

FYS 5429/9429, MARCH 22, 2023

## BASICS of PCA

Data set  $X \in \mathbb{R}^{n \times p}$   
# samples      # features

$$X = \begin{bmatrix} x_0 & x_1 & \dots & x_{p-1} \\ 1 & 1 & & 1 \end{bmatrix}$$

$$\text{cov}[x_i, x_j] = \frac{1}{n} \sum_k (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)$$

$$x_i = \begin{bmatrix} x_{0i} \\ x_{1i} \\ \vdots \\ x_{ni} \end{bmatrix}$$

$$\bar{x}_i = \frac{1}{n} \sum_k x_{ki}$$

$$\tilde{x}_{ki} = x_{ki} - \bar{x}_i \rightarrow x_{ki}$$

$$X = \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \\ x_{20} & x_{21} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ x_0 & x_1 \\ 1 & 1 \end{bmatrix}$$

$$\text{cov}[x_0, x_1] = \frac{1}{n} \sum_k x_{k0} x_{k1}$$

$$= \frac{1}{n} (x_{00}x_{01} + x_{10}x_{11} + x_{20}x_{21})$$

$$= \begin{bmatrix} x_{00} & x_{10} & x_{20} \end{bmatrix} \begin{bmatrix} x_{01} \\ x_{11} \\ x_{21} \end{bmatrix}$$

$$= x_0^T x_1 \frac{1}{n}$$

$$\text{cov}[X] = \frac{1}{n} X X^T$$

$$X \in \mathbb{R}^{n \times p} \quad \uparrow \text{Design matrix}$$

$$\text{cov}[X] = \frac{1}{n} \begin{bmatrix} x_0^T x_0 & x_0^T x_1 \\ x_1^T x_0 & x_1^T x_1 \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_0^2 & \text{cov}[x_0, x_1] \\ \text{cov}[x_1, x_0] & \sigma_1^2 \end{bmatrix}$$

$$\text{corr}[X] = \frac{\text{cov}[x_i, x_j]}{\sqrt{\text{var}[x_i] \text{var}[x_j]}}$$

$$\text{corr}[X] = \begin{bmatrix} 1 & \text{corr}[x_0, x_1] \\ \text{corr}[x_1, x_0] & 1 \end{bmatrix}$$

SVD : singular value decomp

$$X = U \Sigma V^T$$

$$U \in \mathbb{R}^{n \times n} \quad \Sigma \in \mathbb{R}^{n \times p}$$

$$V \in \mathbb{R}^{p \times p}$$

$$U^T U = U U^T = I_{n \times n}$$

$$V V^T = V^T V = I_{p \times p}$$

$$\Sigma = \begin{bmatrix} \lambda_0 & & & \\ & \lambda_1 & & \\ & & \ddots & \\ & & & \lambda_{p-1} \\ & & & & \ddots & 0 \\ & & & & & \ddots & 0 \end{bmatrix}$$

$$\lambda_0 > \lambda_1 > \lambda_2 > \dots > \lambda_{p-1} > 0$$

$$\Sigma = \begin{bmatrix} \tilde{\Sigma} \\ 0 \end{bmatrix}$$

$$\tilde{\Sigma} \in \mathbb{R}^{p \times p}$$

$$\text{cov}[x] = \frac{1}{n} X X^T$$

$$= \frac{1}{n} (V \Sigma^T \underbrace{U^T U}_{I} \Sigma V^T)$$

$$= \frac{1}{n} V \Sigma^T \Sigma V^T \in \mathbb{R}^{p \times p}$$

$$X^T X = V \Sigma^T \Sigma V^T$$

multiply with  $V$  from right

$$(X^T X) V = V \Sigma^T \Sigma \underbrace{V^T V}_1$$

$$= V \Sigma^T \Sigma = V \Sigma^2$$

$$V = \begin{bmatrix} | & | & & | \\ v_0 & v_1 & \dots & v_{p-1} \\ ( & 1 & & 1 \end{bmatrix}$$

$$(X^T X) v_i = v_i \lambda_i^2$$

multiply with  $\frac{1}{n}$

$$\begin{aligned} \left( \frac{1}{n} X^T X \right) v_i &= \frac{1}{n} v_i \lambda_i^2 = \\ \text{Cov}[x] &= v_i \frac{\lambda_i^2}{n} \end{aligned}$$

The covariance matrix  $x$

$$\text{Cov}[\bar{x}] = \frac{1}{n} X^T X$$

$$= E[X^T X]$$

$$S^T \text{Cov}[\bar{x}] S = D$$

$$= \begin{bmatrix} \sigma_0^2 & & \\ & \ddots & \\ & & \sigma_{p-1}^2 \end{bmatrix}$$

$$S^T S = S S^T = \mathbb{1}$$

$$E[S^T X^T X S] =$$

$$S^T E[X^T X] S =$$

$\text{Cov}[y]$  is uncorrelated matrix (is diagonal), along the diagonal we have the variance

$$\text{var}[y], \sigma_0^2/n > \sigma_1^2/n$$

$$\dots \sigma_{p-1}^2/n$$

$$S = \begin{bmatrix} s_0 & s_1 & \dots & s_{p-1} \\ 1 & 1 & & 1 \end{bmatrix}$$

$S = V$  from the SVD

$$\lambda_0^2/n = \sigma_0^2/n \Rightarrow$$

$$\lambda_1^2/n = \sigma_1^2/n$$

$$\text{var}[y_0] > \text{var}[y_1] > \dots > \text{var}[y_{p-1}]$$

The PCA theorem states that we can define an optimal dimension  $m < p$  ( $m \ll p$ ), which defines a cutoff

$$\sum_{i=0}^{m-1} \sigma_i^2 \leq \text{cutoff}$$

$\Rightarrow$  Reduced dimensionality

PCA theorem

We assume there is an orthogonal transformation

$$S S^T = S^T S = \mathbb{1}$$

$$S = \begin{bmatrix} s_0 & s_1 & \dots & s_{p-1} \\ 1 & 1 & & 1 \end{bmatrix}$$

$$S_0^T S_0 = 1$$

we define a Lagrangian

$$S_0^T \text{Cov}[X] S_0 + \lambda_0 (1 - S_0^T S_0)$$

optimize wrt  $S_0$  and  $\lambda_0$   
wrt  $S_0^T$ :

$$\text{Cov}[X] S_0 = \lambda_0 S_0$$

multiply by  $S_0^T \Rightarrow$

$$S_0^T \text{Cov}[X] S_0 = \lambda_0$$

= eigenvalue = variance of  
 $\text{Cov}[X]$  for  $\sigma_0^2/n$

$S_0$  is the first principal component.

Next  $S_1$ ,  $S_1^T S_0 = 0$

$$\begin{aligned} \mathcal{L} = & S_1^T \text{Cov}[X] S_1 + \lambda_1 (1 - S_1^T S_1) \\ & + \kappa S_1^T S_0 \end{aligned}$$

Take derivatives wrt  $\lambda_1, \kappa, S_1$   
wrt  $S_1$ ,

$$\text{cov}[x] S_1 + \gamma/2 S_0 = \lambda S_1$$

multiply from the left with  $S_0^T$

$$\underbrace{(S_0^T \text{cov}[x])}_{\lambda_0 S_0^T} S_1 + \gamma/2 \underbrace{S_0^T S_0}_1$$

$$= \lambda S_0^T S_1 \Rightarrow$$

$$\gamma = 2[\lambda_1 - \lambda_0] \underbrace{S_0^T S_1}_{=0} \Rightarrow \gamma = 0$$

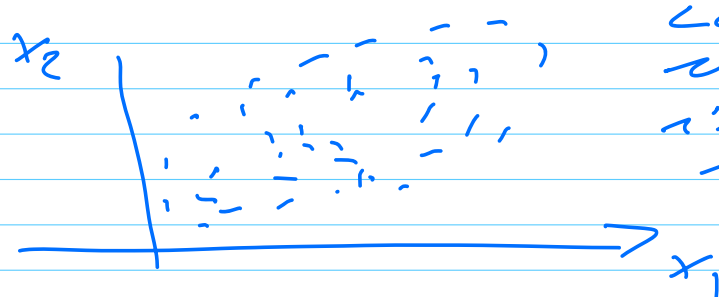
$$S_1^T \text{cov}[x] S_1 = \lambda_1 = \text{var}[y_1]$$

$$\text{cov}[y] = S^T \text{cov}[x] S$$

$$S_i^T S_j = \delta_{ij}, \text{ can construct by induction}$$

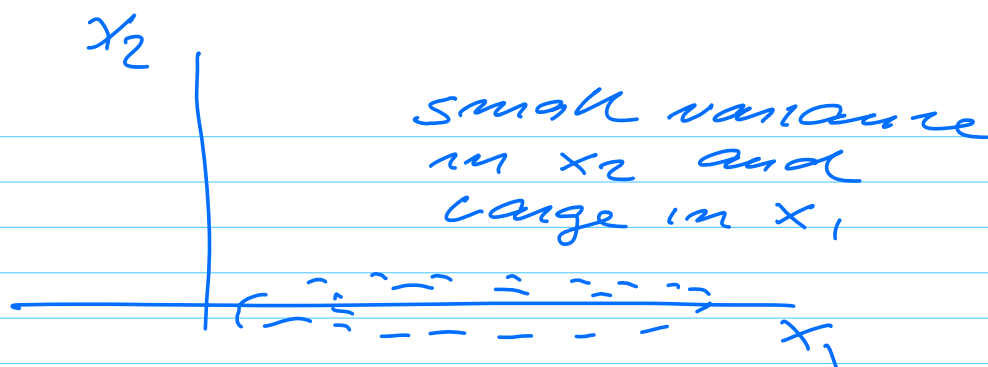
$$\begin{aligned} \mathcal{L} = & S_n^T \text{cov}[x] S_n + \\ & \lambda_n (1 - S_n^T S_n) \\ & + \sum_{j=0}^{n-1} \gamma_j S_n^T S_j \end{aligned}$$

Two-dim set



Large  
variance  
in  $x_1$  and  
in  $x_2$





could drop  $x_2$  as a dof.

Reconstruction error

$$J(x) = \frac{1}{n} \sum_{i=0}^{n-1} (x_i - \tilde{x}_i)^2$$

$\tilde{x}_i$  is our approximation

For autoencoders, with linear dependence

$$\tilde{x} = VWx$$

$$\tilde{J}(x) = \frac{1}{n} \sum_{i=0}^{n-1} (x_i - (VW)x_i)^2$$

we have an approximation

$$\underline{S_0 S_0^T} x = \tilde{x}$$

outer product (matrix)

$$S_0 = \begin{bmatrix} s_{00} \\ s_{10} \end{bmatrix}$$

$$S_0^T S_0 = s_{00}^2 + s_{10}^2$$

$$S_0 S_0^T = \begin{bmatrix} S_{00}^2 & S_{00} S_{10} \\ S_{10} S_{00} & S_{10}^2 \end{bmatrix}$$

The PCA theorem from the reconstruction error gives the smallest reconstruction error if  $S_0$  are the eigenvectors of the cov  $[x]$  with largest variance  $\sigma_0^2$

$$\arg \min_{S_0} \frac{1}{n} \|x - S_0 S_0^T x\|_2^2$$

Think of the autoencoder with linear dependence

$$\arg \min_{V, W} \frac{1}{n} \|x - VWx\|_2^2$$

$$VW \propto S_0 S_0^T$$