

PCA : Principal Component Analysis

→ introduced by Pearson (1901)

$x^i \in \mathbb{R}^D$, D is often very large

$i = \text{sample index}$

→ x_k^i is k^{th} component of x^i $i = 1, \dots, N$

→ Data points are centered (Assumption)

$\frac{1}{N} \sum_i x_k^i = 0$ & k (mean comes out to 0)

$$\langle x_k \rangle = 0$$

$$\langle x \rangle = 0$$

$$\langle x_k \rangle = \frac{1}{N} \sum_i x_k^i$$

(avg. of x)

x = vector (when we drop the lower indices)

→ PCA performs an explicit dimension reduction.

→ Dimensionality reduction means performing the following transformation

$$x^i \in \mathbb{R}^D \longrightarrow y^i \in \mathbb{R}^d$$

$\downarrow \{d \ll D\}$

$$\rightarrow y^i = f_{\pi}(x^i) \quad \left\{ \begin{array}{l} f: \text{function of dimension} \\ \text{parameters} \end{array} \right.$$

→ generally $f_{\pi}()$ is non-linear, but in the case of PCA, $f_{\pi}()$ is a linear transformation.

→ so explicitly func² that transforms original data to reduced representation is a linear transformation of the original data.

→ PCA :- if is a linear function.

$$\rightarrow Y^i = A x^i ; A = \text{matrix}, \quad [Y^i, x^i = \text{vectors}]$$

→ if x^i & y^i have varying dimensions, no A dimensions, A must be a rectangular matrix.

\rightarrow i) $x^i \in$ row vectors (convention)
 $y^i \in$ col " (-)

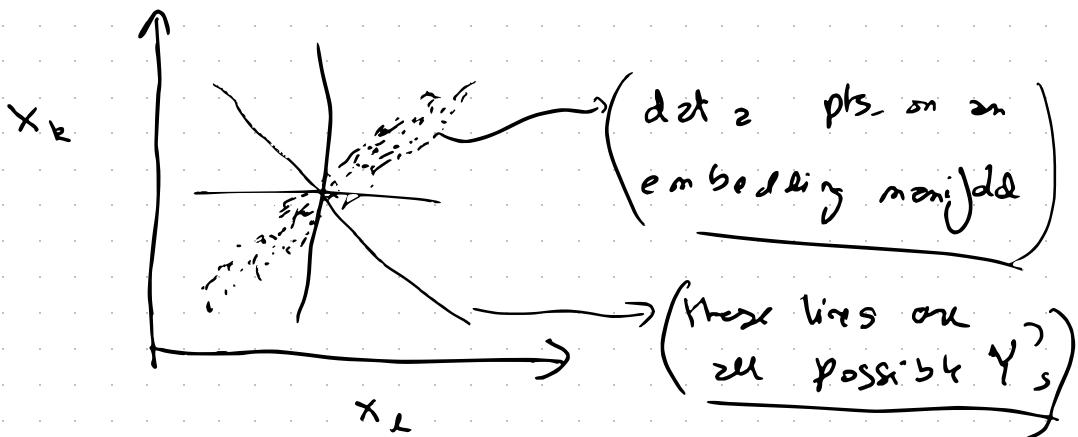
$A \rightarrow$ d rows \times D cols
(y^i dim) (x^i dim)

\rightarrow in PCA, the idea is that we determine
 C such that, $\{Y\}$ describes a
large fraction of the covariance of
the original dataset.

\rightarrow what is covariance? $\langle x \cdot x^t \rangle$ (transpose)
① $C_{D \times D} = \langle x \cdot x^t \rangle - \underbrace{\langle x \rangle \langle x^t \rangle}_{\text{if data is centered}}$
→ covariance matrix

$\therefore C_{k,l} = \langle x_k \cdot x_l \rangle$
→ (k, l) element of cov. matrix

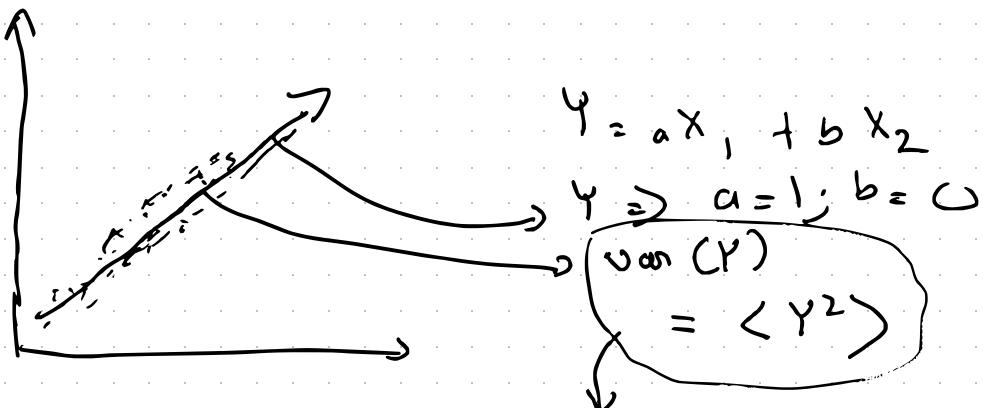
$$\therefore C_{k,l} = \frac{1}{N} \sum x_k^i x_l^i$$



C = symmetric matrix

lets say $\xrightarrow{\text{C}} C = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1.05 \end{pmatrix}$

$\Rightarrow d=1$ (target dimensionality reduction from $2 \rightarrow 1$)
 $\Rightarrow y^i = a x_1^i + b x_2^i$
 $\Rightarrow a^2 + b^2 = 1$ (2 coeffs are normalized)



∴ we use only a single component here.

→ so this is a \approx suitable choice for Y .

→ so how do we find Y ?

- ① We diagonalize the covariance matrix of $\{\underline{X}\}$.

This means we are solving the following eigenvalue eq⁼,

$$CY = \lambda Y$$

$$\rightarrow C = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1.05 \end{pmatrix} \quad : \begin{cases} \lambda_1 = 1.89 \\ \lambda_2 = 0.12 \end{cases}$$

(2 eigenvalues)

→ 2 eigen-vectors :

$$Y_1 = (0.7 \ 0.76)$$

$$Y_2 = (-0.71 \ 0.7)$$

∵ the covariance matrix is symmetric,
 → the left & right eigenvectors are the same
 → " " " " " are mutually
 orthogonal, i.e., they form an orthonormal
basis set.

$\rightarrow \lambda_1 = \text{Variance along a } Y \text{ coordinate}$
 defined by \underline{Y}_1

$\lambda_L = \dots - Y \text{ coordinate} \dots \underline{Y}_2,$

\rightarrow if \underline{Y}_1 is Y coord, Var = 1.89
 - \underline{Y}_2 " " " " " Var = 0.12
 so, now, \underline{Y}_2 is appropriate choice of reduced dimension

① Qualitative pipeline for PCA -

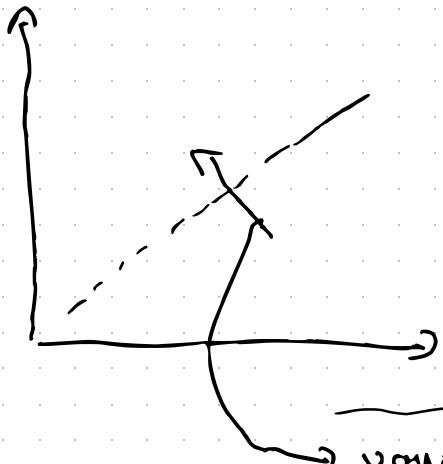
Given = dataset

- ① compute covariance $D \times D$ matrix
- ② find the eigenvalues & eigenvectors
- ③ if we know ' d ', choose as Y coordinates for first ' d ' eigenvectors of C

{ Assume : $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d$ }

\downarrow
 eigenvalues are sorted in descending order,

→ Specific Case : data pts. belong to a hyperplane

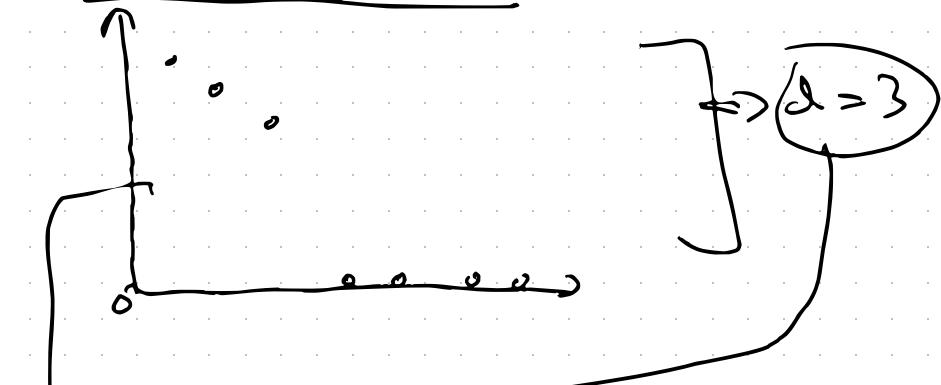


if x_i belongs to a hyperplane of dimension \underline{d} , the rank of $\underline{\underline{C}}$ will be exactly $\underline{\underline{d}}$

volume along this dir $= 0$,
 ∵ the 2 eigenvectors are orthogonal to each other,
 & all data points lie
on the other eigenvector

→ garden, if $D \approx 1000$, all data pts are on a hyperplane of dim ≈ 10
 Then 940 eigenvalues ≈ 0 , $\lambda \neq 0$ → \therefore the 10 eigenvectors will be orthogonal & have no points

→ plot the eigenvalues,



⇒ implies that the embedding manifold
can BE a hyperplane

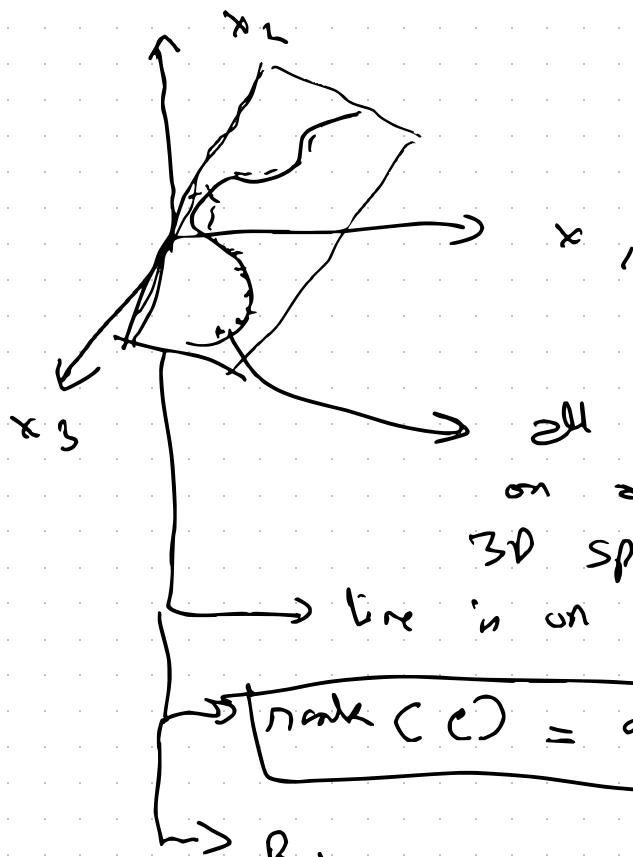
→ This is a necessary but not sufficient condition for embedding manifold to be a hyperplane

→ sufficient condition: $\text{rank}(C) = \text{ID}$
↳ (reason in next page)

→ $\text{rank}(C) \equiv r_c$

→ This only implies $(r_c \geq \text{ID})$

\Rightarrow An other case:



all data ph. lie
on ≈ 20 line in
3D space.

line is on a hyperplane

$$\text{rank}(C) = 2$$

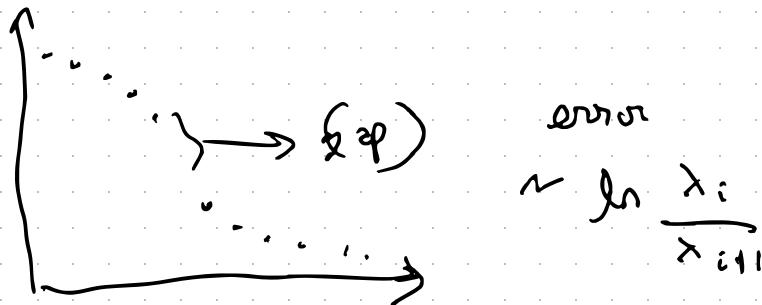
\Rightarrow But, embedding manifold is
NOT a hyperplane, \because it is
a line.

\rightarrow Here, $\text{rank} C = 2$; $ID = 1$,
 \therefore not a hyperplane.

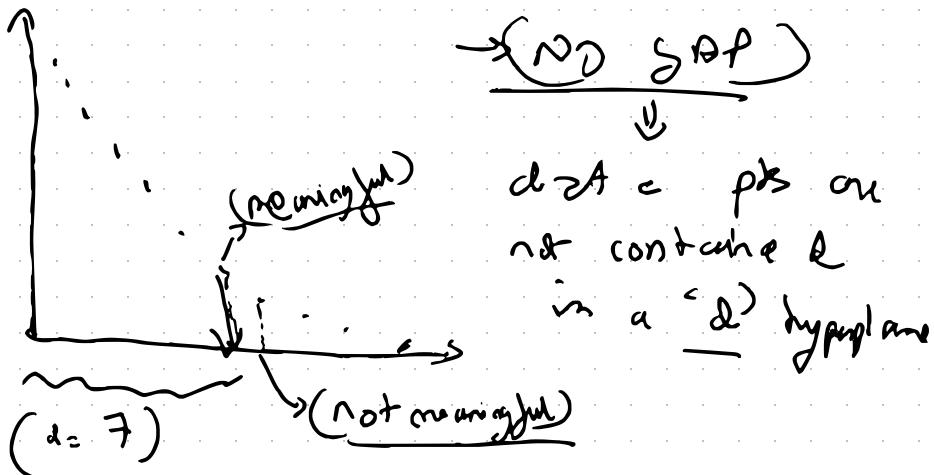
\rightarrow Real-world Scenario (1) ..

(1) General rule in PCA

$g \geq p$ in eigenvalues block



\rightarrow Real-world scenario (2) ..



\rightarrow if we say $d \geq c$, by truncating at $d = 7$, then we say that $d \geq 7$ are "meaningful" & we can't put a meaningful upper bound on the error.

→ This is because we don't have a g of so, we choose arbitrarily a point to cut-off the eigenvalues.

$$\Rightarrow \text{"Total" variance of } X = \overline{\text{Tr } C}$$

$$\text{Tr } C = \sum_{l=1}^D \lambda_l$$

(sum of eigenvalues)

$$ij \quad n_C = d,$$

$$\text{Tr } C = \sum_{l=1}^d \lambda_l$$

let's say truncation level $\hat{d} = \underline{d}$ is

$$\gamma(\hat{d}) = \sum_{l=1}^{\hat{d}} \lambda_l$$

$$\sum_{l=1}^D \lambda_l$$

→ $\eta(\alpha)$ = factor of covariance described by truncated description.

→ a common way of choosing ' α ' is to reach a pre-defined value

$$\int \eta(\alpha) d\alpha \xrightarrow{\text{for } \alpha \text{ such that}} \boxed{\eta(\alpha) \geq 0.45}$$

→ Form PCA derivation,

$$x^i \in \mathbb{R}^D \quad y^i \in \mathbb{R}^d$$

$$Y; \sim 1) Y = Ax, \quad y_k^i = \sum_{l=1}^D A_{kl} x_l^i$$

$$2) \sum_l A_{pl} A_{ml} = \delta_{km}$$

3) Trace of the covariance matrix of Y is maximal.
(if C_{cov} is given)

$$\rightarrow \text{Tr}(\text{Cov}(Y)) = \langle Y, Y^T \rangle$$

$$= \left\langle \underbrace{\sum_{n,p} A_{ln} x_n}_{Y_L} \quad \underbrace{A_{pm} x_p}_{Y_R} \right\rangle$$

$$= \sum_{n,p} A_{ln} A_{pm} C_{np} \quad \leftarrow$$

$$\because Y = Ax \quad \left\langle Y_L Y_m \right\rangle = \sum_{k,n} A_{ln} A_{mk} C_{nk}$$

$$(C_{np} = \langle x_n x_p \rangle)$$

$$\rightarrow \text{Cov}(Y) \rightarrow \text{Tr}(\text{Cov}(Y)) = \sum_k \langle Y_L Y_m \rangle$$

① Condition :- all coeffs of transformation matrix are normalized so that eigenvectors necessarily form an orthonormal basis set.

$$\sum_n A_{in}^2 = 1 \quad \forall i$$

$\xrightarrow{\text{norm}} \mathbf{y}^i = (\mathbf{A})\mathbf{x}^i$

\Rightarrow (for normalization)

$P(A)$ loss function :-

$$\text{Tr}(\mathbf{A}^t \mathbf{A})$$

$$L(\mathbf{A}) = \sum_{l,k} A_{lk} A_{lk} C_{lk}$$

(due to normalization)

(we have ' d ' normalization cond^{ns})

$$+ \sum_{l=1}^d \lambda_l \left(\sum_k A_{lk}^2 - 1 \right)$$

Lagrange multipliers

\rightarrow sum over diff. Lagrange multipliers

\rightarrow globally convex func^{ns} of A ; it has a single maxima.

$$\frac{\partial L}{\partial A_{hk}} = 0 \rightarrow \forall h, k$$

$$\Rightarrow \sum_k A_{hk} (\lambda_k - \gamma_k) \alpha_{hk} = 0$$

$\hookrightarrow \boxed{AC = \gamma A}$ (in matrix formulation)

① The coefficients of the linear transformation defining $Y = A X$ with the cond' that the trace of the covariance of Y is maximal are the 'd' eigenvectors of the matrix C

⇒ TL; DR: Recap of PCA from last lec.

① PCA (from last lecture → contd.)

$$x^i \in \mathbb{R}^D \rightarrow y^i \in \mathbb{R}^S$$

(dataset) $(S \ll D)$
(ideally)

$$y^i = A x^i \quad \left(\begin{array}{l} \text{linear transformation} \\ \downarrow \\ \text{allows } P(A) \end{array} \right)$$

rectangle matrix
 $(S \times D)$

→ each row of this matrix is an eigenvector
 of the covariance matrix.

$$\rightarrow C = \langle x x^t \rangle$$

(out of condition let's define $x^i = \text{col vector}$)

$(D \times D)$

$\rightarrow \underline{\underline{A}} \underline{\underline{C.}}$ has first δ eigenvectors of the

$\rightarrow A$ is chosen s.t., $\text{cov}(Y)$ is as large
 \hookrightarrow possible.

$$\begin{aligned}\text{cov}(Y) &= \langle Y Y^t \rangle \\ &= \langle A X X^t A^t \rangle \\ &= A \cdot C \cdot A^t \\ &\quad \downarrow \\ C &= \text{cov}(X)\end{aligned}$$

$\left[\begin{array}{l} \because Y = Ax \\ \rightarrow A \text{ can be taken} \\ \text{out of } \langle \rangle; \text{ it} \\ \text{doesn't depend on } i^{\text{th}} \\ \text{index} \end{array} \right]$

$$\therefore \boxed{\text{cov}(Y) = A \cdot C \cdot A^t}$$

\Rightarrow Goal & minimize $\text{Tr}(A^t A)$, under
 the constraint that each row of
A matrix is normalized.

$$\therefore \boxed{\sum_k A_{kl}^2 = 1 \forall k}$$

→ To impose these constraints we use
the Lagrange Multiplier term (λ_k)

↓

$$\sum_{kl} \lambda_k (A^2_{kl} - 1)$$

$$\rightarrow \Lambda = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & & \\ 0 & & \lambda_3 & \\ \vdots & & & \ddots \end{pmatrix} = (5 \times 5) \text{ diagonal matrix}$$

all the diagonal elements are Lagrange multipliers. It's a (5×5) matrix.

$$\therefore \sum_{kl} \lambda_k (A^2_{kl} - 1)$$

$$= \text{Tr} (\Lambda A A^t)$$

⇒ loss function to be minimized,

$$L(A) = \text{Tr}(A C A^t - \Lambda A A^t)$$

$$\frac{\partial L}{\partial A^t} = 0 \Rightarrow \frac{\partial}{\partial A^t} \text{Tr}(A C A^t - \Lambda A A^t) = 0$$

$$\Rightarrow \frac{\partial}{\partial A^t} (A C A^t - \Lambda A A^t) = 0$$

$$\rightarrow A C - \Lambda A = 0$$

$$\Leftrightarrow A C = \Lambda A$$

(doing this directly gives us the eigenvalue eq²)

\Rightarrow A must contain eigenvectors of C

\rightarrow We want to maximize $\text{Tr}(C \text{Co}(Y))$

$$\text{Tr}(C \text{Co}(Y)) = \text{Tr}(C A C A^t)$$

$$= \text{Tr}(\Lambda A A^t)$$

\hookrightarrow (if A satisfies

$$AC = \Lambda A$$

eigenval eq²)

$\rightarrow \because C \equiv \text{symm. matrix}$

$$\therefore A A^t = I_{S \times S}$$

(Co considers only S-orthogonal eigenvectors)

$$\begin{aligned}\therefore \text{Tr}(\text{Cov}(Y)) &= \text{Tr}(A \Lambda A^T) \\ &\rightarrow \text{Tr}(\Lambda \underbrace{I}_{5 \times 5})\end{aligned}$$

$$\begin{aligned}&= \text{Tr}(\Lambda) \\ &= \sum_{i=1}^{\delta} \lambda_i\end{aligned}$$

↓

$\delta \equiv \text{no. of chosen eigenvectors}$

→ To maximize fraction of chosen variance

we minimize $\text{Tr}(\text{Cov}(Y))$



⇒ We have to choose the δ -largest eigenvalues,

→ The largest possible $\text{Tr}(\text{Cov}(Y))$ is chosen by taking δ -largest eigenvalues,

TL; DR ↳ (1) compute covariance matrix of the data
(2) find eigenvalues & eigenvectors
(3) see the spectrum & choose truncation criterion up to 5 dimensions \rightarrow either W.R.T. $\approx y \approx p$, or W.R.T. reproducing a given fraction of the variance of the original data.


PCA allows to find explicitly the coordinates on the embedding manifold iff the "is on a hyperplane"


But, in many cases the embedding manifold is NOT on a hyperplane