

PCA

• Motivation

- Usually, datasets have lot more dimensions than required.
- Ex: Movement of spring in 1-d captured by three cameras at different angles.
- Therefore reduce dimensionality of data to get rid of redundant information.
Ex: Remove two extra cameras. Only one camera needed to capture 1-d motion.

• Working

- Suppose the spring oscillating in x-direction
If the cameras were recording three separate dimensions (x, y and z), we could easily remove 2 of them and keep the one that is capturing the 1-d movement

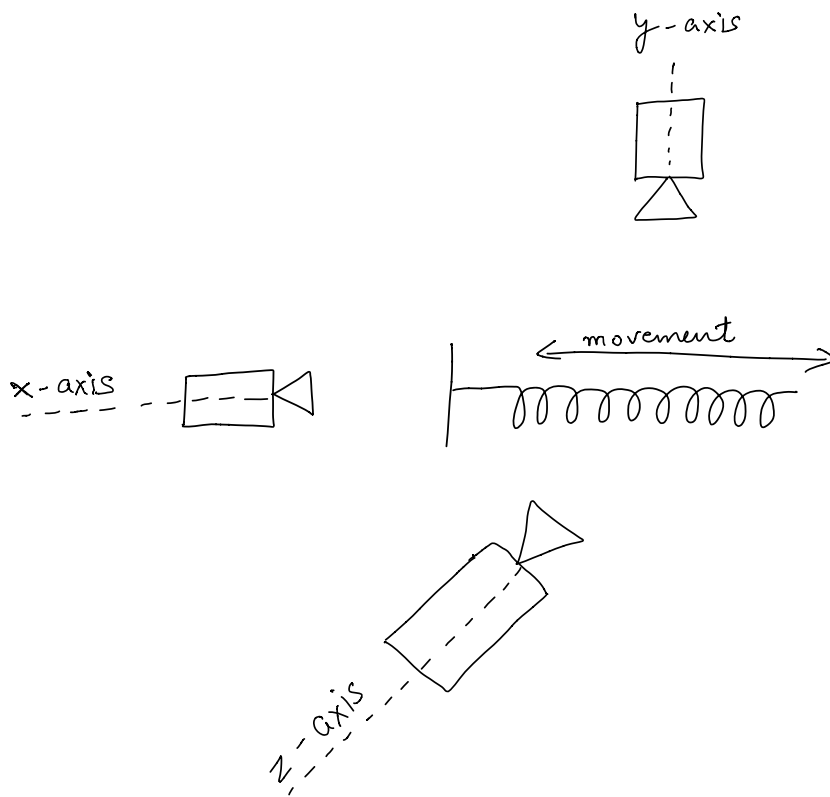


Figure:- Camera x should be deleted.
 One of camera y or z should be deleted to capture the movement of the spring uniquely.

- From the figure, it is clear that two dimensions can be removed to capture non-redundant information. But, only if the dimensions are orthonormal to each other. In real world data sets, it is unlikely that we get data where attributes are orthonormal..

- Change of basis:

- We need to find orthonormal basis such that an entire dimension (column) could be removed from the dataset without loss of information:

$$\begin{array}{ccc} & \text{original} \\ & \text{data} \\ & \uparrow \\ P & X & = & Y \\ \downarrow & & & \downarrow \\ \text{Principal} & & & \text{new data} \\ \text{component} & & & \text{with orthonormal} \\ \text{matrix} & & & \text{columns.} \end{array}$$

- Need to find P that transforms X to Y where Y has orthonormal columns.

- Maximising Variance

- Another nice thing would be to capture maximum information (= maximum variance) in the first dimension of Y , second highest information in second dimension, and so on.

- In other words, we want Y where the covariance matrix is diagonal (= covariance between any two dimensions is 0)

$$P X = Y$$

such that C_Y is diagonal.

$$C_Y = \frac{1}{n} Y Y^T$$

$$= \frac{1}{n} (P X) (P X)^T$$

$$= \frac{1}{n} P X X^T P^T$$

$$= P \left(\frac{1}{n} X X^T \right) P^T$$

$$C_Y = P C_X P^T$$

Diagonalizing C_X :

$$C_X = P (E D E^T) P^T$$

$$= (P E) D E^T P^T$$

$$C_X = (P E) D (P E)^T$$

By selecting $P = E^T \Rightarrow P E = E^T E = I$

$$\Rightarrow \boxed{C_Y = D}$$

P is the (transpose of) eigenvector matrix of C_X .

Moreover, if D is arranged such that $d_1 \geq d_2 \geq \dots$,
then $\text{var}_1 \geq \text{var}_2 \geq \dots$ where var_i is variance
of i^{th} dimension of Y .

Summary

- Need P in $P X = Y$.
 P is principal component matrix.
 Y is new reduced data.
- Need $\text{var}_1 \geq \text{var}_2 \geq \dots$ in columns of Y .
 \Rightarrow covariance between any pair of features of $Y = 0$. $\begin{pmatrix} \text{cov}_{i,j} = 0 \\ i \neq j \end{pmatrix}$
 \Rightarrow covariance matrix of $Y (= \frac{1}{n} Y Y^T)$ is diagonal matrix.
 $\Rightarrow P$ is eigenvector matrix of covariance matrix of X .

Note:- Formula used for covariance matrix ($C_x = \frac{1}{n} X X^T$) works only if the data is mean-centered.

Matrix sizes

f = number of features
 s = number of samples

$$\bullet \quad \begin{matrix} P \\ (f, f) \end{matrix} \times \begin{matrix} X \\ (f, s) \end{matrix} = \begin{matrix} Y \\ (f, s) \end{matrix}$$

$$\bullet \quad \begin{matrix} C_x \\ (f, f) \end{matrix} = \frac{1}{s} \begin{matrix} X \\ (f, s) \end{matrix} \times \begin{matrix} X^T \\ (s, f) \end{matrix}$$

$$\bullet \quad \begin{matrix} P \\ (f, f) \end{matrix} = \begin{matrix} E^{-1} \\ (f, f) \end{matrix} = \begin{matrix} E^T \\ (f, f) \end{matrix}$$