

Algorithms, Geometry, and Optimization

Lecture 4: Jan 22, 2024

Last class

- Estimation via sampling: *what fraction of a population can ski?* (Sampling gives approx answers).
- Main parameters: estimate, error of estimate, confidence in estimate \downarrow estimated value.
- Suppose population had N people \star no N anywhere in bounds.
- Using $\approx \frac{4}{\epsilon^2}$ samples, we can be 90% confident that true fraction is within $\pm \epsilon$ of the estimate $95\% \dots$

Hoeffding's inequality

Theorem. Suppose X_1, X_2, \dots, X_m are independent random variables with the property that $a_i \leq X_i \leq b_i$ always holds. Let $Y = X_1 + X_2 + \dots + X_m$. Then

$$\Pr[|Y - \mathbb{E}[Y]| > t] \leq 2e^{-\frac{t^2}{\sum_{i=1}^m (a_i - b_i)^2}}$$

- Sum of independent random variables is *well concentrated*
- Closely related to the Central Limit Theorem

Dimension reduction

$$d, n \sim 10^9$$

Suppose we are given points $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, where d is large. Can we *embed* these points into a smaller dimensional space, so distances are approximately preserved?

$$x_i \in \mathbb{R}^d \rightsquigarrow \phi(x_i) \in \mathbb{R}^m \quad (m \ll d, \text{ classical applications of rand. algorithms.})$$
$$\|\phi(x_i) - \phi(x_j)\| \approx \|x_i - x_j\| \quad \forall i, j$$

Johnson-Lindenstrauss (JL) lemma

Theorem (1984). There exists an embedding ϕ into $m = O\left(\frac{\log n}{\epsilon^2}\right)$ dimensional space such that:

for every pair x_i, x_j ,
sq. distance in embedded space
sq. distance in original space.

$$(1 - \epsilon) \|x_i - x_j\|^2 \leq \|\phi(x_i) - \phi(x_j)\|^2 \leq (1 + \epsilon) \|x_i - x_j\|^2.$$

- in fact, linear embeddings do it! $\phi(x) = Ax$

(Euclidean dists are special.)

JL-Lemma

Larsen Nelson, 2017

- The bound is tight (cannot be improved)
- As such, stated and proved for the Euclidean norm only
- Weaker versions hold for ℓ_p norms with $p > 1$
- Impossible for $p = 1$ (Brinkman, Charikar '2002).

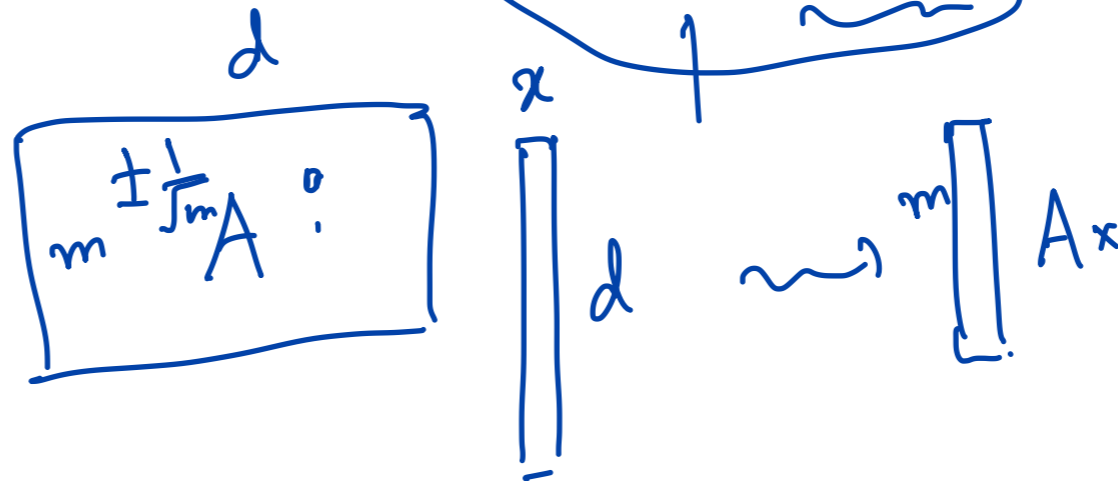
JL-Lemma: Proof Sketch

- Pick A to be an $m \times d$ matrix with each entry A_{ij} uniformly at random!

$$A_{ij} = \begin{cases} \frac{1}{\sqrt{m}} & \text{w.p. } \frac{1}{2} \\ -\frac{1}{\sqrt{m}} & \text{w.p. } \frac{1}{2} \end{cases}$$

- Two choices of distribution: $\mathcal{N}(0, 1/m)$ or $\pm \frac{1}{\sqrt{m}}$

- Set $\phi(x) = Ax$



$$x \rightsquigarrow Ax$$

$$\|Ax\| \approx \|x\|$$

JL-Lemma: Proof Sketch

Lemma. Let x be any vector. Then with probability $\geq 1 - 2e^{-m\epsilon^2/4}$, we have

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

$$m = \frac{8 \log \frac{1}{\epsilon}}{\epsilon^2}$$

$$1 - \frac{1}{4m^2}$$

- This implies the JL lemma. Why?

Actual JL lemma requires the above for $x = x_i - x_j \forall i, j$
 E_{ij} : event that norm preservation holds for pair $(i, j) \rightarrow 1 - \frac{1}{2\binom{n}{2}}$

* Example of Union Bound Argument *

- Can use union bound!

$$\Pr(A_1 \cup A_2 \cup A_3)$$

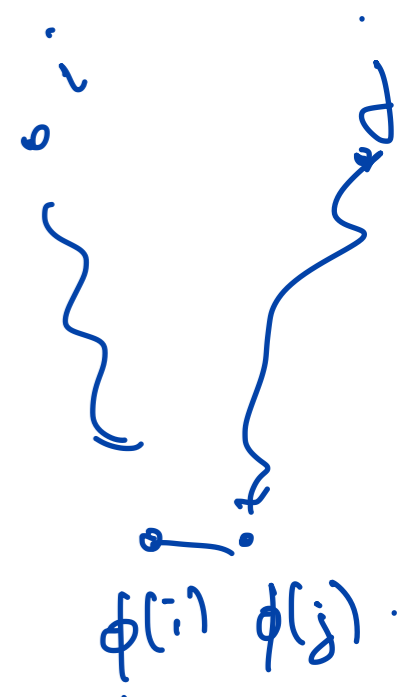
$$\leq \Pr(A_1) + \Pr(A_2) + \Pr(A_3)$$

A_{ij} : event that norm preservation fails to hold for pair (i, j)

$$\Pr(A_{ij}) \text{ for any } i, j \leq \frac{1}{4n^2}$$

$$\Pr\left(\bigcup_{i, j} A_{ij}\right) \leq \binom{n}{2} \cdot \frac{1}{4n^2} \leq \frac{1}{8}$$

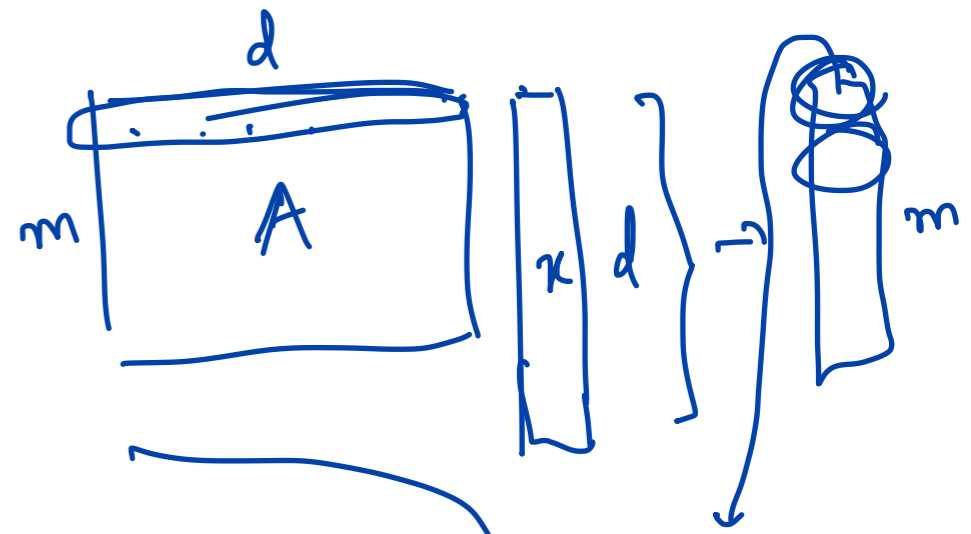
$$\therefore \Pr(\text{all norms are preserved}) \geq \frac{7}{8} \quad (80\%).$$



Expected value

(Lemma wanted to say $\|Ax\|^2 \approx \|x\|^2$)

What is $\mathbb{E}[\|Ax\|^2]$? \rightarrow fixed.
 A_{ij} are random $\pm \frac{1}{\sqrt{m}}$



$$Ax = \begin{bmatrix} \langle A_1, x \rangle \\ \langle A_2, x \rangle \\ \vdots \\ \langle A_m, x \rangle \end{bmatrix}$$

$$\mathbb{E}[\|Ax\|^2] = \mathbb{E}[\langle A_1, x \rangle^2 + \langle A_2, x \rangle^2 + \dots + \langle A_m, x \rangle^2]$$

$\langle A_1, x \rangle$
 \downarrow
 row #1 of A

$$= \mathbb{E}[\langle A_1, x \rangle^2] + \mathbb{E}[\langle A_2, x \rangle^2] + \dots + \mathbb{E}[\langle A_m, x \rangle^2]$$

$$= m \mathbb{E}[\langle A_1, x \rangle^2] = \|x\|^2$$

Want:

$$\mathbb{E} \left[\langle A, x \rangle^2 \right]$$

$$(A_{i1} \ A_{i2} \ \dots \ A_{id})$$

$$\begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_d \end{pmatrix} = x$$

$$= \mathbb{E} \left[(A_{i1} u_1 + A_{i2} u_2 + \dots + A_{id} u_d)^2 \right]$$

$$= \mathbb{E} \left[(A_{i1} u_1)^2 + \dots + (A_{id} u_d)^2 + 2 \sum_{i < j} (A_{i1} u_1)(A_{ij} u_j) \right]$$

$$(x+y+z)^2 = x^2 + y^2 + z^2 + 2xy + 2yz + 2zx$$

$$= \mathbb{E} \left[\frac{1}{m} u_1^2 + \frac{1}{m} u_2^2 + \dots + \frac{1}{m} u_d^2 \right]$$
$$= \frac{1}{m} (u_1^2 + u_2^2 + \dots + u_d^2) = \frac{\|x\|^2}{m}$$

$$\mathbb{E} \left[A_{i1} A_{ij} u_1 u_j \right]$$

\downarrow \downarrow
 $\frac{1}{\sqrt{m}}$ $\frac{1}{\sqrt{m}}$

★★

$$0 \leq X_i \leq 1$$

(Hoeffding...)

Concentration bound

- Define $Y_i = \langle A_i, x \rangle^2$
- What is its range?

Is this concentrated?

$$\mathbb{E} \left[\langle A_1, x \rangle^2 + \langle A_2, x \rangle^2 + \dots + \langle A_m, x \rangle^2 \right]$$
$$= \|x\|^2$$

$a_i \leq Y_i \leq b_i$ (too bad...)

~~$0 \leq Y_i \leq \frac{\|x\|^2}{m}$~~

pretend \rightarrow $\frac{2\|x\|^2}{m}$

$\frac{\|x\|^2}{m}$ in expectation.

Turns out to imply the lemma.

Example: linearity of expectation

- Max 3-SAT problem

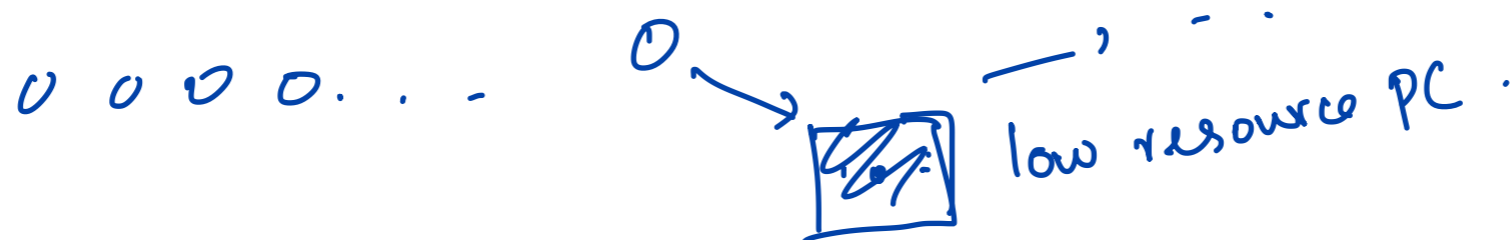
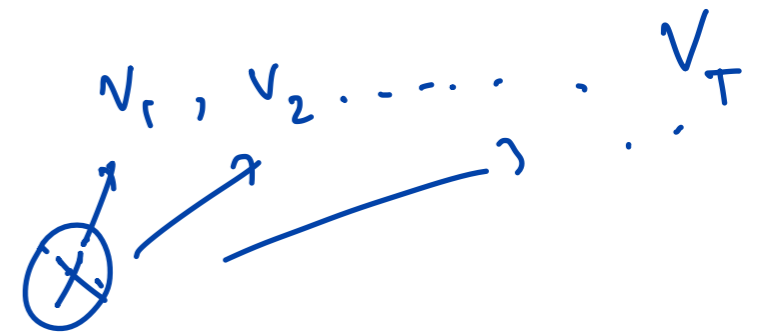
boolean variables x_1, \dots, x_n

$$\underbrace{(x_1 \vee x_2 \vee \bar{x}_3)}_{C_1}, \underbrace{(x_3 \vee \bar{x}_5 \vee x_7)}_{C_2}, \dots$$

Surprising fact: if you set each x_i to ~~a~~ true/false uniformly at random, then $\frac{7}{8}$ fraction of the clauses get satisfied in expectation.

Streaming algorithms: basic model

- Data: values (phone numbers, IP addresses, etc.) arrives one after another
- Full data set is too hard to store
- Router must compute aggregate statistics



Distinct elements problem

- How many distinct values does the stream have?



Suppose answer is k .

- Hash set or equiv: $O(k)$ space, $\log k$ lookup element, processing time.
- Bad if k is very large... $\log k$ space!!

