

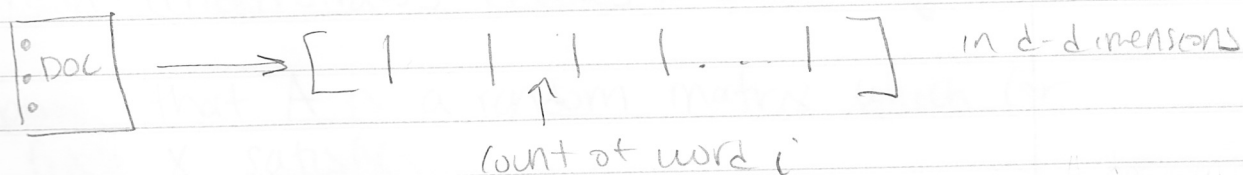
6.854 / 18.415

Last Time: distinct elements and heavy hitters

ie. approximating statistics of a stream w/o storing it

Today: dimensionality reduction

Many types of data are high-dimensional, e.g. documents as bag-of-words



$$\text{Let } \|u-v\|_2^2 = \sum_{i=1}^d (u_i - v_i)^2 \text{ called } \ell_2^2\text{-distance}$$

All-Pairs Distances: Given n documents, compute all pairs of distances (e.g. for hierarchical clustering)
kruskal

Time: $O(n^2 d)$

Can we find representations in lower dimensions that approximately preserve the geometry?

Theorem: [Johnson, Lindenstrauss] For any $\epsilon > 0$ and n points in d -dimensions there is a matrix $A \in \mathbb{R}^{k \times d}$ with $k = \frac{20 \log n}{\epsilon^2}$ such that

$$(1-\epsilon) \|u-v\|_2^2 \leq \|Au - Av\|_2^2 \leq (1+\epsilon) \|u-v\|_2^2$$

for all $u, v \in S$, with $k = \frac{20 \log n}{\epsilon^2}$

JL-property

why does n show up and not d ?

Time via JL: $O\left(\underbrace{\frac{dn \log n}{\epsilon^2}}_{dk \cdot n} + \frac{n^2 \log n}{\epsilon^2}\right)$
 matrix-vec multiply

Moreover we can choose A randomly with ^{independent} Gaussian entries (not only does it exist, but almost every choice works)

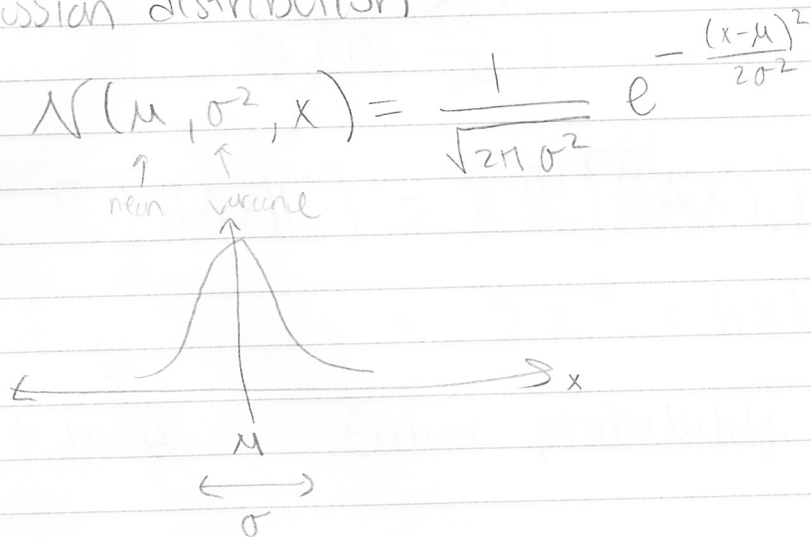
Johnson-Lindenstrauss reduces to preserving norms

Suppose that A is a random matrix which for any fixed x satisfies

$$(*) \quad \Pr\left[(1-\epsilon)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1+\epsilon)\|x\|_2^2\right] \geq 1 - \frac{1}{n^3}$$

then by union bounding over $x = u - v$, for all $u, v \in S$ we get the theorem (since there's $\binom{n}{2}$ of them)

The Gaussian distribution



Ubiquitous in probability and statistics

normalized appropriately

Central Limit Theorem: Sums of independent, random variables converge to a Gaussian

The Chernoff bound shows such sums concentrate, but another explanation is they converge to Gaussian

Fact: If $Z_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $Z_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$

then $Z_1 + Z_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

(think of Z_1, Z_2 as representing sums of indep. r.v.s.)

Now returning to JL, choose entries of A to be $\mathcal{N}(0, \frac{1}{k})$

Then

$$\left((Ax)_i \right)^2 = \left(\sum_{i=1}^d A_{i,j} x_i \right)^2$$

$$\mathcal{N}\left(0, \frac{\sum x_i^2}{k}\right)$$

using fact

$$\begin{aligned} \text{Hence } \mathbb{E}[\|Ax\|_2^2] &= k \mathbb{E}[\left((Ax)_i \right)^2] \\ &= \sum x_i^2 = \|x\|_2^2 \end{aligned}$$

Now to bound the failure probability:

one side:

$$\Pr \left[\|Ax\|_2^2 > (1+\epsilon) \|x\|_2^2 \right]$$
$$= \Pr \left[\sum_{i=1}^k z_i^2 > (1+\epsilon) \|x\|_2^2 \right]$$

where $z_i \sim \mathcal{N}\left(0, \frac{\|x\|_2^2}{k}\right)$, now dividing
thru by $\frac{\|x\|_2^2}{k}$, we get

$$= \Pr \left[\sum_{i=1}^k Y_i^2 > (1+\epsilon)k \right]$$

where $Y_i \sim \mathcal{N}(0, 1)$, and similarly for other side

Lemma:

$$\Pr \left[\sum_{i=1}^k Y_i^2 > (1+\epsilon)k \right] \leq e^{-\frac{k}{4}(\epsilon^2 - \epsilon^3)}$$

$$\Pr \left[\sum_{i=1}^k Y_i^2 < (1-\epsilon)k \right] \leq e^{-\frac{k}{4}(\epsilon^2 - \epsilon^3)}$$

The proof is similar to usual proof of Chernoff
via mgf (see link)

Now if we plug in $k = \frac{20 \log n}{\epsilon^2}$, we can

invoke the lemma to show that w.h.p. $\|Ax\|_2^2 \sim \|x\|_2^2$

then union bound over $x = u - v$

Extensions:

What really happened here is for i.i.d mean-zero, variance one z_1, z_2, \dots, z_d the quantity

$$\left(\sum_{i=1}^d z_i x_i \right)^2$$

is an unbiased estimator for $\|x\|_2^2$, and we boosted it by repeating

Theorem [Achlioptas] ²⁰⁰³ If the entries of A are chosen as

$$\text{for } 1 \leq i \leq k, 1 \leq j \leq d \quad R_{ij} = \sqrt{3} \begin{cases} +1 & \text{with probability } 1/6 \\ 0 & \text{" } 2/3 \\ -1 & \text{" } 1/6 \end{cases}$$

and $A = \frac{1}{\sqrt{k}} R$, then A has the JL-property

multiplying by A is easier in many settings

This idea of finding sparse dimension reduction can be taken much further

Theorem: ²⁰¹⁴ [Kane, Nelson] ^{non i.i.d} there are distributions on $A \in \mathbb{R}^{k \times d}$ with $k = O\left(\frac{\log n}{\epsilon^2}\right)$ and $s = O\left(\frac{\log n}{\epsilon}\right)$ non-zeros per column that have the JL-property

Multiplying sparse vectors by sparse matrices is fast, enables many applications

Fast Johnson-Lindenstrauss:

An alternative approach due to Ailon, Chazelle ²⁰⁰⁶ is

$$A = PHD$$

↑ ↑ ↑
sparse hadamard diagonal

where $H_1 = [1]$ 1x1 hadamard

$$H_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad 2 \times 2 \text{ hadamard}$$

$$H_{2^r} = \begin{bmatrix} H_{2^{r-1}} & H_{2^{r-1}} \\ H_{2^{r-1}} & -H_{2^{r-1}} \end{bmatrix} \quad 2^r \times 2^r \text{ hadamard}$$

and $D = \begin{bmatrix} \pm 1 & & & \\ & \pm 1 & & \\ & & \ddots & \\ & & & \pm 1 \end{bmatrix}$

esp. in the
scattering

Idea: Multiplying a sparse vector by a sparse matrix can result in a poor embedding, so make the vectors dense

FJLT relies on fact that multiplying by H can be done by divide-and-conquer, and uncertainty principles

distribution on

Theorem [Ailon - Hazelle] There is a matrixes
 $A \in \mathbb{R}^{k \times d}$ with $k = O\left(\frac{\log n}{\epsilon^2}\right)$ ⁽¹⁾ that has the
JL-property, ⁽²⁾ and matrix-vector multiplications
can be done in time

$$O(d \log d + \frac{d \log n}{\epsilon^2})$$

provided $n < 2^{d^{1/3}}$

Applications to nearest neighbor (we'll cover this
problem later)