

Lecture 11: Principal Component Analysis (PCA)

©

Last time: Say A is symmetric, $A = A^T$, $D \times D$ matrix.

Then there ~~are~~ is an orthonormal basis

$\vec{v}_1, \dots, \vec{v}_D$ ("eigenvectors") and numbers

$\lambda_1, \dots, \lambda_D$ ("eigenvalues") such that

$A \equiv$ "stretch by factor λ_j in dir. \vec{v}_j "

Saw alg. running in time $O\left(\frac{\log D}{\log \frac{\lambda_{\max}}{\lambda_{\min}}}\right)$ to find

[approximations of] $\lambda_{\max}, \lambda_{\min}$.

Assumed λ_j 's at ≥ 0 . If not, put 1-1 signs.

[Can iterate to find and max, 3rd max, etc.]

Worry: If $\lambda_{\max} = (1+\epsilon)\lambda_{\min}$, denom. is $\log(1+\epsilon) = O(\epsilon)$.

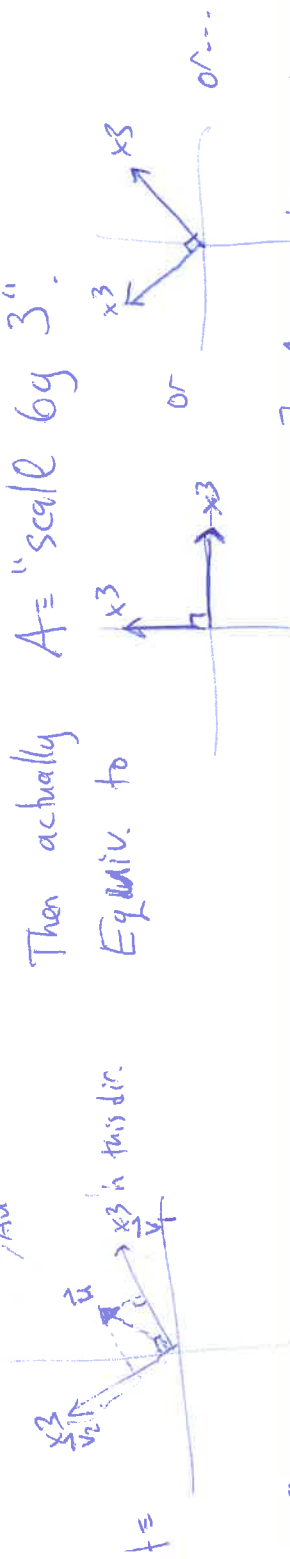
If $\lambda_{\max} = \lambda_{\min}$, denom. is 0!

[I'm here to tell you it's not really a problem.]

Suppose $\lambda_{\max} = \lambda_{\min} = 3$. Indeed, say just $D=2$.

Then actually $A \equiv$ "scale by 3".

Eq. equiv. to



["eigenvectors" \vec{v}_1, \vec{v}_2 not ~~are~~ uniquely defined.] Any two orthonormal

vectors in this 2-D space are valid "eigenvalues of stretch 3,3".

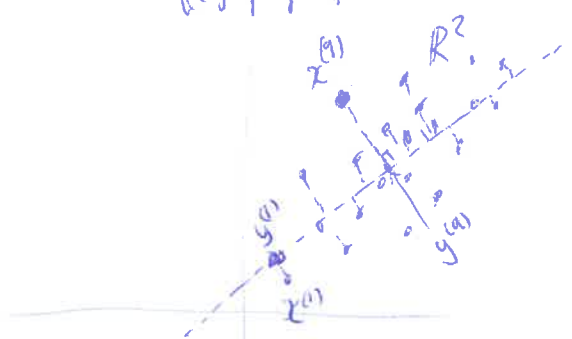
Our alg. will just find a ~~random~~ random pair, and that's okay!

[efficiently, if $\frac{\lambda_{\max}}{\lambda_{\min}}$ large]

On to PCA... Back to "big data" scenario.

Say we have "items" $x^{(1)}, \dots, x^{(n)}$, each with D numerical features, (e.g. people) so $x^{(j)} \in \mathbb{R}^D$.

e.g. $D=2$
(hard to draw $D=3$)



$[x^{(j)}]$

PCA people like the items as row vectors, which annoys me, but oh well.

If you wanted to map these points down to $k=1$ dimension, probably dashed line is best...

Find the "best fitting" line ($k=1$), etc. or 2D-plane ($k=2$) or 3D-plane ($k=3$). Focus on $k=1$ (line) for now.

Goal: minimize $\sum_{j=1}^n (err^{(j)})^2$, where $err^{(j)} = \|x^{(j)} - y^{(j)}\|$

Why square the distance? Well, it makes the math nice.

projection of $x^{(j)}$ onto k -dim space

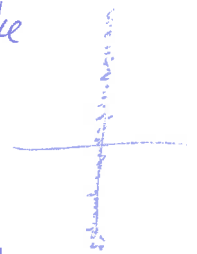
Preprocessing

Linear algebra doesn't like lines/planes not thru origin. So ~~translate~~ translate your data so centered at origin (which turns out to make best-fitting line thru origin too)

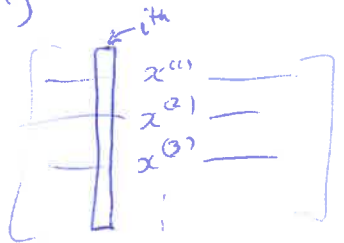
- compute avg. of all $x^{(j)}$'s (rows), then subtract it from each ["center of mass"] [no big deal...]

Now we can assume $avg(x^{(j)}) = [0 \dots 0]$.

Imagine feature 1 is a rating from -5 to 5, and feature 2 is a rating from -10000 to 10000. Data will look like and best-fitting line will be vertical, for a dumb reason.



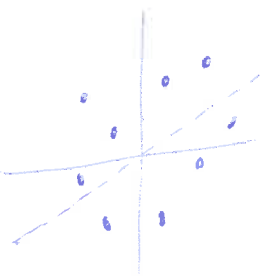
- Scale i^{th} column



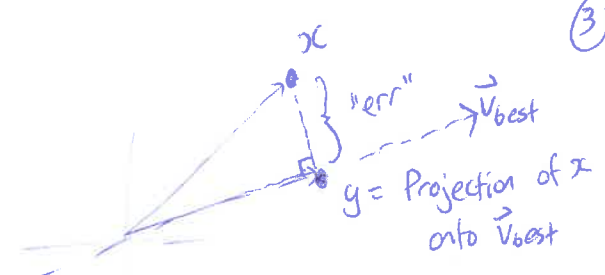
so it's a unit-length vector. Kinda just changing the "units" feature i is measured in.

[[Okay, enough warmup, let's get to it...]]

(3)



[[Zoom in on one point...]]



notes: For any x , $err^2 = \|x\|^2 - \|y\|^2$ (Pythagoras)
 \therefore maxing err^2 small \equiv making $\|y\|^2$ big [$\|x\|^2$ is fixed]

\therefore "best-fitting line" is ~~the~~ equivalently the one in direction \vec{v}_{best} which maximizes $\sum_{j=1}^n \|\text{Proj}_{\vec{v}_{best}}(x^{(j)})\|^2$,

[[also an interesting interpretation: "maximizing variance"]]

$$\|\text{Proj}_{\vec{v}}(\vec{x})\| = \underbrace{|\vec{v} \cdot \vec{x}|}_{\text{unit.}}$$

$$\therefore \|\text{Proj}_{\vec{v}}(\vec{x})\|^2 = (\vec{v} \cdot \vec{x})^2$$

$$\vec{v} \cdot \vec{x} = \|\vec{v}\| \cdot \|\vec{x}\| \cdot \cos\theta$$

$$\text{proj} = \|\vec{x}\| \cos\theta$$

$$\vec{v} \cdot \vec{x} = \begin{bmatrix} -\vec{v}- \\ \hline \hline \end{bmatrix} \begin{bmatrix} | \\ \vec{x}^T \\ | \\ \hline \hline \end{bmatrix}$$

$$= \begin{bmatrix} -\vec{x}- \\ \hline \hline \end{bmatrix} \begin{bmatrix} | \\ \vec{v}^T \\ | \\ \hline \hline \end{bmatrix}$$

$$\therefore (\vec{v} \cdot \vec{x})^2 = \begin{bmatrix} -\vec{v}- \\ \hline \hline \end{bmatrix} \begin{bmatrix} | \\ \vec{x}^T \\ | \\ \hline \hline \end{bmatrix} \begin{bmatrix} -\vec{x}- \\ \hline \hline \end{bmatrix} \begin{bmatrix} | \\ \vec{v}^T \\ | \\ \hline \hline \end{bmatrix}$$

So we want unit ~~vector~~ \vec{v}_{best} maximizing

$$\sum_{j=1}^n \vec{v} x^{(j)T} x^{(j)} \vec{v}^T = \vec{v} \left(\sum_{j=1}^n x^{(j)T} x^{(j)} \right) \vec{v}^T$$

$$\begin{bmatrix} | \\ x \\ | \\ \hline \hline \end{bmatrix} \begin{bmatrix} -x- \\ \hline \hline \end{bmatrix}$$

DxD matrix where

$$A = \sum_{j=1}^n x^{(j)T} x^{(j)}$$

$$= \vec{v} A \vec{v}^T$$

Exercise: If $X = \begin{bmatrix} - & x^{(1)} & - \\ - & x^{(2)} & - \\ & \vdots & \\ & & - \end{bmatrix}$, then $A = X^T X$. (4)

[It's literally just by the defⁿ of how matrix mult. works.]

$$\begin{bmatrix} | & | \\ x^{(1)} & x^{(2)} \\ | & | \end{bmatrix} \begin{bmatrix} - & x^{(1)} & - \\ - & x^{(2)} & - \\ & \vdots & \\ & & - \end{bmatrix}$$

Ta-da, A is a symmetric matrix 😊. $A^T = (X^T X)^T = X^T X^{TT} = X^T X = A$.

[We're so happy, we understand symmetric matrices!]

~~Which vector v makes $v^T A v$ biggest?~~

We know $A \equiv$ stretch by $\lambda_1, \dots, \lambda_D$
in dir. v_1, \dots, v_D (unit basis)

$$\Rightarrow A v_i = \lambda_i v_i \Rightarrow v_i \cdot (A v_i) = v_i^T A v_i = v_i \cdot (\lambda_i v_i) = \lambda_i, \text{ since } v_i \cdot v_i = 1$$

$$\begin{aligned} v_i^T X^T X v_i &= (X v_i)^T (X v_i) \\ &= (X v_i) \cdot (X v_i) = \|X v_i\|^2 \\ \therefore \lambda_i &= \|X v_i\|^2 \geq 0 \\ &\text{(all equals } \geq 0, \text{ extra } \text{😊)} \end{aligned}$$

And now what $\begin{bmatrix} | \\ v_{\text{best}} \\ | \end{bmatrix}$ makes $\begin{bmatrix} - & v_{\text{best}} & - \end{bmatrix} A \begin{bmatrix} | \\ v_{\text{best}} \\ | \end{bmatrix}$ largest?

rotates in dir. of ~~max~~ v_{max}

dot product...

... easy to see it's just v_{max} .

Conclusion: $\vec{v}_{\text{best}} = \text{max eigenvector of } A := X^T X$

[which we learned how to find last time]

Exercise/extension: For original data $X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \end{bmatrix}$, (5)

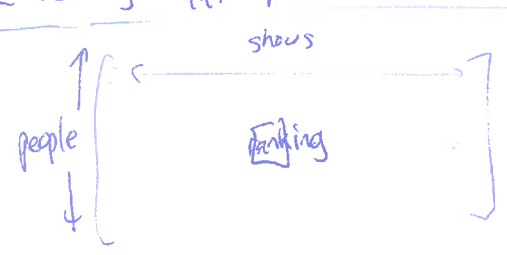
best-fitting k -dimensional subspace (we just did $k=1$),
 = "most important k dirs"
 is given by \sqrt{k} top k largest eigenvalues of $A = X^T X$.
 (which we also discussed finding.)
 eigenvectors of

↳ Computing these is called "PCA: Principal component analysis."

Applications:

① Data visualization: pick $k=2$ or 3 , project data onto k "principal components", plot.
 $y^{(i)} = \begin{bmatrix} \vec{v}_{best1} \cdot x^{(i)} \\ \vec{v}_{best2} \cdot x^{(i)} \\ \vec{v}_{best3} \cdot x^{(i)} \end{bmatrix}$

② Inferring latent features of data



$\vec{v}_{best} = [\dots]$ is how much each show satisfies "Mystery Feature" F
 $\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} x^{(1)} \cdot \vec{v}_{best} \\ x^{(2)} \cdot \vec{v}_{best} \\ \vdots \\ \vdots \end{bmatrix}$ gives how important "F" is to each person

2nd ~~and~~ best direction gives 2nd "mystery feature" F_2 , orthog to 1st, etc.
 (Maybe F 's are "quality" or "funniness" or "length" or.....)

③ Clustering data:

