

3. *Nystrom method*: In class, we have shown that every manifold learning algorithm can be regarded as Kernel PCA on graphs: (1) given  $N$  data points, define a neighborhood graph with  $N$  nodes for data points; (2) construct a positive semidefinite kernel  $K$ ; (3) pursue spectral decomposition of  $K$  to find the embedding (using top or bottom eigenvectors). However, this approach might suffer from the expensive computational cost in spectral decomposition of  $K$  if  $N$  is large and  $K$  is non-sparse, e.g. ISOMAP and MDS.

To overcome this hurdle, Nystrom method leads us to a scalable approach to compute eigenvectors of low rank matrices. Suppose that an  $N$ -by- $N$  positive semidefinite matrix  $K \succeq 0$  admits the following block partition

$$K = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}. \quad (1)$$

where  $A$  is an  $n$ -by- $n$  block. Assume that  $A$  has the spectral decomposition  $A = U\Lambda U^T$ ,  $\Lambda = \text{diag}(\lambda_i)$  ( $\lambda_1 \geq \lambda_2 \geq \dots \lambda_k > \lambda_{k+1} = \dots = 0$ ) and  $U = [u_1, \dots, u_n]$  satisfies  $U^T U = I$ .

- (a) Assume that  $K = XX^T$  for some  $X = [X_1; X_2] \in \mathbb{R}^{N \times k}$  with the block  $X_1 \in \mathbb{R}^{n \times k}$ . Show that  $X_1$  and  $X_2$  can be decided by:

$$X_1 = U_k \Lambda_k^{1/2}, \quad (2)$$

$$X_2 = B^T U_k \Lambda_k^{-1/2}, \quad (3)$$

where  $U_k = [u_1, \dots, u_k]$  consists of those  $k$  columns of  $U$  corresponding to top  $k$  eigenvalues  $\lambda_i$  ( $i = 1, \dots, k$ ).

- (b) Show that for general  $K \succeq 0$ , one can construct an approximation from (2) and (3),

$$\hat{K} = \begin{bmatrix} A & B \\ B^T & \hat{C} \end{bmatrix}. \quad (4)$$

where  $A = X_1 X_1^T$ ,  $B = X_1 X_2^T$ , and  $\hat{C} = X_2 X_2^T = B^T A^\dagger B$ ,  $A^\dagger$  denoting the Moore-Penrose (pseudo-) inverse of  $A$ . Therefore  $\|\hat{K} - K\|_F = \|C - B^T A^\dagger B\|_F$ . Here the matrix  $C - B^T A^\dagger B =: K/A$  is called the (generalized) *Schur Complement* of  $A$  in  $K$ .

$$(1) \quad A = U \Lambda U^T, \quad \Lambda = \text{diag}(\lambda_i)$$

$$\Rightarrow A = U_k \Lambda_k U_k^T.$$

where  $U_k = [u_1, \dots, u_k]$  consists of top- $k$  eigenvalues

$$\Lambda_k = \text{diag}(\lambda_1, \dots, \lambda_k) \quad (\lambda_i \text{ non-zero})$$

$$X = (X_1 \ X_2)^T$$

$$K = X X^T = \begin{pmatrix} X_1 X_1^T & X_1 X_2^T \\ X_2 X_1^T & X_2 X_2^T \end{pmatrix}$$

$$\textcircled{1} \quad X_1 X_1^T = A = U_k \Lambda_k U_k^T$$

$$X_1 X_1^T = U_k \Lambda_k U_k^T$$

$$= U_k \Lambda_k^{-\frac{1}{2}} \Lambda_k^{\frac{1}{2}} U_k^T$$

$$= (U_k \Lambda_k^{\frac{1}{2}}) (U_k \Lambda_k^{\frac{1}{2}})^T$$

$$\text{Def } X_1 \triangleq U_k \Lambda_k^{\frac{1}{2}}. \quad \text{Then } A = X_1 X_1^T$$

$$\textcircled{2} X_1 X_2^T = B$$

$$X_2^T = X_1^T B = (\Lambda_k^{\frac{1}{2}})^T U_k^T B$$

$$= \Lambda_k^{-\frac{1}{2}} U_k^T B$$

$$X_2 = B^T U_k \Lambda_k^{-\frac{1}{2}}$$

$$\textcircled{2} X_1 X_1^T = A$$

$$X_1 X_2^T = B$$

$$X_2 X_1^T = (X_1 X_2^T)^T = B^T$$

$$X_2 X_2^T = B^T U_k \Lambda_k^{-1} U_k^T B$$

$$\text{Def } G \triangleq U_k \Lambda_k^{-1} U_k^T. \quad \text{Then}$$

$$\textcircled{1} G A G = (U_k \Lambda_k^{-1} U_k^T) (U_k \Lambda_k U_k^T) (U_k \Lambda_k^{-1} U_k^T)$$

$$= U_k \Lambda_k^{-1} U_k^T$$

$$= G$$

$$\textcircled{2} A G A = (U_k \Lambda_k U_k^T) (U_k \Lambda_k^{-1} U_k^T) (U_k \Lambda_k U_k^T)$$

$$= U_k \Lambda_k U_k^T$$

$$= A$$

$$\textcircled{3} (A G)^T = (U_k \Lambda_k U_k^T U_k \Lambda_k^{-1} U_k^T)^T = I = A G$$

By definition,  $G$  is the unique Moore-Penrose inverse of  $A$ .  $G = A^+$

Therefore,

$$\|K - \hat{K}\|_F = \sqrt{\|A - \hat{A}\|_F^2 + \|B - \hat{B}\|_F^2 + \|B^T - \hat{B}^T\|_F^2 + \|C - \hat{C}\|_F^2}$$

$$K = \begin{bmatrix} A & B \\ B^T & \hat{C} \end{bmatrix}$$

