq \binom{t}{2} 2 \exp\left(-\frac{\epsilon^2 d}{6}\right)$$

Choose  $t$  so failure prob  $< t$ .

Takeaway: Random vectors  
tend to be far apart in high  $d$ .

$\therefore$  working with high  $d$   
might seem hopeless, but

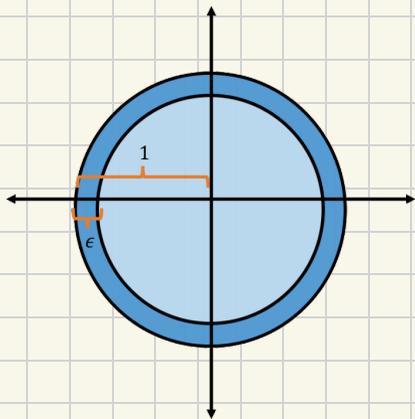
our data is typically not  
random i.e. there is structure

we can use

## High d is weird part 2: Where Random Points Live

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

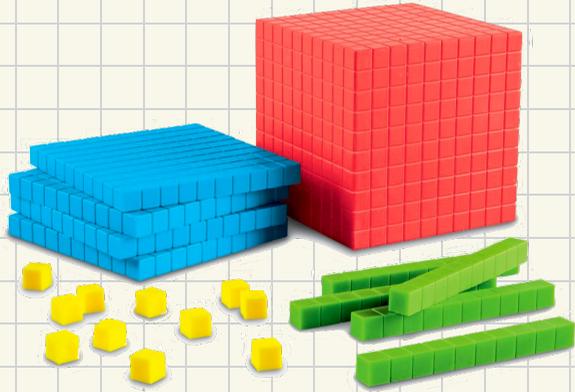
Q: What fraction of volume lies within  $\epsilon$  of surface?



$$\text{Vol}(B_d(R)) = \frac{\pi^{d/2}}{(d/2)!} R^d$$

$$\begin{aligned} \frac{\text{Vol}(B_d(1)) - \text{Vol}(B_d(1-\epsilon))}{\text{Vol}(B_d(1))} &= 1 - (1-\epsilon)^d \\ &= 1 - \left(1 - \epsilon^{1/d}\right)^d \\ &\approx 1 - \frac{1}{e^{\epsilon d}} \end{aligned}$$

All but  $\frac{1}{2e^{\epsilon d}}$  fraction is  $\epsilon$ -close to surface!

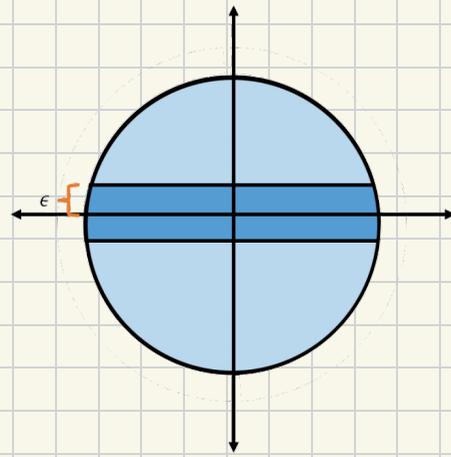


Q: What fraction of cubes  
are near surface?

$d=1$ :

$d=2$ :

$d=3$ :



Q: What fraction of volume  
is  $\epsilon$  close to equator?

A: All but  $\frac{1}{2}$ ced fraction of volume

$\epsilon$ -close to any equator

Draw random points from unit ball

Goal: Show that  $x \sim \beta_d$  has  $|x_i| \leq \epsilon$  w.p.  $1 - \frac{1}{2^c \epsilon^d}$

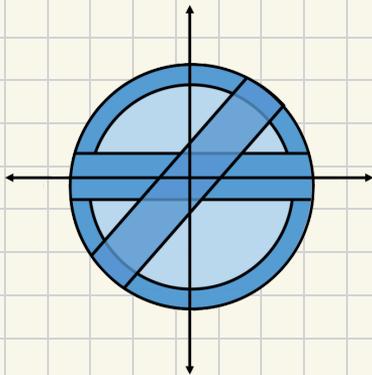
Given  $x$  from interior of unit ball,  $w = \frac{x}{\|x\|_2}$  from surface

If  $|w_i| \leq \epsilon \Rightarrow |x_i| < \epsilon$  since  $\|x\|_2 \leq 1$

New goal: Show that  $w$  from surface has  $|x_i| \leq \epsilon$  w.p.  $1 - \frac{1}{2^c \epsilon^d}$

Let  $g \sim \mathcal{N}(0, I)$ .  $w = \frac{g}{\|g\|_2}$  from surface by rotational invariance

$\mathbb{E} \|g\|_2^2 = \dots$



If:

Then

$$\textcircled{1} \|g\|_2 \geq \sqrt{d/2}$$

$$|w, 1| = \frac{|g, 1|}{\|g\|_2} \leq \frac{\epsilon \sqrt{d/2}}{\sqrt{d/2}} = \epsilon$$

$$\textcircled{2} |g, 1| \leq \epsilon \sqrt{d/2}$$

$$\Pr(\textcircled{1} \text{ or } \textcircled{2}) \leq \Pr(\textcircled{1}) + \Pr(\textcircled{2})$$

$$\begin{aligned} \Pr(\textcircled{1} \text{ and } \textcircled{2}) &= \Pr(\textcircled{1}) + \Pr(\textcircled{2}) - \Pr(\textcircled{1} \text{ or } \textcircled{2}) \\ &\geq \Pr(\textcircled{1}) + \Pr(\textcircled{2}) - 1 \end{aligned}$$

$$\begin{aligned} \Pr(|w, 1| \leq \epsilon) &\geq \Pr(\textcircled{1} \text{ and } \textcircled{2}) \\ &\geq 1 - \Pr(\textcircled{1}^c) - \Pr(\textcircled{2}^c) \end{aligned}$$

$$\begin{aligned} &= (1 - \Pr(\textcircled{1}^c)) + (1 - \Pr(\textcircled{2}^c)) - 1 \\ &= 1 - \Pr(\textcircled{1}^c) + \Pr(\textcircled{2}^c) \end{aligned}$$

$$\Pr(\textcircled{1}^c) = \Pr(\|g\|_2 < \sqrt{d/2}) \leq \frac{1}{2cd} \text{ by Chi-squared concentration}$$

$$\begin{aligned} \Pr(\textcircled{2}^c) &= \Pr(|g, 1| > \epsilon \sqrt{d/2}) \leq \frac{1}{2(c\epsilon \sqrt{d/2})^2} \text{ by Gaussian tail bound} \\ &\geq 1 - \frac{1}{2c^2 \epsilon^2 d/2} - \frac{1}{2cd} \end{aligned}$$

larger for small  $\epsilon$

# Unit Cube vs Sphere

$$C_d = \{x \in \mathbb{R}^d : |x_i| \leq 1 \text{ for all } i \in [d]\}$$

$$\begin{aligned} \frac{\text{Vol}(C_d)}{\text{Vol}(B_d)} &= \frac{2^d (d/2)!}{\pi^{d/2}} \approx \frac{4^{d/2} \sqrt{2\pi}^{d/2} (d/2)^{d/2}}{\pi^{d/2}} \\ &= \sqrt{2\pi}^{d/2} \left(\frac{2d}{\pi e}\right)^{d/2} \approx \sqrt{d}^d \end{aligned}$$

$$\max_{x \in B_d} \|x\|_2^2 = 1$$

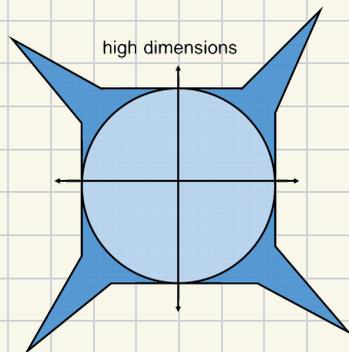
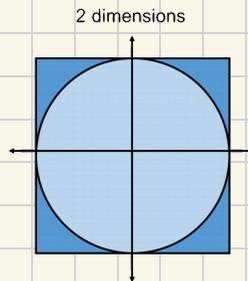
$$\text{vs } \max_{x \in C_d} \|x\|_2^2 = d$$

$$\mathbb{E}_{x \in B_d} [\|x\|_2^2] = 1$$

$$\text{vs } \mathbb{E}_{x \in C_d} [\|x\|_2^2] =$$

Stirling's Approximation

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$$



Next:

Play with high-dimensional data!

Primarily reduce dimension while retaining structure