

# Chapter 8: Principal Component Analysis in high dimensions.

8.1  $V_{j-1}$  denotes a collection of subspaces of dimension  $j-1$ .

$$\max_{u \in V^\perp \cap S^{d-1}} \langle u, Qu \rangle = \arg \max_{v \in S^{d-1}} \langle u, Qu \rangle$$

$$= \arg \max_{v \in S^{d-1}} E[\langle u, Qu \rangle^2]$$

$$= \sum_{j=1}^n \lambda_j(Q) (u_j \otimes u_j) \text{ where } u_j \otimes u_j = u_j^T \cdot u_j$$

$$\min_{V \in V_{j-1}} \sum_{j=1}^n \lambda_j(Q) (u_j \otimes u_j) = \lambda_j(Q)$$

$$= \min_{V \in V_{j-1}} \max_{u \in V^\perp \cap S^{d-1}} \langle u, Qu \rangle$$

$$\lambda_j(Q) = \min_{V \in V_{j-1}} \max_{u \in V^\perp \cap S^{d-1}} \langle u, Qu \rangle$$

(Unitarily Invariant)

$$\|M\| = \|VMV\| \text{ when } d_1 \leq d_2; V_{d_1 \times d_1}; M_{d_1 \times d_2}; U_{d_2 \times d_2}$$

0.2

a) i) Frobenium Norm

Method #1:

$$\|VMU\|_F = \sqrt{|\sum V_i \sum M_i \sum U_i|^2}$$

$$\leq \sqrt{\sum |V_i|^2 \sum |M_i|^2 \sum |U_i|^2}$$

$$\leq \sqrt{\sum |M_i|^2} \text{ when } \sum |V_i|^2 = \sum |U_i|^2 = 1$$

$$\leq \|M\|_F$$

Method #2:

$$\|VMU\|_F = \sqrt{|\sum V_i \sum M_i \sum U_i|^2}$$

$$\leq \sqrt{\text{Tr}(V_i)^2 \cdot \text{Tr}(M_i)^2 \cdot \text{Tr}(U_i)} \leq \|M\|_F$$

ii) Nuclear Norm:

$$\|VMU\|_{Nuc} = \text{Tr}(\sqrt{VMU}) = \sum_i^{\min} \lambda(\sqrt{V \cdot M \cdot V})$$

"This norm describes the minimal form of a matrix."

iii)  $\ell_2$ -operator norm:

$$\|VMU\| = \sqrt{|\sum V_i \sum M_i \sum U_i|^2}$$

$$\leq \sqrt{\text{Tr}(V_i)^2 \cdot \text{Tr}(M_i)^2 \cdot \text{Tr}(U_i)^2}$$

$$\leq \sqrt{\text{Tr}(M_i)^2} \quad \text{when } V_i \text{ and also } U_i \text{ are orthonormal}$$

$$\leq \|M\|_2^2$$

"This is a special case of the  $\ell_p$ -norms."

iv)  $\ell_\infty$ -operator norm:

$$\|VMU\|_\infty = \sqrt[2]{|\sum V_i \sum M_i \sum U_i|^\infty}$$

$$\neq \|M\|_\infty^2$$

Other Method:

... uses:

$$V^T M U = [V_1, \dots, V_{d_1}] \begin{bmatrix} \lambda_1(M) & & 0 \\ & \ddots & \\ 0 & & \lambda_{d_1 \times d_2}(M) \end{bmatrix} \begin{bmatrix} U_1 \\ \vdots \\ U_{d_2} \end{bmatrix}$$

$$\text{where } V^T V = [V_1, \dots, V_{d_1}] \begin{bmatrix} V_1 \\ \vdots \\ V_{d_1} \end{bmatrix} = I$$

$$U^T U = [U_1, \dots, U_{d_2}] \begin{bmatrix} U_1 \\ \vdots \\ U_{d_2} \end{bmatrix}$$

e.g. the norms work with sums, traces, or matrices.



## (Symmetric Gauge Function)

$$\rho(x_1, \dots, x_{d_1}) = \rho(z_1 x_{\pi(1)}, \dots, z_{d_1} x_{\pi(d_1)}) \text{ for all binary strings } z \in \{-1, 1\}^{d_1}$$

and permutations,  $\pi$   
on  $\{1, \dots, d_1\}$

$$\begin{aligned} \text{b) } \rho(\text{III M III}) &= \rho(\text{Tr}(\sqrt{M})) \\ &= \rho\left(\sum_i^{d_1} \gamma_i(M)\right) \\ &= \rho(\gamma_1(M), \dots, \gamma_{d_1}(M)) \\ &= \rho(\gamma(M)) \in \mathbb{R}^{d_1} \end{aligned}$$

c) Frobenius Norm:

$$\rho(\text{III VMU III}_F) = \rho(\sqrt{|\sum_i^{d_1} v_i m_i u_i|^2}) \leq \rho(\sqrt{|\sum_i^{d_1} v_i|^2 \sum_i^{d_1} m_i \sum_i^{d_1} u_i|^2}) \leq \rho(\text{III M III})$$

Nuclear Norm:

$$\rho(\text{III VMU III}_{\text{Nuc}}) = \rho(\text{Tr}(\sqrt{VMU})) \leq \rho(\text{Tr}(\sqrt{V}) \text{Tr}(\sqrt{M}) \text{Tr}(\sqrt{U})) \leq \rho(\text{III M III})$$

$l_2$ -norm:

$$\rho(\text{III VMU III}_F) = \rho(\sqrt{|\sum_i^{d_1} v_i m_i u_i|^2}) \leq \rho(\sqrt{|\sum_i^{d_1} v_i|^2 \sum_i^{d_1} m_i \sum_i^{d_1} u_i|^2}) \leq \rho(\text{III M III})$$

$l_\infty$ -norm:

$$\rho(\text{III VMU III}_\infty) = \rho\left(\sqrt[{\frac{d_1 \times d_2}{d_1}}]{|\sum_i^{d_1} v_i m_i u_i|^\infty}\right) \leq \rho\left(\sqrt[{\frac{d_1 \times d_2}{d_1}}]{|\sum_i^{d_1} v_i|^\infty \sum_i^{d_1} m_i \sum_i^{d_2} u_i^\infty}\right) \leq \rho(\text{III M III})$$

"The norms example a distance measure to a  
nearest matrix, a minimal matrix, or a Cartesian  
matrix"

## (Weyl's Inequality)

$$\max_{j=1 \dots d} |\gamma_j(Q) - \gamma_j(R)| \leq \|Q - R\|_2$$

$$\text{When } Q = R + P$$

8.3 Proof by Example:  $Q = R + P$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1.01 \end{bmatrix} + \epsilon \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1.01 \end{bmatrix}$$

Eigenvalues for each matrix are:

$$\lambda(Q) = \left\{ \frac{1}{2}((a+1) \pm \sqrt{(a-1)^2 + 4\epsilon^2}) \right\} \quad \text{where } a=1.01$$

$$\lambda(R) = 1, a$$

$$\lambda(P) = \pm \epsilon$$

$$\begin{aligned} \max_{j=1,2} |\lambda(Q) - \lambda(R)| &= \left| \frac{1}{2}((a-1) - \sqrt{(a-1)^2 + 4\epsilon^2}) - a \right| \\ &= \frac{1}{2} |(1-a) - \sqrt{(a-1)^2 + 4\epsilon^2}| \\ &\leq \epsilon \quad \text{when } a > \epsilon \end{aligned}$$

8.4.  $V \in \mathbb{R}^{d \times n}$

$$\begin{aligned} \mathbb{E} \|V^T X\|_2^2 &= \mathbb{E} \left[ \sqrt{\left( \sum_i^{d \times n} V_i X_i \right)^2} \right] \\ &= \mathbb{E} \left[ \left( \sum_i^{d \times n} V_i X_i \right)^2 \right] \\ &= \sum_i^{d \times n} \mathbb{E} [ (V_i X_i)^2 ] \\ &= \sum_i^{d \times n} \mathbb{E} [ \langle V, X \rangle^2 ] \end{aligned}$$

8.5.  $\Theta \in S^{d \times d}$

$$\Theta \in \mathbb{R}^d : \{ \Theta^T \}_{t=0}^{\infty} ; \Theta^{t+1} = \frac{Q \Theta^t}{\|Q \Theta^t\|_2}$$

Power Iteration: The largest eigenvector is a good representation for a dataset.





Successive eigenvalues appear from subtracting the "largest" eigenvalue from  $Q$ :

$$Q_1 = \begin{bmatrix} 2-12 \\ 1-5 \end{bmatrix}; \quad Q_2 = \begin{bmatrix} 2-\lambda_1 & -12-\lambda_1 \\ 1-\lambda_1 & -5-\lambda_1 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 10 \\ 3 & 3 \end{bmatrix}$$

Similar steps follow from part a,  $\theta^1 = Q_2 \theta^0$   
 $\theta^2 = Q_2 \theta^1$

$$\theta^3 = Q_2 \theta^2$$

$\vdots$

$$\theta^d = Q_2 \theta^{d-1}$$

$$\text{with } \lambda_2 = \frac{Q_2 \theta^* \theta^*}{\theta^t \cdot \theta^t}$$

### 0.6 Gaussian Mixture Model:

$$f(x|\theta, \sigma^2 I_d) = \kappa \phi(x|\theta^*, \sigma^2 I_d) + (1-\kappa) \phi(x|\theta^*, \sigma^2 I_d)$$

where  $\phi(x|\theta, \sigma^2 I_d)$  is a Gaussian

$\kappa \in (0,1)$ ;  $\sigma > 0$ ;  $\theta^* = \text{mean}$

a) From corollary 6.20,  $\frac{\|\hat{\Sigma} - I\|}{\sqrt{n}} \leq 1 + c' \sqrt{\frac{d}{n}}$  for Gaussians with  $d$  dimensions and  $n > d$

$$\frac{\|\hat{\theta} - \theta\|}{\sqrt{n}} \leq 1 + c' \sqrt{\frac{d}{n}} \text{ where } n > c, \sigma^2(1+\sigma^2)d$$

$$\|\hat{\theta} - \theta\| \leq \sqrt{n} (1 + c') \sqrt{\frac{d}{n}}$$

$$\leq \sqrt{c_1} \sigma (1 + \sigma^2) (1 + c') \sqrt{\frac{d}{n}}$$

$$\leq c_2 \sigma (1 + \sigma^2) \sqrt{\frac{d}{n}} \text{ when } c_2 = (\sqrt{\frac{n}{d}} + c') \sqrt{c_1}$$



b) A classifier separates at a boundary, in this case at  $\hat{\theta}$ .

$$X_i = \begin{cases} -1 & \theta < \hat{\theta} \\ 1 & \theta > \hat{\theta} \end{cases}$$

If  $X_i = -1$ ,  $\varphi(X_i) = -\theta^*$

else  $X_i = 1$ ,  $\varphi(X_i) = \theta^*$

c) A Gaussian mixture model with zero variance has no covariance identity multiple.

$$\begin{aligned} \Gamma &= \mathbb{E}[X \otimes X] \\ &= \theta^* \otimes \theta^* \\ &\quad + \theta^* \otimes \theta^* + \sigma^2 I_d \end{aligned}$$

Other mixture models exist without identity multiples, as with Binomials, Gammas, and Poissons. Their repetitive moments never contain covariance.

0.7  $\theta^* \in \mathbb{R}^d$ ;  $n = \# \text{ samples}$ ;  $\{X_i, y_i\}_{i=1}^n = \{y_i + X_i\}_{i=1}^n$  where  $y_i = \langle X_i, \theta^* \rangle$

$= \sum_{i=1}^n \theta_i^* X_i + X_i$  where  $X_i \sim N(0, \sigma^2)$

$Z = \sum_{i=1}^n \theta_i^* X_i + X_i$

$\mathbb{E}[Z] = \mathbb{E}[Z \otimes Z]$

$= \mathbb{E}\left[\left(\sum_{i=1}^n \theta_i^* X_i + X_i\right)\left(\sum_{i=1}^n \theta_i^* X_i + X_i\right)\right]$

$= \mathbb{E}\left[\sum_{i=1}^n \theta_i^* X_i \cdot \sum_{i=1}^n \theta_i^* X_i\right] + 2 \mathbb{E}\left[\left(\sum_{i=1}^n \theta_i^* X_i\right) X_i\right] + \mathbb{E}[X_i \cdot X_i]$

$= \theta_i \otimes \theta_i \cdot \underbrace{\mathbb{E}\left[\sum_{i=1}^n X_i \sum_{i=1}^n X_i\right]}_{\approx 1.004} + \underbrace{2 \cdot \theta \mathbb{E}\left[\sum X_i\right] \mathbb{E}[X_i]}_{\text{mean} = 0} + \underbrace{\mathbb{E}[X_i]^2}_{\text{Variance}}$

$= \theta_i \otimes \theta_i + \sigma^2 I_d$

8.8 Equation 8.25a:  $\hat{\Theta} = \underset{\|\Theta\|_1=1}{\operatorname{argmin}} \{ \langle \Theta, \hat{\Sigma} \Theta \rangle \}$

a) Equation 8.25b "does" a better job at the relationship between Scotlass's equation and a function/constraint system.

$$\hat{\Theta} = \underset{\|\Theta\|_1=1}{\operatorname{argmax}} \langle \Theta, \hat{\Sigma} \Theta \rangle - \lambda \|\Theta\|_1$$

= 0

Form	Function	Bound	Type of Function	Number of Minima
Equation	$\lambda = \max_{\Theta} \hat{\Sigma} \Theta$	$\ \Theta\ _1 = 1$	"convex"	1
Function/ Constraint	$\max_{\Theta \in S} \operatorname{tr}(\hat{\Sigma} \Theta)$	$\operatorname{trace}(\Theta) = 1$	$\sum_{j,k=1}^d  \Theta_{j,k}  \leq R^2$	$\operatorname{rank}(\Theta) = 1$

b) A rank constraint for a convex function defines one minimum with rank of one. A function/constraint without rank one satisfies multiple minima and a non-convex type.

8.9  $\hat{\Theta} = \underset{\substack{\Theta \in S \\ \operatorname{tr}(\Theta)=1}}{\operatorname{max}} \left\{ \operatorname{trace}(\hat{\Sigma} \Theta) - \lambda_n \sum_{j,k=1}^d |\Theta_{j,k}| \right\}$

$$\hat{U} = \begin{cases} \operatorname{sign}(\hat{\Theta}_{j,k}) & \text{if } \hat{\Theta}_{j,k} \neq 0 \\ \in [-1, 1] & \text{otherwise} \end{cases}$$

$$\hat{\Theta} = \max_{\Theta \in S} \left\{ \operatorname{trace}(\hat{\Sigma} \Theta) - \lambda_n \sum_{j,k=1}^d |\Theta_{j,k}| \right\} = 0$$

$$= \max_{\Theta \in S} \left\{ \operatorname{trace}(\hat{\Sigma} \Theta)^2 - \lambda \right\}$$

$$= \max_{U \in S} \max_{\Theta \in S} \hat{U} \left\{ \operatorname{trace}(\hat{\Sigma} \Theta)^2 - \lambda \right\}$$

$$= \max_{U \in S} \hat{U} \left\{ \operatorname{trace}(\hat{\Sigma})^2 \right\} \Rightarrow \hat{\Theta} = \operatorname{trace}(\hat{\Sigma} \Theta)^2 = \operatorname{trace}(\hat{\Sigma})^2 = 0$$

=  $\Theta^T \Theta$