

# Supplementary Materials: Scalable Multi-view Unsupervised Feature Selection with Structure Learning and Fusion

Anonymous Authors

## 1 APPENDIX A. THE DETAILED OPTIMIZATION FOR $\alpha$

By fixing other variables except  $\alpha$ , we can optimize  $\alpha$  by rows:

$$\min_{\alpha_i \mathbf{1}=1, \alpha_i \geq 0} \|\mathbf{u}_i - \sum_{v=1}^V \alpha_i^v \tilde{\mathbf{u}}_i^v\|_2^2. \quad (1)$$

Eq. (1) can be equally transformed into the following form:

$$\begin{aligned} \|\mathbf{u}_i - \sum_{v=1}^V \alpha_i^v \tilde{\mathbf{u}}_i^v\|_2^2 &= \left\| \sum_{v=1}^V \alpha_i^v \mathbf{u}_i - \sum_{v=1}^V \alpha_i^v \tilde{\mathbf{u}}_i^v \right\|_2^2 \\ \left\| \sum_{v=1}^V \alpha_i^v (\mathbf{u}_i - \tilde{\mathbf{u}}_i^v) \right\|_2^2 &= \|\alpha_i \mathbf{D}_i\|_2^2 \end{aligned} \quad (2)$$

where  $\mathbf{D}_i = [\mathbf{d}_i^1; \dots; \mathbf{d}_i^V] \in \mathbb{R}^{V \times c}$ , and  $\mathbf{d}_i^v = \mathbf{u}_i - \tilde{\mathbf{u}}_i^v \in \mathbb{R}^{1 \times c}$ . Therefore,  $\alpha_i$  can be solved as:

$$\min_{\alpha_i \mathbf{1}=1, \alpha_i \geq 0} \alpha_i \mathbf{D}_i \mathbf{D}_i^T \alpha_i^T. \quad (3)$$

Since  $\mathbf{D}_i \mathbf{D}_i^T$  is semi-definite, Eq. (3) is a quadratic convex programming problem, which can be solved efficiently [2]. Specifically, Eq. (3) can be solved by tackling its counterpart:

$$\min_{\alpha_i \geq 0, \alpha_i \mathbf{1}=1, z} \alpha_i \mathbf{D}_i \mathbf{D}_i^T z^T + \frac{\mu}{2} \|\alpha_i - z + \frac{\tau}{\mu}\|_2^2, \quad (4)$$

where  $z \in \mathbb{R}^{1 \times V}$  denotes a slack variable,  $\mu > 0$  is a penalty parameter, and  $\tau \in \mathbb{R}^{1 \times V}$  is a Lagrangian multiplier. Eq. (4) can be iteratively optimized by the augmented Lagrangian multiplier method. The solution steps are as follows:

*Step 1. Update  $z$ :* When  $\alpha_i$  is fixed, Eq. (4) is an unconstrained optimization problem. By setting the derivative of Eq. (4) w.r.t.  $z$  to zero, we update  $z$  by:

$$z = \alpha_i - \frac{1}{\mu} (\alpha_i \mathbf{D}_i \mathbf{D}_i^T - \tau). \quad (5)$$

*Step 2. Update  $\alpha_i$ :* When  $z$  is fixed with its current value of  $z$  (i.e.,  $z^*$ ),  $\alpha_i$  can be updated by minimizing the following problem:

$$\min_{\alpha_i \mathbf{1}=1, \alpha_i \geq 0} \|\alpha_i - z^* + \frac{1}{\mu} (\tau + z^* \mathbf{D}_i \mathbf{D}_i^T)\|_2^2, \quad (6)$$

which can be solved with a closed-form solution [1].

*Step 3. Update  $\tau$  and  $\mu$ :* In each iteration, we update the Lagrange multipliers  $\tau$  and the penalty parameter  $\mu$  as follows:

$$\begin{aligned} \tau &= \tau + \mu (\alpha_i - z) \\ \mu &= \rho \mu. \end{aligned} \quad (7)$$

where  $\rho$  is a constant update rate. In this way,  $\alpha_i$  can be adaptively updated according to the aforementioned steps.

## 2 APPENDIX B. THE DETAILED OPTIMIZATION FOR $S$

By fixing other variables except  $S$ , we have the following problem:

$$\min_{S \mathbf{1}=1, S \geq 0} \lambda \|\mathbf{u}_i - \mathbf{u}_j\|_2^2 s_{ij} + \beta \|S\|_F^2. \quad (8)$$

Noting that each row of  $S$  (i.e.,  $s_i$ ) is uncorrelated with others, hence Eq. (8) can be optimized for each row independently as follows:

$$\min_{s_i \mathbf{1}=1, s_i \geq 0} \|s_i + \frac{1}{2\beta_i} \mathbf{d}_i\|_2^2, \quad (9)$$

where  $\mathbf{d}_i$  is a row vector with  $d_{ij} = \lambda \|\mathbf{u}_i - \mathbf{u}_j\|_2^2$ . The Lagrangian function of the above function is:

$$\mathcal{L}(s_i, \theta_i, \zeta_i) = \frac{1}{2} \|s_i + \frac{1}{2\beta_i} \mathbf{d}_i\|_2^2 - \theta_i (s_i \mathbf{1} - 1) - s_i \zeta_i,$$

where  $\theta_i \in \mathbb{R}$  and  $\zeta_i \in \mathbb{R}^{n \times 1}$  are Lagrangian multipliers. According to the KKT condition, the optimal solution of  $s_i$  is:

$$s_{ij} = \left( -\frac{d_{ij}}{2\beta_i} + \theta_i^* \right)_+,$$

where  $\theta_i^*$  denotes the optimal value equipped for the optimal solution of  $s_i$  and  $(x)_+ = \max(x, 0)$ . Since the local structure contains more useful and detailed information about the data compared to the global structure, it is suitable to construct a sparse graph to focus on a small number of neighbors [4]. With  $\mathbf{d}_i$  being sorted from small to large (i.e.,  $\tilde{\mathbf{d}}$ ), there is  $s_{i1} \geq s_{i2} \geq \dots \geq s_{in}$ . Assuming that each sample has  $f$ -nearest neighbors (i.e.,  $s_i$  has  $f$  nonzero elements), we derive:

$$\begin{cases} s_{i,f} > 0 \\ s_{i,f+1} = 0 \end{cases} \implies \begin{cases} -\frac{\tilde{d}_{i,f}}{2\beta_i} + \theta_i^* > 0 \\ -\frac{\tilde{d}_{i,f+1}}{2\beta_i} + \theta_i^* \leq 0 \end{cases}.$$

Due to  $s_i \mathbf{1} = 1$ , we obtain:

$$\sum_{j=1}^f \left( -\frac{\tilde{d}_{i,j}}{2\beta_i} + \theta_i^* \right) = 1 \implies \theta_i^* = \frac{1}{f} + \frac{1}{2f\beta_i} \sum_{j=1}^f \tilde{d}_{i,j}.$$

Based on the above analysis, the inequality on  $\beta_i$  can be derived as:

$$\frac{f}{2} \tilde{d}_{i,f} - \frac{1}{2} \sum_{j=1}^f \tilde{d}_{i,j} < \beta_i \leq \frac{f}{2} \tilde{d}_{i,f+1} - \frac{1}{2} \sum_{j=1}^f \tilde{d}_{i,j}.$$

when  $\beta_i = \frac{f}{2} \tilde{d}_{i,f+1} - \frac{1}{2} \sum_{j=1}^f \tilde{d}_{i,j}$ , it satisfies that  $s_{i,f+1} = 0$  and  $s_i$  has  $f$  nonzero elements exactly. With the optimal  $\beta_i$  and  $\theta_i^*$ , the solution of  $s_i$  is derived as:

$$s_{ij} = \left( \frac{\tilde{d}_{i,f+1} - \tilde{d}_{i,j}}{f \tilde{d}_{i,f+1} - \sum_{j=1}^f \tilde{d}_{i,j}} \right)_+.$$

According to [3],  $\beta = \sum_{i=1}^n \frac{f \tilde{d}_{i,f+1} - \sum_{j=1}^f \tilde{d}_{i,j}}{2n}$  is set to the mean of  $\beta_1, \beta_2, \dots, \beta_n$ .

## REFERENCES

- [1] Jin Huang, Feiping Nie, and Heng Huang. 2015. A new simplex sparse learning model to measure data similarity for clustering. In *International Joint Conference on Artificial Intelligence*. 3569–3575.
- [2] Bingbing Jiang, Chenglong Zhang, Yan Zhong, Yi Liu, Yingwei Zhang, Xingyu Wu, and Weiguo Sheng. 2023. Adaptive collaborative fusion for multi-view semi-supervised classification. *Information Fusion* 96 (2023), 37–50.
- [3] Feiping Nie, Xiaoqian Wang, and Heng Huang. 2014. Clustering and projected clustering with adaptive neighbors. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining*. 977–986.
- [4] Hebing Nie, Qun Wu, Haifeng Zhao, Weiping Ding, and Muhammet Deveci. 2023. Flexible Adaptive Graph Embedding for semi-supervised dimension reduction. *Information Fusion* (2023), 101872.