

Date: 28/10/2023

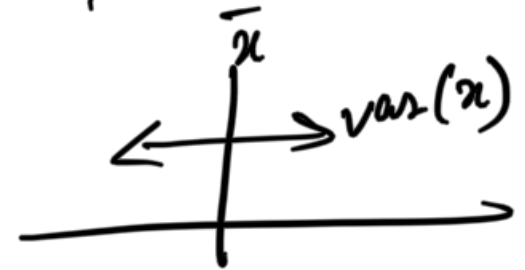
PRINCIPAL COMPONENT ANALYSIS (PCA).

What is PCA?

- Technique of dimensionality reduction of observed data.
- dimensionality reduction will depend accuracy. Cons.
- But DR \rightarrow simplify the computation Pros.
- Maintain a trade-off between

Accuracy ~

↖
Computation Complexity ~



Definition

— $E V$, $E V$, $S V$, $S Vector$.

Let, x & y are two r.v.

$$Var(x) = E[(x - \bar{x})^2], \quad \bar{x} = E(x)$$

$$\underline{Cov(x, y)} = E[(x - \bar{x})(y - \bar{y})], \quad \bar{y} = E(y)$$

Covariance matrix

$$= \begin{pmatrix} \underbrace{Cov(x, x)}_{var(x)} & Cov(x, y) \\ Cov(y, x) & \underbrace{Cov(y, y)}_{var(y)} \end{pmatrix}$$

Let, an observation data/vector of 'm' - dimension (features) / parameters

$$\underline{y_{m \times 1}} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

$$y = \underline{\mu} + Wx$$

data (observation) matrix with 'n' samples.

Original Data Matrix

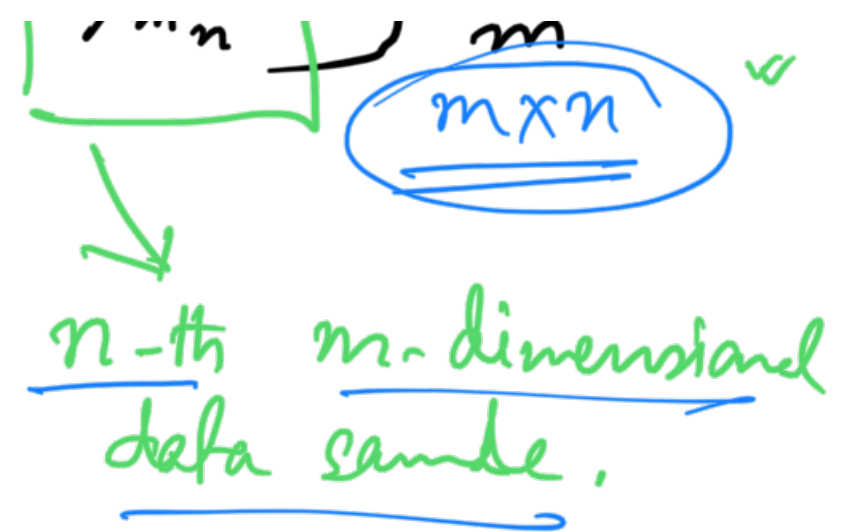
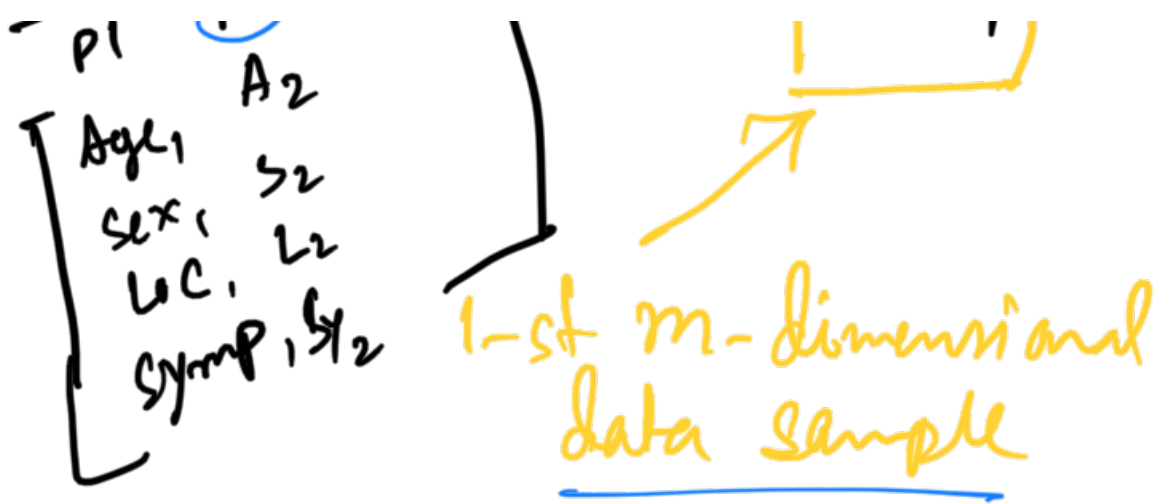
$$\underline{Y}_{m \times n} = \begin{bmatrix} y_{1,1} & y_{1,2} & \dots & y_{1,n} \\ y_{2,1} & y_{2,2} & \dots & \vdots \\ y_{3,1} & y_{3,2} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ y_{m,1} & \dots & \dots & y_{m,n} \end{bmatrix}$$

Data

Ex: Patient data
p2 p3

Dimension.

$m=3$
 (x_1, y_1, z_1)



STEPS of PCA

⑥ After getting raw data, prepare the observation matrix.

① After getting the data (observation) matrix, find the mean in all dimensions (for each

$$\bar{y}_1 = \frac{\sum_{i=1}^n y_{1i}}{n} = \frac{y_{11} + y_{12} + \dots + y_{1n}}{n}$$

$$\vdots$$

$$\bar{y}_m = \frac{\sum_{i=1}^n y_{mi}}{n}$$

② STANDARDIZATION ✓

Subtracting mean from each data across each dimension.

Stat. data matrix

Data points.

$$\begin{bmatrix} (y_{11} - \bar{y}_1) & (y_{12} - \bar{y}_1) & \dots & (y_{1n} - \bar{y}_1) \\ \vdots & \vdots & \ddots & \vdots \\ (y_{m1} - \bar{y}_1) & (y_{m2} - \bar{y}_1) & \dots & (y_{mn} - \bar{y}_1) \end{bmatrix}_{m \times n}$$

[illegible]

(3) Calculate all possible cov across the dimensions, and form the covariance matrix. n, y, z

most computation intensive.

$$C = COV = \begin{pmatrix} \text{cov}(y_1, y_1) & \text{cov}(y_1, y_2) & \dots & \text{cov}(y_1, y_m) \\ \text{cov}(y_2, y_1) & & & \\ \vdots & & & \\ \text{cov}(y_m, y_1) & & & \text{cov}(y_m, y_m) \end{pmatrix}$$

i.e. finding the dependencies among different dimensions

Covariance matrix is always SPD. $m \times m$

(4) Find the Eigen vector (v_i) ✓
and Eigen values (λ_i) of cov matrix

$$C = V \Lambda V^T \text{ (EVD)}$$

Choose components in non-decreasing order of eigen-values
s.t.

$$C = \underbrace{v_1 \lambda_1 v_1^T + v_2 \lambda_2 v_2^T + \dots + v_m \lambda_m v_m^T}$$

where $\lambda_1 > \lambda_2 > \dots > \lambda_m$ ✓

⑤ Find Featured vectors:

After arranging eigen-vectors according to non-increasing order of eigen-values, one can keep all the eigen-vectors or can discard less significant eigen-vectors (associated with smaller eigen-values) and form the remaining ones a matrix. This matrix (collection of eigen-vectors of the components that are

decided to be kept)

ex! we have
= λ_1 λ_2 ... λ_k ... λ_m : eigen-values

associated
eigen vectors

$$\begin{pmatrix} v_{11} \\ v_{12} \\ \vdots \\ v_{1m} \end{pmatrix}$$

i.e.

v_1

$$\begin{pmatrix} v_{21} \\ v_{22} \\ \vdots \\ v_{2m} \end{pmatrix}$$

v_2

$$\begin{pmatrix} v_{k1} \\ v_{k2} \\ \vdots \\ v_{km} \end{pmatrix}$$

v_k

$$\begin{pmatrix} v_{m1} \\ v_{m2} \\ \vdots \\ v_{mm} \end{pmatrix}$$

v_m

as $\lambda_1 > \lambda_2 > \dots > \lambda_k > \dots > \lambda_m$, we can
keep upto certain eigen-vector, say 'K'.

Then, the Feature Vector/matrix

$$= \begin{pmatrix} v_{11} & v_{21} & \dots & v_{k1} \\ v_{12} & v_{22} & & v_{k2} \\ \vdots & \vdots & & \vdots \\ v_{1m} & v_{2m} & & v_{km} \end{pmatrix} \quad \underline{m \times k}$$

and $k < m$

This is the steps of dimensionality reduction.

⑥ Resulting data:

$$\begin{aligned} & \text{Final / Recasted Data Matrix (Set)} \quad \underline{k \times n} \quad \begin{array}{l} \text{data with PCs.} \\ \text{New-dimension} \end{array} \\ & = \text{Featured Vector}^T \times \text{Original Data Matrix} \quad \underline{m \times n} \\ & \quad \underline{k \times m} \end{aligned}$$

Diagram: A matrix with dimensions m (rows) and k (columns) is shown on the left, with elements v₁₁, v₁₂, ..., v_{1m}, v₂₁, v₂₂, ..., v_{2m}, ..., v_{k1}, v_{k2}, ..., v_{km}. Arrows point from this matrix to the 'Final / Recasted Data Matrix (Set)' and the 'Featured Vector'.

Therefore dimension of data reduced from 'm' to 'k'

NOTE:

- * PCA transforms data linearly into new properties that are not correlated with each other (as eigen-vectors are orthonormal to each other)
i.e. PCs are uncorrelated as each of them associated with eigen vectors.
- * Principal components are constructed in such a manner that the first (few) component(s) account(s) for the largest possible variance(s) in the data set.

SPD: SYMMETRIC POSITIVE DEFINITE
Matrix.

An $n \times n$ symmetric real (centrosymmetric) matrix M is called positive definite.

$$\text{If } x^T M x > 0 \quad \forall \text{ non-zero } x \text{ in } \mathbb{R}^n$$

$$x^T M x > 0 \quad \forall x \in \mathbb{R}^n \setminus \{0\}.$$

Ex:-

3-dimensional (variable/features) data set
with 4 observations (samples).

Dimension $\begin{Bmatrix} l \\ h \\ b \end{Bmatrix}$ $\begin{bmatrix} 7 & 4 & 6 & 8 \\ 4 & 1 & 3 & 2 \\ 3 & 8 & 7 & 9 \end{bmatrix}$

observation.

Consider the PCA and
reconst upto 2-dimension
(l & h).