

## Supplementary material for “Factor Analysis on a Graph”

### A Proof for Theorem 1

We discuss the relation between the graph connectivity and our kernel  $\widehat{\Sigma}$ , by using covariance decomposition of Jones and West (2005), which was originally proposed for analyzing paths on a graphical model. The  $(i, j)$  element of the covariance matrix can be decomposed as a weighted sum of products of conditional correlations of consecutive node pairs on all possible paths between  $i$  and  $j$ .

**Theorem 4** (Jones and West (2005)). *Let  $\mathcal{P}_{ij}$  be a set of paths between nodes  $i$  and  $j$  on the graph. A path  $\mathcal{P} \in \mathcal{P}_{ij}$  is defined by a set of nodes ordered from  $i$  to  $j$ , i.e.,  $\mathcal{P} := \{(p_1, \dots, p_m) | p_1 = i, p_m = j, m \leq d\}$ . We then have*

$$\Sigma_{ij} = (-1)^{m+1} \Theta_{p_1, p_2} \Theta_{p_2, p_3} \cdots \Theta_{p_{m-1}, p_m} \frac{\det(\Theta_{\setminus \mathcal{P}})}{\det(\Theta)} \quad (10)$$

According to the decomposition (10) and  $\det(\Sigma_{\mathcal{P}}) = \det(\Theta_{\setminus \mathcal{P}}) / \det(\Theta)$ , we obtain Theorem 1 in the main text.

### B Optimality Condition of Factor Loading Matrix

The optimality condition of factor loading matrix  $\mathbf{A}$  is as follows:

**Lemma 1** (e.g., Jöreskog 1967). *Assuming that we already have  $\widehat{\Psi}$ , defined as the maximum likelihood estimate for  $\Psi$ , then the maximum likelihood solution for  $\mathbf{A}$  satisfies the following equation:*

$$(\widehat{\Psi}^{-1/2} \widehat{\Sigma} \widehat{\Psi}^{-1/2}) (\widehat{\Psi}^{-1/2} \mathbf{A}) = (\widehat{\Psi}^{-1/2} \mathbf{A}) (\mathbf{I} + \mathbf{A}^\top \widehat{\Psi}^{-1} \mathbf{A}).$$

Suppose that  $\mathbf{A}^\top \widehat{\Psi}^{-1} \mathbf{A}$  is a diagonal matrix (This can be achieved by post-multiplying  $\mathbf{A}$  by an orthogonal matrix, which does not affect the value of the likelihood), the equation can be regarded as an eigenvalue decomposition by which we obtain the estimator  $\widehat{\mathbf{A}}$  for  $\mathbf{A}$  as follows:

$$\widehat{\mathbf{A}} := \widehat{\Psi}^{1/2} \mathbf{U}_k (\mathbf{\Lambda}_k - \mathbf{I})^{1/2}, \quad (11)$$

where  $\mathbf{\Lambda}_k := \text{diag}(\lambda_1, \dots, \lambda_k)$  for the first  $k$  largest eigenvalues  $\lambda_1, \dots, \lambda_k$  of  $\widehat{\Psi}^{-1/2} \widehat{\Sigma} \widehat{\Psi}^{-1/2}$ , and  $\mathbf{U}_k \in \mathbb{R}^{d \times k}$  is a set of the corresponding eigenvectors.

We now derive the first order condition of  $\mathbf{A}$  in the above lemma. The derivative of the objective function of factor analysis in the main text in terms of  $\mathbf{A}$  is

$$\begin{aligned} & 2(\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A} - (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \widehat{\Sigma} (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A} \\ & = 2(\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} (\mathbf{A}\mathbf{A}^\top + \Psi - \widehat{\Sigma}) (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A}. \end{aligned}$$

The first order condition is written as

$$(\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} (\mathbf{A}\mathbf{A}^\top + \Psi - \widehat{\Sigma}) (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A} = \mathbf{0}.$$

Multiplying this equation by  $(\mathbf{A}\mathbf{A}^\top + \Psi)$  from the left, we obtain

$$(\mathbf{A}\mathbf{A}^\top + \Psi - \widehat{\Sigma}) (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A} = \mathbf{0}.$$

Using Woodbury formula, we see

$$\begin{aligned} & (\mathbf{A}\mathbf{A}^\top + \Psi)^{-1} \mathbf{A} \\ & = (\Psi^{-1} - \Psi^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A}^\top \Psi^{-1} \mathbf{A})^{-1} \mathbf{A}^\top \Psi^{-1}) \mathbf{A} \\ & = \Psi^{-1} \mathbf{A} (\mathbf{I} - (\mathbf{I} + \mathbf{A}^\top \Psi^{-1} \mathbf{A})^{-1} \mathbf{A}^\top \Psi^{-1} \mathbf{A}) \\ & = \Psi^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A}^\top \Psi^{-1} \mathbf{A})^{-1} (\mathbf{I} + \mathbf{A}^\top \Psi^{-1} \mathbf{A} - \mathbf{A}^\top \Psi^{-1} \mathbf{A}) \\ & = \Psi^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A}^\top \Psi^{-1} \mathbf{A})^{-1}. \end{aligned}$$

By this transformation, the first order condition can be written as

$$\begin{aligned} (\mathbf{A}\mathbf{A}^\top + \Psi - \widehat{\Sigma})\Psi^{-1}\mathbf{A}(\mathbf{I} + \mathbf{A}^\top\Psi^{-1}\mathbf{A})^{-1} &= \mathbf{0} \\ (\mathbf{A}\mathbf{A}^\top + \Psi - \widehat{\Sigma})\Psi^{-1}\mathbf{A} &= \mathbf{0} \\ \mathbf{A}\mathbf{A}^\top\Psi^{-1}\mathbf{A} + \mathbf{A} - \widehat{\Sigma}\Psi^{-1}\mathbf{A} &= \mathbf{0}. \end{aligned}$$

To derive the second equation in the above, we multiplied through by  $(\mathbf{I} + \mathbf{A}^\top\Psi^{-1}\mathbf{A})$  from the right. Arranging this equation, we obtain

$$\widehat{\Sigma}\Psi^{-1}\mathbf{A} = \mathbf{A}(\mathbf{I} + \mathbf{A}^\top\Psi^{-1}\mathbf{A}^\top).$$

Multiplying this equation by  $\Psi^{-1/2}$  from the left, we finally see

$$(\Psi^{-1/2}\widehat{\Sigma}\Psi^{-1/2})(\Psi^{-1/2}\mathbf{A}) = (\Psi^{-1/2}\mathbf{A})(\mathbf{I} + \mathbf{A}^\top\Psi^{-1}\mathbf{A}^\top).$$

## C Spectral Relaxation of Weighted Kernel $k$ -means

Then the objective function of weighted kernel  $k$ -means is defined by

$$\sum_{i=1}^k \sum_{j \in \mathcal{C}_i} \widehat{\psi}_j^{-1} \|\phi_j - \mu_i\|_2^2, \quad (12)$$

where  $\mathcal{C}_i$  for  $i = 1, \dots, k$  is an index set of the  $i$ -th cluster, and  $\mu_i$  is a centroid of the  $i$ -th cluster. In this objective function (12), the squared error between each  $\phi_i$  and its centroid is weighted by  $\widehat{\psi}_i^{-1}$ , which means that if the  $i$ -th dimension of the factor analysis error term  $\epsilon$  has a smaller variance, a corresponding  $\phi_i$  is penalized more strongly. Using an indicator matrix  $\mathbf{Z}$ , in which the  $(i, j)$  element takes 1 if the  $i$ -th instance belongs to the  $j$ -th cluster or takes 0 otherwise, this function can be re-written as:

$$\text{trace} \left\{ (\Phi - \mathbf{Z}\mathbf{M}^\top)^\top \widehat{\Psi}^{-1} (\Phi - \mathbf{Z}\mathbf{M}^\top) \right\}, \quad (13)$$

where  $\mathbf{M} := [\mu_1, \dots, \mu_k]$ .

We consider a spectral relaxation of this weighted kernel  $k$ -means. Given a cluster assignment  $\mathcal{C}_i$ , the centroid which minimizes the squared error is the weighted average of the instances:  $\sum_{j \in \mathcal{C}_i} \widehat{\psi}_j^{-1} \phi_j / \sum_{j \in \mathcal{C}_i} \widehat{\psi}_j^{-1}$ . Then, the set of centroids can be written as

$$\mathbf{M} = \Phi^\top \widehat{\Psi}^{-1} \mathbf{Z}\mathbf{C},$$

where  $\mathbf{C} = \text{diag}(1/\sum_{j \in \mathcal{C}_1} \widehat{\psi}_j^{-1}, \dots, 1/\sum_{j \in \mathcal{C}_k} \widehat{\psi}_j^{-1})$ . Substituting this into (13), the objective function can be transformed into

$$\begin{aligned} &\text{trace} \left\{ (\Phi - \mathbf{Z}\mathbf{M}^\top)^\top \widehat{\Psi}^{-1} (\Phi - \mathbf{Z}\mathbf{M}^\top) \right\} \\ &= \text{trace} \left( \Phi^\top \widehat{\Psi}^{-1} \Phi - \Phi^\top \widehat{\Psi}^{-1} \mathbf{Z}\mathbf{C}\mathbf{Z}^\top \widehat{\Psi}^{-1} \Phi \right). \\ &= \text{trace} \left( \Phi^\top \widehat{\Psi}^{-1} \Phi - \mathbf{C}^{1/2} \mathbf{Z}^\top \widehat{\Psi}^{-1} \Phi \Phi^\top \widehat{\Psi}^{-1} \mathbf{Z}\mathbf{C}^{1/2} \right). \end{aligned}$$

Here, we used  $\mathbf{Z}^\top \widehat{\Psi}^{-1} \mathbf{Z} = \mathbf{C}^{-1}$ , and the first term is now constant. Defining  $\mathbf{V}_k := \widehat{\Psi}^{-1/2} \mathbf{Z}\mathbf{C}^{1/2}$ , which leads  $\mathbf{V}_k^\top \mathbf{V}_k = \mathbf{C}^{1/2} \mathbf{Z}^\top \widehat{\Psi}^{-1} \mathbf{Z}\mathbf{C}^{1/2} = \mathbf{C}^{1/2} \mathbf{C}^{-1} \mathbf{C}^{1/2} = \mathbf{I}$ , the following spectral relaxation of the weighted kernel  $k$ -means can be derived:

$$\begin{aligned} &\max_{\mathbf{V}_k \in \mathbb{R}^{d \times k}} \text{trace} \left( \mathbf{V}_k^\top \widehat{\Psi}^{-1/2} \widehat{\Sigma} \widehat{\Psi}^{-1/2} \mathbf{V}_k \right) \\ &\text{s.t.} \quad \mathbf{V}_k^\top \mathbf{V}_k = \mathbf{I}. \end{aligned} \quad (14)$$

## D Proof for Theorem 2

Theorem 2 can be derived from the optimality condition for the factor loading matrix  $\mathbf{A}$  written in supplementary appendix B.

Since the set of eigenvectors corresponding to the  $k$  largest eigenvalues of  $\widehat{\Psi}^{-1/2} \widehat{\Sigma} \widehat{\Psi}^{-1/2}$  is an optimal solution to (14), we see  $\mathbf{V}_k = \mathbf{U}_k$ . We therefore obtain the relation  $\widehat{\mathbf{Z}} = \widehat{\mathbf{A}}\mathbf{D}$ , where  $\mathbf{D} := \mathbf{C}^{-1/2}(\mathbf{\Lambda}_k - \mathbf{I})^{-1/2}$ , which is a diagonal matrix.

## E Proof for Theorem 3

Replacing  $\mathbf{V}_k$  in  $\widehat{\mathbf{Z}}$  (5) (written the main text) by  $\mathbf{V}_k\mathbf{Q}$  keeps the objective of kernel  $k$ -means (14) optimal, and we obtain  $\widehat{\mathbf{Z}} = \widehat{\Psi}^{1/2} \mathbf{V}_k \mathbf{Q} \mathbf{C}^{-1/2}$ . Then, we see  $\widehat{\mathbf{Z}}\mathbf{C}^{1/2} = \widehat{\mathbf{A}}(\mathbf{\Lambda}_k - \mathbf{I})^{-1/2}\mathbf{Q}$  (Note that  $\mathbf{C}$  is diagonal). The invariance of the likelihood can be easily seen by  $\widehat{\mathbf{A}}_{\text{rot}}(\mathbf{Q}^\top(\mathbf{\Lambda}_k - \mathbf{I})\mathbf{Q})\widehat{\mathbf{A}}_{\text{rot}}^\top = \widehat{\mathbf{A}}\widehat{\mathbf{A}}^\top$ .

## F Formulation of Lap-PCA

Let  $\mathbf{W}$  be an adjacency matrix of the graph  $\mathcal{G}$  in which the  $(i, j)$  element is  $W_{ij} = 1$  if  $(i, j) \in \mathcal{E}$ , and  $W_{ij} = 0$  otherwise.

For factor analysis and PCA, In our case, the graph structure can be incorporated into the matrix  $\mathbf{A}$  by the following formulation: in factor analysis or PCA by

$$\min_{\mathbf{A} \in \mathbb{R}^{d \times k}, \Psi \in \mathcal{D}_+^d} (1 - \alpha) \ell(\mathbf{A}, \Psi) + \alpha \sum_{k'=1}^k \sum_{(i,j) \in \mathcal{E}} W_{ij} (A_{ik'} - A_{jk'})^2, \quad (15)$$

where  $\ell$  is a loss function (negative log-likelihood),  $\mathbf{L} \in \mathbb{R}^{d \times d}$  is the graph Laplacian matrix (see, e.g., Chung, 1997, for detail), and  $\alpha \in [0, 1]$  is a regularization parameter. For PCA,  $\Psi$  has an additional constraint  $\Psi = \sigma^2 \mathbf{I}$ .

In experiments, we chose the best regularization parameter  $\alpha$  in (15) out of  $\{0.25, 0.5, 0.75\}$  in terms of each result.

## References

- M. Ashburner. Gene ontology: Tool for the unification of biology. *Nature Genetics*, 25:25–29, 2000.
- M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15(6):1373–1396, 2003.
- M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434, 2006.
- D. M. Boyd and N. B. Ellison. Social network sites: Definition, history, and scholarship. *Journal of Computer-Mediated Communication*, 13(1):210–230, 2007.
- E. G. Cerami and et al. Pathway commons, a web resource for biological pathway data. *Nucleic Acids Research*, 39:685–690, 2011.
- F. R. K. Chung. *Spectral Graph Theory*. American Mathematical Society, 1997.
- R. Cohen and S. Havlin. *Complex networks: Structure, Robustness and Function*. Cambridge University Press, 2010.
- I. S. Dhillon, Y. Guan, and B. Kulis. Kernel  $k$ -means: Spectral clustering and normalized cuts. In *Proc. of the 10th ACM SIGKDD*, pages 551–556. ACM, 2004.
- A. B. Goldberg, X. Zhu, and S. J. Wright. Dissimilarity in graph-based semi-supervised classification. In M. Meila and X. Shen, editors, *Proceedings of the 11th International Conference on Artificial Intelligence and Statistics*, volume 2, pages 155–162. JMLR.org, 2007.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. J. Smola. A kernel method for the two-sample problem. In *Advances in NIPS 19*, pages 513–520. MIT Press, 2007.

- H. H. Harman. *Modern Factor Analysis*. The university of chicago press, 1960.
- T. Ideker, O. Ozier, B. Schwikowski, and A. F. Siegel. Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics*, 18(suppl 1):S233–S240, 2002.
- B. Jiang, C. Ding, B. Luo, and J. Tang. Graph-Laplacian PCA: Closed-form solution and robustness. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3492–3498, 2013.
- I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, 2002.
- B. Jones and M. West. Covariance decomposition in undirected gaussian graphical models. *Biometrika*, 92(4): 779–786, 2005.
- K. G. Jöreskog. Some contributions to maximum likelihood factor analysis. *Psychometrika*, 32(4):443–482, 1967.
- H. F. Kaiser. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200, 1958.
- M. Kanehisa and S. Goto. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Research*, 28(1): 27–30, 2000.
- B. Karrer and M. E. J. Newman. Stochastic blockmodels and community structure in networks. *PHYSICAL REVIEW E*, 83:016107, Jan 2011.
- R. Kosala and H. Blockeel. Web mining research: A survey. *SIGKDD Explorations Newsletter*, 2(1):1–15, June 2000.
- S. Kullback and R. A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1): 79–86, 1951.
- C. Li and H. Li. Network-constrained regularization and variable selection for analysis of genomic data. *Bioinformatics*, 24(9):1175–1182, 2008.
- M. Meila and J. Shi. A random walks view of spectral segmentation. In *Proceedings of the 8th International Workshop on Artificial Intelligence and Statistics*. Morgan Kaufmann, 2001.
- M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582, 2006.
- A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in NIPS 14*, pages 849–856. MIT Press, 2001.
- J. S.-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge university press, 2004.
- M. Saerens, F. Fouss, L. Yen, and P. Dupont. The principal components analysis of a graph, and its relationships to spectral clustering. In *Proc. of the 15th ECML*, volume 3201 of *LNCS*, pages 371–383. Springer Berlin Heidelberg, 2004.
- G. Sales, E. Calura, D. Cavalieri, and C. Romualdi. graphite - a bioconductor package to convert pathway topology to gene network. *BMC Bioinformatics*, 13(1):20, 2012.
- T. Sandler, J. Blitzer, P. P. Talukdar, and L. H. Ungar. Regularized learning with networks of features. In *Advances in NIPS 21*, pages 1401–1408, 2008.
- S. E. Schaeffer. Survey: Graph clustering. *Computer Science Review*, 1(1):27–64, Aug. 2007.
- J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- A. J. Smola and I. R. Kondor. Kernels and regularization on graphs. In *Proceedings of the Annual Conference on Computational Learning Theory*, pages 144–158, 2003.
- R. Tibshirani and J. Taylor. The solution path of the generalized lasso. *The Annals of Statistics*, 39(3):1335–1371, 2011.
- M. J. van de Vijver and et al. A gene-expression signature as a predictor of survival in breast cancer. *New England Journal of Medicine*, 347(25):1999–2009, 2002.
- U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- J. Whittaker. *Graphical Models in Applied Multivariate Statistics*. Wiley Publishing, 1990.
- L. Wu, X. Ying, X. Wu, A. Lu, and Z.-H. Zhou. Spectral analysis of  $k$ -balanced signed graphs. In *Advances in Knowledge Discovery and Data Mining - 15th Pacific-Asia Conference (PAKDD)*, volume 6635 of *Lecture Notes in Computer Science*, pages 1–12. Springer, 2011.

- S. Yang, L. Yuan, Y.-C. Lai, X. Shen, P. Wonka, and J. Ye. Feature grouping and selection over an undirected graph. In *Proc. of the 18th ACM SIGKDD*, pages 922–930. ACM, 2012.
- H. Zha, X. He, C. Ding, H. Simon, and M. Gu. Spectral relaxation for k-means clustering. In *Advances in NIPS 14*, pages 1057–1064. MIT Press, 2001.
- M. Zheng, J. Bu, C. Chen, C. Wang, L. Zhang, G. Qiu, and D. Cai. Graph regularized sparse coding for image representation. *Image Processing, IEEE Transactions on*, 20(5):1327–1336, 2011.
- D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf. Learning with local and global consistency. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems (NIPS) 16*. MIT Press, 2004.
- X. Zhu, Z. Ghahramani, and J. D. Lafferty. Semi-supervised learning using Gaussian fields and harmonic functions. In T. Fawcett and N. Mishra, editors, *Proceedings of the 20th International Conference on Machine Learning (ICML)*, pages 912–919. AAAI Press, 2003.