

# Principal Component Analysis (PCA)

## Dimensionality Reduction

explicit structure correlation  $\rightarrow$  noise  $\rightarrow$  compact representation of data  
High dimension data  $\rightarrow$  compression algorithm  $\rightarrow$  low dimension data

## Statistics of datasets - compact way to describe data

Average Mean  $E(D) = \frac{1}{N} \sum D_i$   
Spread 1D Variance  $E(x - \mu)^2$   
nD Covariance  $E(x - \mu_x)(y - \mu_y) \rightarrow$  Covariance Matrix  $D \times D$

## Linear Transformation

$$E(aD + b) = aE(D) + b$$

$$\text{Var}(aD + b) = a^2 \text{Var}(D)$$

$$\text{Var}(AD + b) = A \text{Var}(D) A^T$$

## Inner product

Dimensionality Reduction  $\rightarrow$  Compression

Orthogonality Distance vector  
Angle vector space

## Inner Product

Length/norm  
finite number of entries

$$\|x\| = \sqrt{\langle x, x \rangle}$$

$$\cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

similarity

orthogonal  $\theta = 90^\circ$ , nothing in common except origin

## function

Discrete  
Continuous

$$\langle u, v \rangle = \int_a^b u(x) v(x) dx \quad \text{orthogonality} = 0$$

## random variable

$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$  if  $X$  &  $Y$  are uncorrelated random variable

$$\langle x, y \rangle = \text{Cov}(X, Y)$$

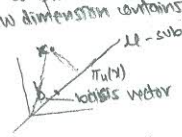
$$\|x\| = \sigma(x)$$

$$\cos \theta = \frac{\langle x, y \rangle}{\sigma(x) \sigma(y)}$$

## Orthogonal Projection

High dimension data  $\rightarrow$  low dimension data

few dimension contains most information  
to 1-D subspace



orthogonal projection on subspace  $U$  with basis vector  $b$

$$1) \pi_U(x) \in U \Rightarrow \exists \lambda \in \mathbb{R} : \pi_U(x) = \lambda b$$

$$2) \langle b, \pi_U(x) - x \rangle = 0 \rightarrow \text{orthogonal}$$

$$\Leftrightarrow \lambda = \frac{\langle b, x \rangle}{\|b\|^2}$$

$$\pi_U(x) = \frac{\langle b, x \rangle}{\|b\|^2} b = \frac{b b^T}{\|b\|^2} x \quad \text{in terms of dot product}$$

projection matrix

to n-D subspace

dimension vector to m dim subspace

$$1) \pi_U(x) = \sum_{i=1}^m \lambda_i b_i = B \lambda$$

$$2) \langle \pi_U(x) - x, b_i \rangle = 0 \text{ for } i=1, 2, \dots, m \Rightarrow \lambda = (B^T B)^{-1} B^T x$$

$$\Leftrightarrow \lambda = (B^T B)^{-1} B^T x$$

$$\pi_U(x) = B (B^T B)^{-1} B^T x$$

ONB  $\pi_U(x) = B B^T x$

Projection Matrix

$$B = [b_1 \ b_2 \ \dots \ b_m]$$

basis vector

## PCA

### Derivation

Data set  $X = \{x_1, x_2, \dots, x_n\} \quad x_i \in \mathbb{R}^D$

1)  $x_n = \sum_{i=1}^M b_i \alpha_i$  - every vector in  $\mathbb{R}^D$  can be represented as linear combination of basis vector

2)  $b_i$  in  $\mathbb{R}^D$  - orthogonal basis

3)  $B = [b_1 \ b_2 \ \dots \ b_m]$

$\tilde{x} = B B^T x$  - orthogonal projection of  $x$  on subspace spanned by  $B$

$\tilde{x}$  is orthogonal projection of  $x$  on subspace spanned by  $M$  basis vector  $b_j$  where  $j=1$  to  $M$

$$\tilde{x}_n = \sum_{j=1}^M b_j \alpha_j \quad \Leftrightarrow \quad x_n = \sum_{j=1}^M b_j \alpha_j + \sum_{j=M+1}^D b_j \alpha_j$$

$$x_n = \left( \sum_{j=1}^M b_j b_j^T \right) x_n + \left( \sum_{j=M+1}^D b_j b_j^T \right) x_n$$

$$\Leftrightarrow x_n - \tilde{x}_n = \left( \sum_{j=M+1}^D b_j b_j^T \right) x_n = \sum_{j=M+1}^D (b_j^T x_n) b_j$$

Displacement vector

$$J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \sum_{j=M+1}^D b_j^T \left( \frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_j$$

data covariance (S)

$$= \sum_{j=M+1}^D b_j^T S b_j = \text{trace} \left( \sum_{j=M+1}^D b_j b_j^T S \right)$$

minimize loss matrix - average square reconstruction error

= minimize variance of data projected onto subspace orthogonal to principal subspace

= retain as much as variance after projection as possible.

orthonormal basis for m-dim subspace

Optimize  $\rightarrow$  Lagrangian  $\rightarrow$  gradient w.r.t  $\lambda$  & basis vector = 0

$$b_j^T S b_j = \lambda_j \quad b_j^T S b_j = \lambda_j \quad J \text{ is minimum}$$

$\lambda$  eigen value is minimum

choose basis vectors that span ignored subspace to be eigen vectors of S that belong to smallest eigen values.

= principal subspace - spanned by eigen vectors belonging to M largest eigen values of S.

## PCA

find lower dimension representation

$$\tilde{x}_n \leftarrow x_n$$

fewer basis vector

Assume - Centered Data  $E[x] = 0$

- ONB  $b_1, \dots, b_D$

$$\tilde{x}_n = \sum_{i=1}^M b_i \alpha_i + \sum_{i=M+1}^D b_i \beta_i \in \mathbb{R}^D$$

coordinate code

span principal subspace

ignore in PCA

Given  $x$ , find  $B$  in ONB  $b_i$ , average square reconstruction error is minimized

$$J = \frac{1}{N} \|x_n - \tilde{x}_n\|^2 \quad \frac{\partial J}{\partial b_i} = 0 \quad \frac{\partial J}{\partial \alpha_i} = 0$$

$$\frac{\partial J}{\partial \alpha_i} = -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i \quad \frac{\partial J}{\partial b_i} = b_i^T (x_n - \tilde{x}_n) = 0$$

$$\Leftrightarrow B^T (x_n - \tilde{x}_n) = 0 \quad \Leftrightarrow B^T x_n = B^T \tilde{x}_n$$

optimal coordinate of  $\tilde{x}_n$  w.r.t our basis

orthogonal projection of coordinate of original data onto  $i$ th basis vector that spans our principal subspace.

eigen vector are orthogonal to each other

eigen vector with largest eigen value points in the direction of data with largest variance and variance in that direction is given by corresponding eigen value

# PCA - Algorithm

## Steps 1) Data Normalization

a)  $x_i^* = x_i - \bar{x}$  mean normalization [not necessary] numerical difficulties  $\rightarrow$  covariance

b)  $x_i^* = x_i / \sigma_x$  [unit free] variance in all dimension = 1, while correlation is intact

## 2) Covariance Matrix (S)

Eigen values  
Eigen vectors

3) Project to principal subspace that is spanned by the eigenvectors that belong to largest eigenvalues

## Data covariance Matrix

in D dimension =  $D \times D \rightarrow$  computational intensive to calculate

Dataset  $x_1, x_2, \dots, x_N \in \mathbb{R}^D$

$$S = \frac{1}{N} X^T X$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \in \mathbb{R}^{N \times D}$$

$N < D$  m. of datapoints < dimensionality of data

eigen value  
eigen vector

$O(n^3)$

Covariance Matrix (S)

$$\text{Rank}(S) = N$$

$D - (N+1)$  eigen values = 0

$\rightarrow$  full rank matrix

rows & columns are linearly dependent

Convert  $D \times D$  covariance matrix into

$\rightarrow$  full rank  $N \times N$  covariance matrix without eigen values 0.

$$S b_i = \lambda_i b_i$$

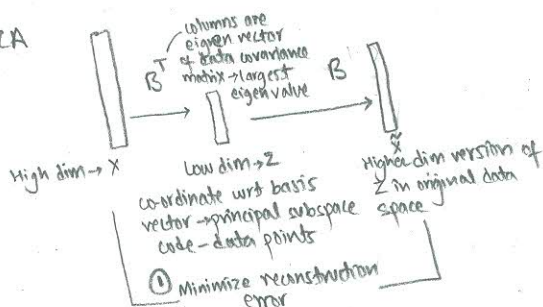
$$\frac{1}{N} X^T X b_i = \lambda_i b_i \Leftrightarrow$$

$$\underbrace{\frac{1}{N} X^T X}_{\substack{\in \mathbb{R}^{N \times N} \\ \text{covariance matrix} \\ \text{non zero eigen values}}} \underbrace{X b_i}_{c_i} = \underbrace{\lambda_i}_{\text{eigen value}} \underbrace{X b_i}_{c_i}$$

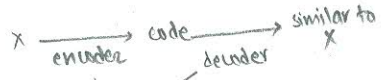
Recover original eigen vector  $\rightarrow$  PCA

$$\frac{1}{N} X^T X X^T c_i = \lambda_i X^T c_i$$

## Other Interpretation of PCA

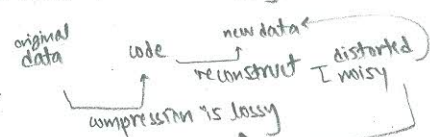


## Auto encoder



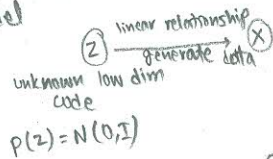
linear - linear mapping - minimize squared encoding loss  $\sim$  PCA  
non-linear - non linear mapping - Deep autoencoder - Deep Neural Network

## Information Theory



information contained in data  
 $\rightarrow$  minimizing variance of data when projected on principal subspace

## Latent Variable Model



$$P(z) = N(0, I)$$

$$x = Bz + \mu + \epsilon \quad \epsilon \sim N(0, \sigma^2 I) \text{ isotropic error}$$

$$\text{likelihood of model } P(x/z) = N(x | Bz + \mu, \sigma^2 I)$$

$$\text{Marginal likelihood } P(x) = \int P(x/z) P(z) \cdot dz = N(x | \mu, BB^T + \sigma^2 I)$$

5) Use Maximum Likelihood Estimation (MLE)  $\rightarrow (\mu, B, \sigma)$  parameter

$\mu$  = mean of data

$B$  = matrix of eigen vectors that contains largest eigen values

Find low dim code of data point.

Bayes Theorem to invert linear relationship between  $z$  &  $x$

$$P(z/x) = \frac{P(x/z) \times P(z)}{P(x)}$$

likelihood  $\times$  prior  $\rightarrow$  marginal likelihood