

LECTURE 23

Q Let $A \in \mathbb{R}^{m \times n}$, $m \geq n$, rank r .

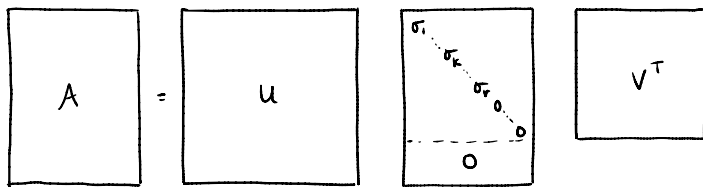
Let $k < r$. Recall the rank k approximation $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$

What is the residual $\|A - A_k\|_2$? (Note: 2-norm!)

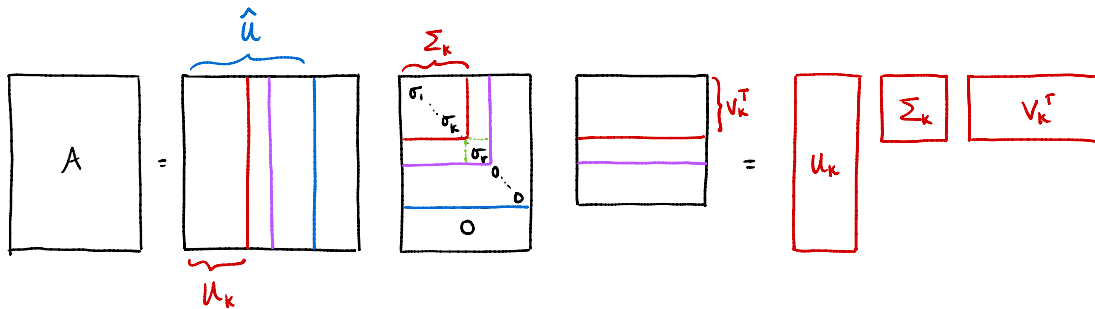
A. $\|A - A_k\|_2 =$ "biggest stretch factor of $A - A_k$ " $= \sigma_{k+1}$. $A - A_k = \sum_{i=k+1}^r \sigma_i u_i v_i^T$

Recall $A = \sum_{i=1}^r \sigma_i u_i v_i^T$ rank $r \leq n \leq m$

Visualization of low-rank approximation:



So $A - A_k$ has singular values $\sigma_{k+1}, \dots, \sigma_r$.



Aside Uniqueness of SVD? $A = \sum_{i=1}^r \sigma_i u_i v_i^T$ $A \in \mathbb{R}^{m \times n}$, rank r

- Σ is determined by A , since we require $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_r > 0$.
- For $1 \leq i \leq r$, we may replace $u_i v_i^T$ with $(-u_i)(-v_i)^T$ $(-1)^2 = 1$.
 - we can't scale them b/c they're supposed to be length 1!
 - $u_i v_i^T$ are linearly independent in a space of matrices since their sum is an exactly rank r matrix A .
- For $j \geq r+1$, we have more freedom
 - permutation of basis vectors u_j, v_j for $j \geq r+1$.
 - also mult both u_j, v_j by (-1) , as above

Condition # and choosing the rank k for approximation:

Recall $K(A) = \sigma_1 / \sigma_r$, $K(A) \approx 10^k \Rightarrow$ lose k digits of precision

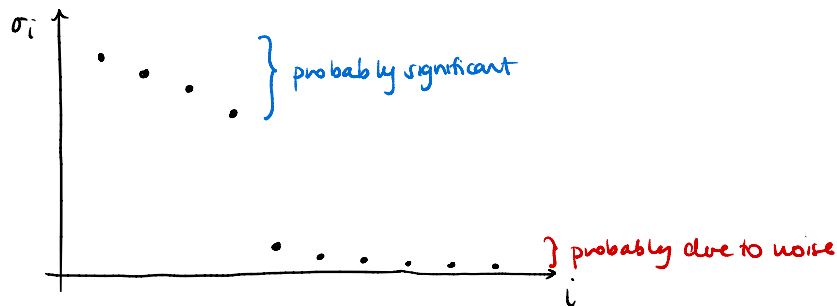
$$k = \log_{10} K(A)$$

Imagine we are given a matrix $A \in \mathbb{R}^{10 \times 10}$ that really should be lower rank but has noise + so is actually nonsingular

↳ nonsingular is "generic": $\det(A) = 0$ is way more special than $\det(A) \neq 0$.

We want to capture the essence of the data in lower rank.

Plot the singular values:



Compare: $K(A) = \sigma_1 / \sigma_{10}$ large $K(A_4)$ is much smaller.

So we might choose to focus on the approximation A_4 .



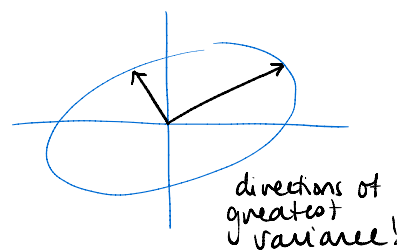
PCA = Principal Component Analysis is the same concept as SVD.

PCA is a data analysis technique that uses orthogonal transformations to convert a set of observations of maybe correlated variables into a set of linearly uncorrelated variables called principal components.

"important directions",
ie better linear combo of
your original variables
to study.

eg. The sign outside MSB (cloud of data points roughly in the shape of an ellipse in 2D)

high-dimensional data
presumably normally
distributed $\xrightarrow{\text{PCA}}$ fit a 2D ellipsoid
(ie. ellipse) to data



$$A = [a_1 \ a_2 \ \dots \ a_{100}] \xrightarrow{\text{SVD}} A_2 = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T$$

Key idea Obtain lower-dimensional visualization of a high-dimensional set of data, while preserving as much of the variance as possible.

eg. 2000 students, 8 homework scores b/w 0 and 100.

Many names: Invented by Pearson, Hotelling (1901, 1930s)

- Hotelling transform
- Karhunen-Loève transform (KLT) in signal processing
- proper orthogonal decomposition (POD) in mechE
- eigenvalue decoup of the symm. matrix $A^T A \dots$
- empirical orthogonal functions (EOF), spectral decomposition...

eg. (small example) Measure ^T temperature, ^P air pressure, ^H humidity
* I am not a meteorologist

⇒ clusters of dots in 3D. (Say, 2000 data points in \mathbb{R}^3)

Change of variables: $(T, P, H) \rightsquigarrow (T - \bar{T}, P - \bar{P}, H - \bar{H})$ so that we are centered
around $(0, 0, 0) \in \mathbb{R}^3$

$\tilde{X} = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{2000}] \in \mathbb{R}^{3 \times 2000}$ (very likely rank 3) $d=3, n=2000$

Form the (sample) covariance matrix ← will explain later

$$\mathbb{S} := \frac{1}{n-1} \tilde{X} \tilde{X}^T = \Phi \Lambda \Phi^T \quad (\text{symmetric matrix} \Rightarrow \text{SVD is eigenvalue decomp.})$$

$\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3)$ (decreasing)

note: λ_i here are scaled versions of $\sigma_i(X)$: $\lambda_i = \sigma_i^2/n$

$$\tilde{X} = USV^T \Rightarrow \tilde{X} \tilde{X}^T = USV^T V S^T U^T = U \cdot S \cdot U^T \quad (\text{the basis } \Phi \text{ is very related to } U)$$

⇒ $\Phi = [\varphi_1, \varphi_2, \varphi_3] \in \mathbb{R}^d = \mathbb{R}^3$ an orthonormal basis of \mathbb{R}^d .

- $\varphi_i \perp \varphi_j$ when $i \neq j \Rightarrow$ uncorrelated
- φ_i would be (a normalized version of) some linear combo of T, P, H. (unnormalized lengths depend on the units used)

↳ "feature extraction": maybe $\varphi_1 \sim T + H$.

⇒ T, H vary together very strongly!

cf. Exercise 3 on HW07! Face: pixels near eyes should be related.
in these, d could be $\gg n$.

Some background statistics (what is this S?)

Given: Set of observations $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ each $\bar{x}_i \in \mathbb{R}^d$ eg. $d = 128 \times 128$ pixels!
n data points

Let $X := [x_1 \dots x_n]$

Mean/average of data set: $\bar{X} = \frac{1}{n} \sum_{j=1}^n x_j$ simultaneously averaging all components, independently

Define the centered data matrix $\tilde{X} = [x_1 - \bar{x}_1 \quad x_2 - \bar{x}_2 \quad \dots \quad x_n - \bar{x}_n]$
 $= [\tilde{x}_1 \quad \dots \quad \tilde{x}_n]$

Aside Covariance b/w variables v, w w/ data $(a_i, b_i) \quad i \in \{1, \dots, N\}$ is

$$\text{cov}_{v,w} = \frac{\sum (a_i - \bar{a})(b_i - \bar{b})}{N-1} \leftarrow \begin{array}{l} \text{typo in} \\ \text{board in} \\ \text{class!} \end{array}$$

(Sample) covariance matrix $S = \frac{1}{n} \tilde{X} \tilde{X}^T \in \mathbb{R}^{d \times d}$

$\begin{array}{l} d \gg n \\ \text{in eg.} \end{array}$

$$S = \frac{1}{n} [\tilde{x}_1 \quad \tilde{x}_2 \quad \dots \quad \tilde{x}_n] \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T$$

$n-1$

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} \quad \frac{1}{n} \vec{v} \vec{v}^T = \frac{1}{n} \begin{bmatrix} v_1^2 & v_1 v_2 & v_1 v_3 \\ v_2 v_1 & v_2^2 & v_2 v_3 \\ v_3 v_1 & v_3 v_2 & v_3^2 \end{bmatrix} \Rightarrow \text{eg. } S_{12} = \frac{1}{n} \sum_{i=1}^n x_{i1} x_{i2}$$

the average v_1, v_2 value among all the vectors v .

\sim covariance
ie mutual correlation b/w
1st and 2nd entries of each
vector