



Dependence-Guided Multi-view Clustering

Xia Dong^{1,3,4}, Danyang Wu^{1,3,4}, Feiping Nie^{1,3,4}, Rong Wang^{2,3,4}, and Xuelong Li^{1,3,4}, *Fellow, IEEE*

Poster Number: 3547



¹ School of Computer Science ² School of Cybersecurity

³ Center for OPTical IMagery Analysis and Learning (OPTIMAL) ⁴ Northwestern Polytechnical University (Nwpu), Xi'an, China
xiadongpgh@gmail.com, danyangwu41x@mail.nwpu.edu.cn, feipingnie@gmail.com, wangrong07@tsinghua.org.cn, li@nwpu.edu.cn

Problem Formulation

1. The dependence between unified embedding \mathbf{P} and embedding of each view $\mathbf{F}^{(v)}$ is measured by **Hilbert Schmidt Independence Criterion (HSIC)** [1]. 2. A cluster indicator matrix \mathbf{Y} is recovered from the unified embedding \mathbf{P} via estimating a rotation matrix $\mathbf{R} \in \mathbb{R}^{C \times C}$. Considering the orthogonality of \mathbf{P} , we transform \mathbf{Y} into its orthogonal counterpart $\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}}$. Then cluster indicator matrix \mathbf{Y} can be recovered by finding a rotation matrix \mathbf{R} to minimize the squared Euclidean distance between \mathbf{PR} and $\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}}$. 3. To guarantee the diversity of different views, a weight is introduced for each view. Overall, the objective function of the proposed DGMC is as follows:

$$\max_{\mathbf{P}, \mathbf{Y}, \mathbf{R}} \sum_{v=1}^V \alpha^{(v)} \text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}) - \lambda \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2 \quad \text{s.t. } \mathbf{P}^\top \mathbf{P} = \mathbf{I}, \mathbf{Y} \in \text{Ind}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}, \quad (1)$$

where **Ind** defines the set of cluster indicator matrices, λ is a balance parameter, $\alpha^{(v)}$ is the pre-weight of the v -th view. Nevertheless, it is hard to determine weights without prior knowledge.

Contributions

A novel approach named Dependence-Guided Multi-view Clustering (DGMC) is proposed. The main contributions of this paper are summarized as follows.

- The proposed model enhances the dependence between unified embedding learning and clustering, as well as increases the dependence between unified embedding and embedding of each view.
- A joint framework for unified embedding learning and clustering is proposed.
- A unified embedding can be learned from different views in Reproducing Kernel Hilbert Spaces (RKHSs) to capture the high-order and non-linear dependence among these views.
- Implicit-weight learning mechanism enhances the diversity of different views.

An Equivalent Model

In this paper, an implicit-weight learning mechanism is introduced to smartly learn $\alpha^{(v)}$. To this end, we give the Remark that problem (1) is equivalent to the following problem (2) with implicit weights:

$$\max_{\mathbf{P}, \mathbf{Y}, \mathbf{R}} \sum_{v=1}^V (\text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^r - \lambda \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2 \quad (2)$$

s.t. $\mathbf{P}^\top \mathbf{P} = \mathbf{I}, \mathbf{Y} \in \text{Ind}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}$,

where $r > 1$ controls the curvature of weighted-learning curve. Suppose the solutions of problems (1) and (2) are $\langle \mathbf{P}_*, \mathbf{Y}_*, \mathbf{R}_* \rangle$ and $\langle \mathbf{P}_0, \mathbf{Y}_0, \mathbf{R}_0 \rangle$, respectively. According to the KKT condition *w.r.t* \mathbf{P} of problem (1), it can be easily verified that $\langle \mathbf{P}_0, \mathbf{Y}_0, \mathbf{R}_0 \rangle = \langle \mathbf{P}_*, \mathbf{Y}_*, \mathbf{R}_* \rangle$ if $\alpha^{(v)}$ is calculated as $\alpha^{(v)} = r \cdot (\text{Tr}(\mathbf{H}\mathbf{P}_0 \mathbf{P}_0^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^{r-1}$.

Optimization

1. When \mathbf{P} and \mathbf{R} are fixed, problem (1) becomes

$$\min_{\mathbf{Y} \in \text{Ind}} \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2. \quad (3)$$

Further, by simply deriving problem (3), we have

$$\max_{\mathbf{Y} \in \text{Ind}, \mathbf{G} = \mathbf{PR}} \text{Tr}((\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} \mathbf{Y}^\top \mathbf{G}). \quad (4)$$

Let $\mathbf{Q} = \mathbf{Y}^\top \mathbf{Y}$, then $q_{lj} = \sum_{i=1}^N y_{li} y_{ij}$. Since $\mathbf{Y} \in \text{Ind}$, $y_{li} y_{ij} = 0$ and $q_{lj} = 0$ hold if $l \neq j$. Thus $\mathbf{Q}^{-\frac{1}{2}}$ is a diagonal matrix with the l -th diagonal element as $(y_l^\top y_l)^{-\frac{1}{2}}$. Then problem (4) becomes

$$\max_{\mathbf{Y} \in \text{Ind}} \sum_{l=1}^C y_l^\top g_l (y_l^\top y_l)^{-\frac{1}{2}}. \quad (5)$$

Since this problem is independent among different rows, we can solve \mathbf{Y} row by row. Given an optimal $\tilde{\mathbf{Y}}$, to update the i -th row \mathbf{y}^i , all we need to consider is the increment of the objective function value from $y_{il} = 0$ to $y_{il} = 1$. Therefore, Problem (5) can be solved by **coordinate descent method**.

2. When \mathbf{P} and \mathbf{Y} are fixed, problem (1) becomes

$$\min_{\mathbf{R}^\top \mathbf{R} = \mathbf{I}} \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2, \quad (6)$$

which is the **classical orthogonal procrustes problem**. Then a closed-form solution to \mathbf{R} is $\mathbf{R} = \mathbf{UV}^\top$, where \mathbf{USV}^\top is the SVD of $(\mathbf{P}^\top \mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}})$.

3. When \mathbf{Y} and \mathbf{R} are fixed, problem (1) becomes

$$\max_{\mathbf{P}^\top \mathbf{P} = \mathbf{I}} \text{Tr}(\mathbf{P}^\top \mathbf{A}\mathbf{P}) - \lambda \text{Tr}(\mathbf{P}^\top \mathbf{B}), \quad (7)$$

where $\mathbf{A} = \mathbf{H} \sum_{v=1}^V w^{(v)} \mathbf{F}^{(v)} \mathbf{F}^{(v)\top} \mathbf{H}$, $\mathbf{Z} = \mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}}$, $\mathbf{B} = \mathbf{ZR}^\top$. Problem (7) is the **typical Quadratic Problem on the Stiefel Manifold (QPSM)**, which can be solved by an efficient algorithm [2].

Theoretical Support

The Lagrangian function of problem (2) is as follows:

$$\mathcal{L} = \sum_{v=1}^V (\text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^r - \lambda \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2 - \Delta(\mathbf{P}, \mathbf{Y}, \mathbf{R}, \Lambda), \quad (8)$$

where $\Delta(\mathbf{P}, \mathbf{Y}, \mathbf{R}, \Lambda)$ is the penalty term for the constraints in problem (2). Let the derivative of problem (8) *w.r.t* \mathbf{P} be zero, and according to the chain-rule, we have

$$\sum_{v=1}^V \frac{\partial (\text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^r}{\partial \text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top})} \frac{\partial \text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top})}{\partial \mathbf{P}} - \lambda \frac{\partial \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2}{\partial \mathbf{P}} - \frac{\partial \Delta(\mathbf{P}, \mathbf{Y}, \mathbf{R}, \Lambda)}{\partial \mathbf{P}} = 0. \quad (9)$$

By denoting $w^{(v)} \stackrel{\text{def}}{=} \frac{\partial (\text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^r}{\partial \text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top})} = r \cdot (\text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top}))^{r-1}$ (Eq. (10)), then Eq. (9) becomes

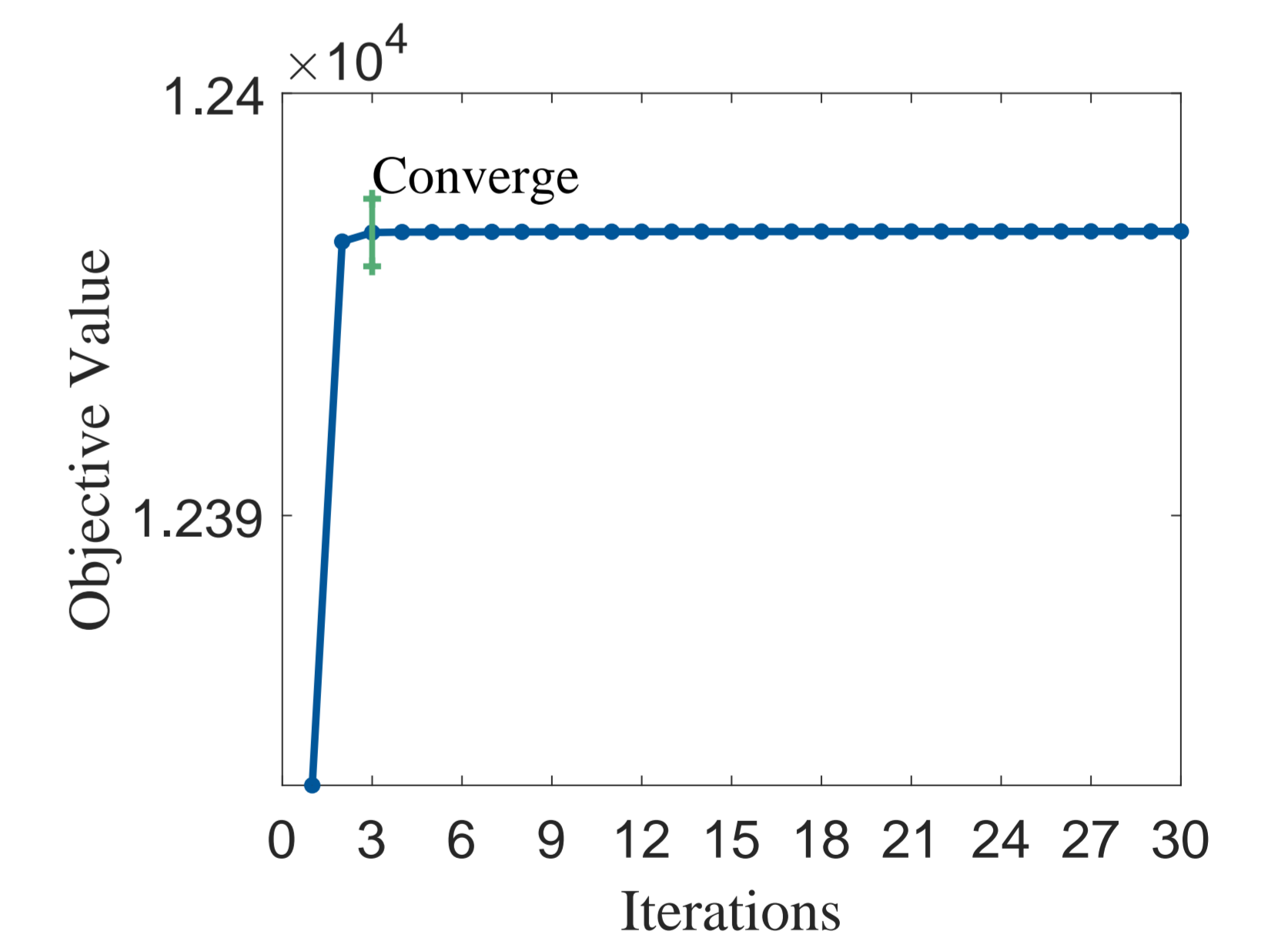
$\sum_{v=1}^V w^{(v)} \frac{\partial \text{Tr}(\mathbf{HPP}^\top \mathbf{H}\mathbf{F}^{(v)} \mathbf{F}^{(v)\top})}{\partial \mathbf{P}} - \lambda \frac{\partial \|\mathbf{Y}(\mathbf{Y}^\top \mathbf{Y})^{-\frac{1}{2}} - \mathbf{PR}\|_F^2}{\partial \mathbf{P}} - \frac{\partial \Delta(\mathbf{P}, \mathbf{Y}, \mathbf{R}, \Lambda)}{\partial \mathbf{P}} = 0$, which is the KKT condition *w.r.t* \mathbf{P} of prob-

lem (1) by letting $\alpha^{(v)} = w^{(v)}$. Overall, problem (2) can be optimized by iteratively performing the two steps: **Step 1.** update $w^{(v)}$ by Eq. (10); **Step 2.** update $\langle \mathbf{P}, \mathbf{Y}, \mathbf{R} \rangle$ by letting $\alpha^{(v)} = w^{(v)}$ and solving problem (1). Since there are three optimization variables in problem (1), we adopt alternative iterative optimization strategy to optimize them.

Experiment Results

The following Table gives the clustering results of the proposed DGMC and related four methods in terms of four metrics on three real-world datasets. The proposed DGMC shows large advantages over other methods in all the cases. Moreover, The following Figure shows the convergence curve of the proposed Algorithm on SCENE dataset. We can see that the proposed Algorithm can converge within 3 iterations.

Dataset	Metric	AMGL	MEA	MLAN	MVGL	DGMC
MSRC	ACC	0.8571	0.8714	0.6952	0.8714	0.9048
	NMI	0.7623	0.7834	0.6565	0.7731	0.8102
	ARI	0.7081	0.7199	0.5332	0.7152	0.7861
	F	0.7494	0.7597	0.6089	0.7560	0.8160
BBCSport	ACC	0.7206	0.4963	0.7279	0.7169	0.9062
	NMI	0.6867	0.2347	0.7146	0.6858	0.8230
	ARI	0.5832	0.1492	0.6069	0.5857	0.8510
	F	0.7088	0.4546	0.7244	0.7098	0.8881
SCENE	ACC	0.5856	0.6031	0.5335	0.3251	0.6763
	NMI	0.5017	0.5042	0.4596	0.2093	0.5282
	ARI	0.3821	0.3907	0.3158	0.0728	0.4538
	F	0.4820	0.4901	0.4377	0.2680	0.5243



Algorithm Description

Algorithm to solve problem (2).

Input: Initialized $\mathbf{P}, \mathbf{R}, \mathbf{X}^{(v)}, \mathbf{F}^{(v)}, V, \lambda, r$

while not converge do

1. Update $w^{(v)}$ via Eq. (10).
2. Update \mathbf{Y} via solving problem (3).
3. Update \mathbf{R} via solving problem (6).
4. Update \mathbf{P} via solving problem (7).

Output: \mathbf{Y}

Reference

- [1] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *Proceedings of the International Conference on Algorithmic Learning Theory*, pages 63–77. Springer, 2005.
- [2] Feiping Nie, Rui Zhang, and Xuelong Li. A generalized power iteration method for solving quadratic problem on the stiefel manifold. *Science China Information Sciences*, 60(11):112101, 2017.