

A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks

Akifumi Okuno^{1,2}, Tetsuya Hada³, and Hidetoshi Shimodaira^{1,2}

¹Graduates school of Informatics, Kyoto University, Kyoto, Japan, ²RIKEN Center for Artificial Intelligence Project (AIP), Tokyo, Japan, ³Recruit Technologies Co., Ltd., Tokyo, Japan.



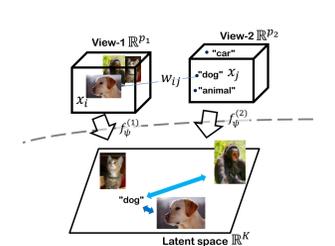
Introduction to Multi-view feature learning

Different types of data, such as user information, posted images, and their tags, etc. in Social Networking Service (SNS) is called **multi-view data**.



The number of view is $D \in \mathbb{N}$, and i -th data v_i belongs to view- $d_i \in \{1, 2, \dots, D\}$. i -th data v_i has p_{d_i} -dimensional vector representation $\mathbf{x}_i \in \mathbb{R}^{p_{d_i}}$, which we call **data vector**. The strength of association between $\mathbf{x}_i, \mathbf{x}_j$ is denoted as $w_{ij} = w_{ji} \geq 0$, and is called **matching weight**. $\{x_i\}$ and $\{w_{ij}\}$ are observed variables in our setting.

These multi-view data is hard to analyze, because their dimension may differ depending on the view. To get over the problem, data vector is transformed into **feature vector** $\mathbf{y}_i := f_{\psi}^{(d_i)}(\mathbf{x}_i) \in \mathbb{R}^K$ by continuous maps $f_{\psi}^{(d)} : \mathbb{R}^{p_d} \rightarrow \mathbb{R}^K$, ($d = 1, 2, \dots, D$) so that the inner product similarity of $\{\mathbf{y}_i\}$ in \mathbb{R}^K represents matching weights $\{w_{ij}\}$. The dimension K is the same among different views.



We call the procedure in general to transform multi-view data vectors $\{\mathbf{x}_i\}$ into feature vectors $\{\mathbf{y}_i\}$ as **multi-view feature learning (MvFL)**.

However, every existing methods for MvFL has at least one of following **drawbacks**: (i)Limited to **linear** setting, (ii) Limited to **one-to-one** association, (iii)Requires **high computational cost**.

Proposed probabilistic model

To realize MVFL simultaneously solving above mentioned problems (i)–(iii), we propose a novel probabilistic model of link weights:

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} \text{Po}(\exp(\mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}))), \quad (1)$$

where $\mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) := \alpha^{(d_i, d_j)} \exp(\langle f_{\psi}^{(d_i)}(\mathbf{x}_i), f_{\psi}^{(d_j)}(\mathbf{x}_j) \rangle)$ and $\boldsymbol{\theta} := (\boldsymbol{\psi}, \{\alpha^{(d,e)}\})$ is to be estimated.

As the continuous map $f_{\psi}^{(d)} : \mathbb{R}^{p_d} \rightarrow \mathbb{R}^K$, we use

- **Neural network (NN)**. Any deterministic NN can be used, but we basically consider multi layer perceptron (MLP).
- **Linear transformation (LT)**. By applying identity activation function, MLP reduces to LT; LT is a special case of NN.

Probabilistic Multi-view Graph Embedding

To find the optimal $\boldsymbol{\theta} = (\boldsymbol{\psi}, \{\alpha^{(d,e)}\})$, we consider to maximize the log-likelihood based on (1),

$$\ell(\boldsymbol{\theta}) := \sum_{1 \leq i < j \leq n} (w_{ij} \log \mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) - \mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta})),$$

by gradient-descent (GD). However, using a plain GD requires summing up $O(n^2)$ terms for computing the gradient. To reduce the high computational complexity, we alternatively utilize

$$\frac{\partial}{\partial \boldsymbol{\theta}} \left\{ \sum_{(i,j) \in \mathcal{W}_n} w_{ij} \log \mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) - \lambda \sum_{(i,j) \in \mathcal{I}_n} \mu(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) \right\}$$

as an approximation of the gradient, where $\lambda \geq 0$ is a tuning parameter and $\mathcal{I}_n, \mathcal{W}_n$ contain elements re-sampled from $\mathcal{I}_n := \{(i, j) \mid 1 \leq i < j \leq n\}$, $\mathcal{W}_n := \{(i, j) \in \mathcal{I}_n \mid w_{ij} > 0\}$ uniformly so that they contain m_1, m_2 elements, respectively.

Using $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\psi}}, \{\hat{\alpha}^{(d,e)}\})$ obtained by minibatch SGD, we have optimal feature vector

$$\mathbf{y}_i := f_{\hat{\boldsymbol{\psi}}}^{(d_i)}(\mathbf{x}_i) \in \mathbb{R}^K. \quad (2)$$

We call the above mentioned procedure to obtain feature vectors as **Probabilistic Multi-view Graph Embedding (PMvGE)**.

Comparison

(Nv): Number of views, (MM): Many-to-many, (NL): Non-linear, (Ind): Inductive, (Lik): Likelihood-based. PMvGE has all the properties. Nv = D represents that the method can deal with arbitrary number of views.

	(Nv)	(MM)	(NL)	(Ind)	(Lik)
CCA (Hotelling 1936)	2			✓	
Deep CCA (Andrew et al. 2013)	2		✓	✓	
MCCA (Kettenring 1971)	D			✓	
SGE (Belkin et al. 2001)	0	✓			
LINE (Tang et al. 2015)	0	✓			✓
LPP (He et al. 2004)	1	✓		✓	
CvGE (Huang et al. 2013)	2	✓		✓	
CDMCA (Shimodaira 2016)	D	✓		✓	
DeepWalk (Perozzi et al. 2014)	0	✓			✓
SBM (Holland et al. 1983)	1	✓		✓	✓
GCN (Kipf et al. 2017)	1	✓	✓	✓	✓
GraphSAGE (Hamilton et al. 2017)	1	✓	✓	✓	✓
IDW (Dai et al. 2018)	1	✓	✓	✓	✓
PMvGE (Proposed)	D	✓	✓	✓	✓

Representation power of PMvGE with NN

Theorem 1

$f_{\psi}^{(d)} : [-M, M]^{p_d} \rightarrow [-M', M']^K$, ($d = 1, 2, \dots, D$) are continuous maps and $g_* : [-M', M']^{2K} \rightarrow \mathbb{R}$ is a positive-definite kernel for some $M, M' > 0$ and $K_* \in \mathbb{N}$. For arbitrary $\varepsilon > 0$, by specifying large T, K , there exist MLPs $f_{\psi}^{(d)} : \mathbb{R}^{p_d} \rightarrow \mathbb{R}^K$ with T hidden units and ReLU or sigmoid activation s.t.

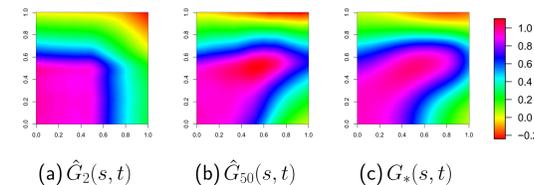
$$\left| g_*(f_{\psi}^{(d)}(\mathbf{x}), f_{\psi}^{(e)}(\mathbf{x}')) - \langle f_{\psi}^{(d)}(\mathbf{x}), f_{\psi}^{(e)}(\mathbf{x}') \rangle \right| < \varepsilon,$$

$$\forall (\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}, \forall d, e.$$

We visualize Theorem 1. With cosine similarity g_* , we define

$$G_*(s, t) := g_*(f_*(s\mathbf{e}_1), f_*(t\mathbf{e}_2)), \hat{G}_K(s, t) := \langle f_{\psi}(s\mathbf{e}_1), f_{\psi}(t\mathbf{e}_2) \rangle,$$

where $f_*(\mathbf{x}) = (x_1, \cos x_2, \exp(-x_3), \sin(x_4 - x_5))$, $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{R}^5$ are directions, $f_{\psi} : \mathbb{R}^5 \rightarrow \mathbb{R}^K$ is ReLU-based MLP with $T = 10^3$ hidden units and K output units.



Figures (a)–(c) show that $\hat{G}_K(s, t)$ with sufficiently large T, K approximates $G_*(s, t)$ well.

- Our new paper [Okuno and Shimodaira (2018). “On representation power of neural network-based graph embedding and beyond.” arXiv:1805.12332] will be presented at TADGM workshop in ICML2018!!

PMvGE with Linear Transformation approximately generalizes various existing methods

As the closest approach, Cross-Domain Matching Correlation Analysis (CDMCA; Shimodaira, 2016) obtains linear transformation (LT) $\mathbf{y}_i := \boldsymbol{\psi}^{(d_i)\top} \mathbf{x}_i$ by

$$\arg \max_{\{\boldsymbol{\psi}^{(d)}\}} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \langle \boldsymbol{\psi}^{(d_i)\top} \mathbf{x}_i, \boldsymbol{\psi}^{(d_j)\top} \mathbf{x}_j \rangle,$$

with a quadratic constraint on $\boldsymbol{\psi}$, where $\boldsymbol{\psi}$ is a concatenation of $\{\boldsymbol{\psi}^{(d)}\}_{d=1}^D$. The obtained matrix is denoted as $\hat{\boldsymbol{\psi}}_{\text{CDMCA}}$. Here, we approximate PMvGE with LT. $\tilde{\ell}_Q(\boldsymbol{\psi})$ denotes a quadratic approximation of the log-likelihood $\ell(\boldsymbol{\theta})$ around $\boldsymbol{\psi}$ with $\alpha^{(d,e)} = 1, \forall (d, e)$. The maximizer of $\tilde{\ell}_Q(\boldsymbol{\psi})$ is denoted as $\hat{\boldsymbol{\psi}}_{\text{Apr.PMvGE}}$.

Then, there exist $\gamma_1, \dots, \gamma_K \geq 0$ such that

$$\hat{\boldsymbol{\psi}}_{\text{Apr.PMvGE}} = \hat{\boldsymbol{\psi}}_{\text{CDMCA}} \text{diag}(\gamma_1^{1/2}, \dots, \gamma_K^{1/2}).$$

Since CDMCA generalizes various existing methods such as HIMFAC (Nori et al., 2012), CvGE (Huang et al., 2013), CCA, PCA, thus PMvGE with LT approximately generalizes these methods as well.

Experiments

Dataset: We use Cora citation dataset (Sen et al., 2008, C-cite) for task 1 and 2, and Animal with Attribute dataset (Lampert et al., 2009, AwA) for task 3.

- **C-cite** is ($D=1$)-view dataset consisting of 2,708 nodes with 1,433-dim bag-of-words vectors. Each node has a class label of 7 classes.
- **AwA** is ($D=2$)-view dataset, which consists of resampled 2,500 images with 4,096-dim DeCAF features (Donahue et al., 2014), and 85 attributes with 300-dim GloVe (Pennington et al., 2014) features.

For each dataset, we preliminary compute feature vectors. Then several existing methods are applied for the following tasks:

- **Task-1:** Nodes in C-cite is classified into 7-classes using multi-class logistic regression. Results are evaluated by classification accuracy.
- **Task-2:** We made 7 clusters of nodes in C-cite by k -means clustering. Results are evaluated by Normalized Mutual Information (NMI).
- **Task-3:** For each query image, we rank attributes according to the cosine similarity of feature vectors across views. Results are evaluated by Average Precision (AP).

Results are listed in the form [Ave. \pm std.dev] over 10 times experiments.

		(A)	(B)
Task 1 (D=1)	ISOMAP	54.5 \pm 1.78	54.8 \pm 2.43
	LLE	30.2 \pm 1.91	31.9 \pm 2.62
	SGE	47.6 \pm 1.64	-
	MDS	29.8 \pm 2.25	-
	DeepWalk	54.2 \pm 2.04	-
Task 2 (D=1)	GraphSAGE	60.8 \pm 1.73	57.1 \pm 1.61
	PMvGE	74.8 \pm 2.55	71.1 \pm 2.10
	SBM	4.37 \pm 1.44	2.81 \pm 0.10
	ISOMAP	13.0 \pm 0.36	14.3 \pm 1.98
	LLE	7.40 \pm 3.40	9.47 \pm 3.00
Task 3 (D=2)	SGE	1.41 \pm 0.34	-
	MDS	2.81 \pm 0.10	-
	DeepWalk	16.7 \pm 1.05	-
	GraphSAGE	19.6 \pm 0.93	12.4 \pm 3.00
	PMvGE	35.9 \pm 0.88	30.5 \pm 3.90
Task 3 (D=2)	CCA	45.5 \pm 0.20	42.4 \pm 0.30
	Deep CCA	41.4 \pm 0.30	41.2 \pm 0.35
	SGE	43.5 \pm 0.39	-
	DeepWalk	71.3 \pm 0.57	-
	PMvGE	71.5 \pm 0.48	70.5 \pm 0.53

Conclusion

We proposed Probabilistic Multi-view Graph Embedding (PMvGE) that simultaneously realizes (i)Non-linear, (ii)many-to-many, (iii)computationally efficient multi-view feature learning. Experiments demonstrated that PMvGE outperformed existing methods.

References

- Okuno, A. and Shimodaira, H. (2018). On representation power of neural network-based graph embedding and beyond. In *Proceedings of the ICML2018 workshop on Theoretical Foundations and Applications of Deep Generative Models (TADGM)*.
- Shimodaira, H. (2016). Cross-validation of matching correlation analysis by resampling matching weights. *Neural Networks*, 75:126–140.
- Andrew, G., Arora, R., Biles, J., and Livescu, K. (2013). Deep Canonical Correlation Analysis. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1247–1255.