

Lecture 9: Dimensionality Reduction

[[We live in an era of "Big Data".

Not unusual for an "object" (image, person, video,...) to have millions, billions of numerical "features", and/or to have millions, billions of items. Doing one thing a billion times is maybe okay, algorithmically. Thousands of such things? Not so. Billion x billion? Don't make me laugh.]]

Say you have n data items, $x_1^{(i)}, \dots, x_n^{(i)}$ $n = 2^{20 \dots 30}$ lots!]]
Each has D numerical features. [pixel colors, song preferences...]
 \hookrightarrow an item $x_i \cong$ vector / point in \mathbb{R}^D . [Idea: nD is very large; too large.]

~~More~~ "Big Data" / "High-Dim. Data" ["Curse of dimensionality"]

[Any number of "geometric" algs. you might want to run on data...]
• clustering • nearest neighbor queries • PCA [next time]

Idea: (try to) map data to low dimension,
 $x_i^{(i)} \mapsto y_i^{(i)} \in \mathbb{R}^K$, $K \ll D$

so that ... relationships... preserved

\hookrightarrow distances
• lengths
• angles
• dot products } [these are all kinda the same...]

Recall of inner/dot product
of $u, v \in \mathbb{R}^m$

$$\langle u, v \rangle = u \cdot v = \sum_{i=1}^m u_i v_i$$

$$u \cdot u = \sum_i u_i^2 = \|u\|^2$$

Fact ("cosine law"): $u \cdot v = \|u\| \cdot \|v\| \cdot \cos \theta$
 θ angle between u, v .

$$\text{dist}(u, v)^2 = \|u - v\|^2 = (u - v) \cdot (u - v) = u \cdot u - 2u \cdot v + v \cdot v$$

[So: all interrelated.]

② Potential objection: seems impossible to preserve all this in lower dimensions.

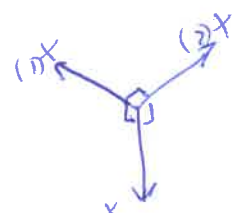
Worry: Say $\|x^{(i)}\|^2 = 1$ $A: f_j$ (all data vectors are "unit vectors")

And $\|x^{(i)} - x^{(j)}\| = \sqrt{2}$

$$\Rightarrow \|x^{(i)} - x^{(j)}\|^2 = 2 = x^{(i)} \cdot x^{(i)} - 2x^{(i)} \cdot x^{(j)} + x^{(j)} \cdot x^{(j)} = 2$$

$\Rightarrow x^{(i)} \cdot x^{(j)} = 0$

$\Rightarrow x^{(i)}, x^{(j)} = 0$ \Rightarrow perpendicular (angle = $\frac{\pi}{2} = 90^\circ$)



∴ (high-dim. analogue of)

Basic linear algebra \Rightarrow n perpendicular vecs need to be in n dims. \therefore cannot achieve $K < n$.

But. Say you don't mind "10% error".

$\rightarrow y^{(i)}$'s with $\|y^{(i)}\|^2 \leq 0.1$ $A: f_j$

$\|y^{(i)}\|^2 \leq \|y^{(i)} - y^{(j)}\|^2 \leq \|y^{(j)}\|^2 \leq 0.1$ $A: f_j$

Possible to have n such $y^{(i)}$'s in $\dots \approx \log_{10} n$ dims!

Exponentially fewer dims!

$\|K \approx 3000$ for $n = 1 \text{ Billion}\|$

$\|$ processing or data item is now $300,000 \times$ faster! $\|$

$\|$ we'll show it's possible for any $x^{(i)}$'s but let's first see why for these special $x^{(i)}$'s How to find

exponentially many "almost-orthogonal" vectors in K dims?

$\|$ When you don't know how to find something, pick it at random! $\|$

"random ± 1 bits"

Idea: Let each $y^{(i)} = \begin{bmatrix} +1 \\ +1 \\ \vdots \\ +1 \\ -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$

En. then $\|y^{(i)}\|^2 = \sum_{k=1}^K (\pm 1)^2 = K$, supposed to be ≈ 1 .

So actually, $y^{(i)} = \frac{1}{\sqrt{K}} \begin{bmatrix} +1 \\ +1 \\ \vdots \\ +1 \\ -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$, random y_i for $i = 1, 2, 3, \dots$

Q: What is $y^{(1)} \cdot y^{(2)}$?
 (for any $y^{(1)}, y^{(2)}$)
 It's like $\frac{1}{\sqrt{K}} \left(\frac{1}{\sqrt{K}} + \frac{1}{\sqrt{K}} + \dots + \frac{1}{\sqrt{K}} - \frac{1}{\sqrt{K}} - \frac{1}{\sqrt{K}} - \dots - \frac{1}{\sqrt{K}} \right)$
 all these bits random

"distributed as" $\sim \frac{1}{K} \left((-1) + (-1) + \dots + (-1) + (+1) + \dots + (+1) \right)$
 = avg. of K random ± 1 bits

In expectation, $E[y^{(1)} \cdot y^{(2)}] = 0$ [exactly]

But we have to worry about fluctuations. What be this exactly? Is it, like 80% of the time at most ± 1 , or what?

Let B_1, \dots, B_K be random ± 1 . Let $A = \text{avg}(B_1, \dots, B_K)$

std dev $[A] = ?$ // key quantity to understand

fact: $\text{Var}[cX] = c^2 \text{Var}[X]$
 $\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$ if independent // as B_i 's are

$$\text{Var}[A] = \text{Var}\left[\frac{1}{K}(B_1 + \dots + B_K)\right]$$

$$= \frac{1}{K^2} (\text{Var}[B_1] + \dots + \text{Var}[B_K])$$

$$\text{Var}[B_i] = E[B_i^2] - E[B_i]^2 = 1 - 0^2 = 1$$

$$\therefore \text{Var}[A] = \frac{1}{K} \cdot \frac{1}{K} = \frac{1}{K^2}$$

intuitively/heuristically we expect A is "mean \pm a few std dev"

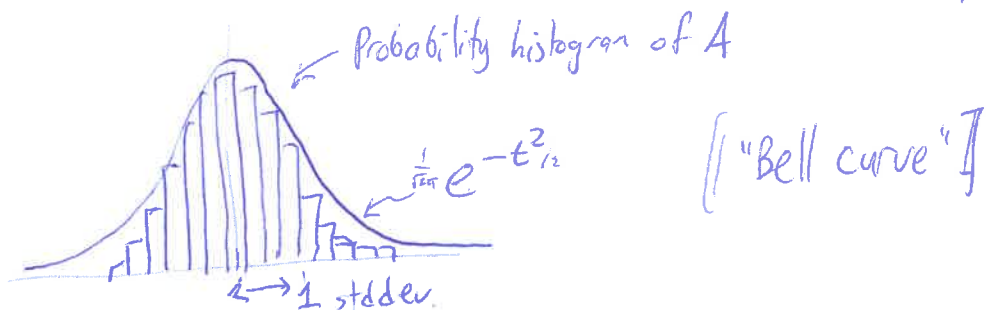
~~distributed~~ A distributed like dot prod. of $y^{(1)}, y^{(2)}$, want it ≈ 0.1

So want $\frac{1}{\sqrt{K}} \approx \epsilon \Leftrightarrow K \approx \frac{1}{\epsilon^2}$ // that was 100, for $\epsilon = 10\%$

OK, but that's just one dot product $y^{(1)} \cdot y^{(2)}$. Want it $\approx \frac{1}{\sqrt{K}}$ for all (i, j)

Q: What is $\Pr[A \geq 20 \cdot \frac{1}{\sqrt{k}}]$? (I just made up "20".) (4)

Heuristic fact: Linear combo of random bits, (indeed, any "nice" r.v.'s)
 $c_1 B_1 + \dots + c_k B_k$,
 acts like a Gaussian random var. (of same mean, stdev.)



$$\Rightarrow \Pr[|A| \geq t \cdot \text{stdev}] \leq e^{-\Theta(t^2)}$$

\therefore for $t = \text{Const} \cdot \sqrt{\log n}$, can make expon. like $-3 \ln(n)$.

$$\Rightarrow \Pr[|A| \geq C \cdot \sqrt{\log n} \cdot \frac{1}{\sqrt{k}}] \leq e^{-3 \ln(n)} = \frac{1}{n^3}$$

\therefore any particular pair $y^{(i)}, y^{(j)}$ has Prob $\leq \frac{1}{n^3}$ of having
 $|\text{dot prod}| \geq C \cdot \sqrt{\log n} \cdot \frac{1}{\sqrt{k}}$

There are $\binom{n}{2} \leq n^2$ pairs \therefore By Union Bound,

$$\Pr[\text{any pair has } |\text{dot prod}| \geq C \cdot \sqrt{\log n} \cdot \frac{1}{\sqrt{k}}] \leq n^2 \cdot \frac{1}{n^3} \leq \frac{1}{n} \quad \left[\begin{array}{l} \text{like} \\ 1 \text{ in a} \\ \text{billion} \end{array} \right]$$

\therefore "w.h.p.", all pairs have $|\text{dot prod}| \leq O(\sqrt{\log n} \cdot \frac{1}{\sqrt{k}})$ (with high probability) [solving] $= \epsilon$ if $k = O(\frac{\log n}{\epsilon^2})$

[OK, but that was for a specific set of data items $x^{(1)} \dots x^{(n)}$, an orthonormal basis.]

In general, say $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^D$, let $\epsilon > 0$ be given.

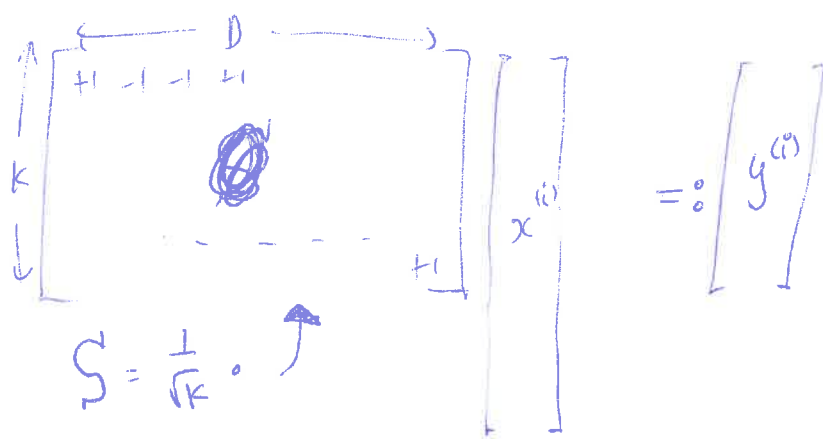
$$\text{Set } k = C \cdot \frac{\log n}{\epsilon^2}$$

(a moderately large constant)

Define:

$$S = \frac{1}{\sqrt{k}} \cdot (k \times D \text{ matrix of random } \pm 1\text{'s})$$

"random projection"



3

Theorem ("J.L. Lemma"): With high prob, (failure prob. $\ll \frac{1}{n}$),

Johnson-Lindenstrauss

$$(1-\epsilon) \|x^{(i)}\| \leq \|y^{(i)}\| \leq (1+\epsilon) \|x^{(i)}\| \quad \textcircled{1} \quad (\text{"lengths preserved"})$$

$$(1-\epsilon) \|x^{(i)} - x^{(j)}\| \leq \|y^{(i)} - y^{(j)}\| \leq (1+\epsilon) \|x^{(i)} - x^{(j)}\| \quad \textcircled{2} \quad (\text{"distances preserved"})$$

Proof sketch: Suffices to do ①. Why? Say we know ① true.

~~Length preserved~~ Given x 's, define $z^{(ij)} := x^{(i)} - x^{(j)}$

Like $\binom{n}{2}$ new data points

$$\textcircled{1} \Rightarrow \|S z^{(ij)}\| \approx \|z^{(ij)}\|$$

$$\Rightarrow \|S(x^{(i)} - x^{(j)})\| \approx \|x^{(i)} - x^{(j)}\| = \|Sx^{(i)} - Sx^{(j)}\| \Rightarrow \textcircled{2}$$

OK, "n" is now $\binom{n}{2}$, but $\log\left(\binom{n}{2}\right) = O(\log n)$

$$\textcircled{1} \Rightarrow (1-\epsilon)^2 \|x^{(i)}\|^2 \leq \|Sx^{(i)}\|^2 \leq (1+\epsilon)^2 \|x^{(i)}\|^2$$

$$(1 \pm \epsilon)^2 \approx 1 \pm 2\epsilon$$

$$e^{\pm 2\epsilon} \approx e^{\pm 2\epsilon} \quad \left\{ \begin{array}{l} \text{again} \\ \rightarrow O(\cdot) \end{array} \right.$$

Drop (i) for notational simplicity.

$$\Leftrightarrow (1-2\epsilon) \sum x_j^2 \leq Sx \cdot Sx \leq (1+2\epsilon) \sum x_j^2$$

$$Sx = \frac{1}{\sqrt{k}} \begin{bmatrix} B_{11}x_1 + B_{12}x_2 + \dots + B_{1D}x_D \\ B_{21}x_1 + \dots + B_{2D}x_D \\ \vdots \\ B_{k1}x_1 + \dots + B_{kD}x_D \end{bmatrix}$$

Dot with itself...

$$S = \frac{1}{\sqrt{k}} \begin{bmatrix} B_{11} & \dots & B_{1D} \\ \vdots & & \vdots \\ B_{k1} & \dots & B_{kD} \end{bmatrix}$$

$$Sx \cdot Sx = \frac{1}{K} \left(B_{11}x_1 + \dots + B_{1D}x_D \right)^2 + \frac{1}{K} \left(B_{21}x_1 + \dots + B_{2D}x_D \right)^2 + \dots$$

(K similar, independent terms)

looks like $(+x_1 - x_2 - x_3 + x_4 - \dots - x_D)^2$

$$= x_1^2 + x_2^2 + \dots + x_D^2 + \text{cross-terms}$$

$$= \|x\|^2 + \text{cross-terms}$$

$\pm x_1 x_2 \pm x_1 x_3 \pm \dots$
random bits

$$\therefore Sx \cdot Sx = \frac{1}{K} \|x\|^2 + \frac{1}{K} (\text{D sets of indep. cross-terms})$$

$$\therefore \mathbb{E}[Sx \cdot Sx] = \|x\|^2 + 0 \quad \text{②}$$

exercise: ~~all~~ [all cross terms] $\leq \frac{3}{\sqrt{K}} \cdot \|x\|^2$

\therefore heuristically, $Sx \cdot Sx \sim \|x\|^2 \pm \frac{\text{const.}}{\sqrt{K}} \cdot \|x\|^2$
 $\uparrow \epsilon$ if $K = O(1/\epsilon^2)$

$$\text{Pr}[Sx \cdot Sx \text{ not in } (1 \pm \epsilon)\|x\|^2] \leq \frac{1}{N^5} \therefore$$

Small Enough for union bound.

□

Painful math
 [but similar bell curve heuristic]

Remarks: Storing S , computing Sx is somewhat costly.

Known: S can be sparse: only ϵ frac. of entries ± 1 , rest 0.
 OK if S not random, just comes from "random enough" hash function

Alternate method [Ailon-Cuzzelle] using

Fast Fourier Transform:

Can compute Sx in $\approx O(D \log D) + O(K \log D)$ time instead of $O(KD)$.