

京都大学 KYOTO UNIVERSITY



ISI WSC 2019 – ISI World Statistics Congress 2019/08/23 Journal of Data Science, Statistics and Visualisation Session

Multiview graph embedding as a generalization of canonical correlation analysis

Hidetoshi Shimodaira Kyoto University / RIKEN AIP

In collaboration with Kazuki Fukui, Tetsuya Hada, Geewook Kim, Akifumi Okuno, Takamasa Oshikiri

Slides available <u>http://stat.sys.i.kyoto-u.ac.jp/research/</u>

"SIPS-IPS" is on the board here??



Today's Talk

- CCA is generalized to multiview graph embedding
- Further generalization to non-linear multiview graph embedding with neural networks
- A theory of representation learning: "new" similarity models beyond inner product (IPS, SIPS, WIPS)

CCA and Multiview graph embedding

Canonical Correlation Analysis (CCA) has been used in pattern recognition



[Socher+ 2010], [Bruni+ 2014], [Lazaridou+ 2014], etc.

CCA maximizes the correlations or minimizes the distances



Hotelling (1936). Biometrika: "RELATIONS BETWEEN TWO SETS OF VARIATES"

Equivalent formulations of CCA

$$\max_{\boldsymbol{A}^{(1)},\boldsymbol{A}^{(2)}} \operatorname{tr} \boldsymbol{A}^{(1)\top} \boldsymbol{X}^{(1)\top} \boldsymbol{X}^{(2)} \boldsymbol{A}^{(2)}$$

subject to $A^{(1)\top}X^{(1)\top}X^{(1)}A^{(1)} = A^{(2)\top}X^{(2)\top}X^{(2)}A^{(2)} = I_K/2$

$$\begin{split} & \min_{\boldsymbol{A}^{(1)}, \boldsymbol{A}^{(2)}} \| \boldsymbol{x}_i^{(1)} \boldsymbol{A}^{(1)} - \boldsymbol{x}_i^{(2)} \boldsymbol{A}^{(2)} \|^2 \\ & \text{subject to} \quad \boldsymbol{A}^{(1)\top} \boldsymbol{X}^{(1)\top} \boldsymbol{X}^{(1)} \boldsymbol{A}^{(1)} + \boldsymbol{A}^{(2)\top} \boldsymbol{X}^{(2)\top} \boldsymbol{X}^{(2)} \boldsymbol{A}^{(2)} = \boldsymbol{I}_K \end{split}$$

$$\max_{A \in \mathbb{R}^{N \times K}} A^{\top} H A \qquad A = \begin{bmatrix} A^{(1)} \\ A^{(2)} \end{bmatrix} \qquad H = \begin{bmatrix} X^{(1)\top} X^{(1)} \end{bmatrix}$$

subject to $A^{\top} G A = I_K$
$$G = \begin{bmatrix} X^{(1)\top} X^{(1)} \\ X^{(2)\top} X^{(2)} \end{bmatrix}$$

This formulation is generalized to graph embedding

Many-to-many relations



ml-CCA [Ranjan+ 2015], cluster CCA [Rasiwasia+ 2014], etc.

CCA is generalized to graphembedding (two or more views)



Cross-domain matching correlation analysis (CDMCA) generalizes several multivariate analysis methods such as CCA, multiset-CCA, and PCA. Shimodaira (2016). Neural Networks. "Cross-validation of matching correlation analysis by resampling matching weights"

CDMCA (2-view case)

$$A = \begin{bmatrix} A^{(1)} \\ A^{(2)} \end{bmatrix} \qquad W = \begin{bmatrix} W^{(11)} & W^{(12)} \\ W^{(21)} & W^{(22)} \end{bmatrix} \qquad M = \begin{bmatrix} M^{(1)} \\ M^{(2)} \end{bmatrix} \qquad m_i = \sum_{j=1}^N w_{ij}$$

CCA is given by
$$W = \begin{bmatrix} I_n \\ I_n \end{bmatrix} \qquad M = \begin{bmatrix} I_n \\ I_n \end{bmatrix}$$

$$\max_{A \in \mathbb{R}^{N \times K}} A^{\top} H A \qquad H = X^{\top} W X = \begin{bmatrix} X^{(1)\top} W^{(11)} X^{(1)} & X^{(1)\top} W^{(12)} X^{(2)} \\ X^{(2)\top} W^{(21)} X^{(1)} & X^{(2)\top} W^{(22)} X^{(2)} \end{bmatrix}$$

subject to $A^{\top} G A = I_K \qquad G = X^{\top} M X = \begin{bmatrix} X^{(1)\top} M^{(1)} X^{(1)} & X^{(2)\top} M^{(2)} X^{(2)} \end{bmatrix}$

CDMCA is solved as a generalized eigenvalue problem

$$oldsymbol{G}^{-1/2 op}oldsymbol{H}oldsymbol{G}^{-1/2 op}oldsymbol{H}oldsymbol{G}^{-1/2}oldsymbol{u}_{oldsymbol{i}}=\lambda_{i}oldsymbol{u}_{i}\qquad\lambda_{1}\geq\lambda_{2}\geq\cdots$$
 $oldsymbol{\hat{A}}=oldsymbol{G}^{-1/2}(oldsymbol{u}_{1},\ldots,oldsymbol{u}_{K})$

MNIST handwritten digits

3-views: $\{0,1,...,9\}$ - images - {prime,even,odd}



11

Bilingual word representations

5-views



T. Oshikiri, K. Fukui, H. Shimodaira (2016). Cross-Lingual Word Representations via Spectral Graph Embeddings, The 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016.

estocolmo stockholm

10

5

PC1

Flickr image-tag search and "linguistic regularities"



Fukui, Oshikiri, Shimodaira (2017). Spectral Graph-Based Method of Multimodal Word Embedding. TextGraphs-11, the Workshop on Graph-based Methods for Natural Language Processing, ACL 2017.

Graph embedding with Neural Networks

PMvGE: Probabilistic Multi-view Graph Embedding



Okuno, Hada, Shimodaira (2018). A probabilistic framework for multi-view feature learning with manyto-many associations via neural networks. The International Conference on Machine Learning, ICML 2018.

Optimization of PMvGE

Log-likelihood:
$$\ell_n(\boldsymbol{\theta}) = \sum_{(i,j)\in\mathcal{P}_n} w_{ij} \langle \boldsymbol{y}_i, \boldsymbol{y}_j \rangle - \sum_{(i,j)\in\mathcal{I}_n} \exp(\langle \boldsymbol{y}_i, \boldsymbol{y}_j \rangle)$$
$$\mathcal{P}_n = \{(i,j): w_{ij} > 0, \ 1 \le i < j \le n\}$$
$$\mathcal{I}_n = \{(i,j): 1 \le i < j \le n\}$$

In each iteration, Stochastic Gradient Decent with mini-batch uses:

$$\ell_{\min}(\boldsymbol{\theta}) = \sum_{(i,j)\in\mathcal{P}_{\min}} w_{ij} \langle \boldsymbol{y}_i, \boldsymbol{y}_j \rangle - \tau \sum_{(i,j)\in\mathcal{I}_{\min}} \exp(\langle \boldsymbol{y}_i, \boldsymbol{y}_j \rangle)$$

 $\mathcal{P}_{\min}, \mathcal{I}_{\min}$ are mini-batch sampling of $\mathcal{P}_n, \mathcal{I}_n$

MLE estimates the true similarity up to a constant factor

$$\lim_{n \to \infty} \exp(\langle \boldsymbol{y}_i, \boldsymbol{y}_j \rangle) = \text{const} \times \mathbb{E}(w_{ij} | \boldsymbol{x}_i, \boldsymbol{x}_j)$$

Okuno and Shimodaira (arXiv 2019) Hyperlink Regression via Bregman Divergence

PMvGE generalizes CCA and CDMCA to non-linear transformations

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

Representation learning

Similarities are expressed via vector representations

Representation Learning on graphs aims to learn useful vector representations of nodes (e.g., words, users) in a graph-structured data.



users in a social network

$$\mathbb{E}(w_{ij}|\boldsymbol{x}_i, \boldsymbol{x}_j) = \exp(s_{ij})$$
Similarity: s_{ij}
Attributes: \boldsymbol{x}_i

user embeddings in a vector space \mathbb{R}^{K}

$$s_{ij} = \langle \boldsymbol{y}_i, \boldsymbol{y}_j
angle$$

Embedding: y_i (vector representation)

NN-based similarity learning: Siamese (twin) neural networks



Bromley, Guyon, Le Cun, Sackinger, Shah (1994). "Signature verification using a "Siamese" time delay neural network." NIPS.

Why many methods are based on inner product of vector representations?

The true underlying similarity

S(

Inner Product Similarity (IPS) "Neural networks + Inner product"

$$oldsymbol{x},oldsymbol{x}') \hspace{1cm} h_{ ext{IPS}}(oldsymbol{x},oldsymbol{x}') \coloneqq \langle oldsymbol{f}_{ ext{NN}}(oldsymbol{x}),oldsymbol{f}_{ ext{NN}}(oldsymbol{x}')
angle$$

$$s_{ij} = s(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) \qquad h_{ij} = h_{\text{IPS}}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j})$$

$$\begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1n} \\ s_{21} & s_{22} & \cdots & s_{2n} \\ \vdots & \ddots & \vdots \\ s_{n1} & s_{n2} & \cdots & s_{nn} \end{bmatrix} \stackrel{\text{approx.}}{\approx} \begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1n} \\ h_{21} & h_{22} & \cdots & h_{2n} \\ \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \cdots & h_{nn} \end{bmatrix}$$

Inner product is a key to good vector representations

The distributive law

$$\langle \boldsymbol{a} + \boldsymbol{b}, \boldsymbol{c}
angle = \langle \boldsymbol{a}, \boldsymbol{c}
angle + \langle \boldsymbol{b}, \boldsymbol{c}
angle$$

"linguistic regularities"

 $\langle \boldsymbol{y}_{ ext{king}} - \boldsymbol{y}_{ ext{man}} + \boldsymbol{y}_{ ext{woman}}, \boldsymbol{y}'
angle = \langle \boldsymbol{y}_{ ext{king}}, \boldsymbol{y}'
angle - \langle \boldsymbol{y}_{ ext{man}}, \boldsymbol{y}'
angle + \langle \boldsymbol{y}_{ ext{woman}}, \boldsymbol{y}'
angle$

will be maximized by $~~m{y}'=m{y}_{ ext{queen}}$

Inner Product Similarity (IPS) can approximate any Positive Definite (PD) similarity

Theorem 1 (Okuno, Kim, Shimodaira 2018) $s_{\rm PD}(\boldsymbol{x}, \boldsymbol{x}')$ is a positive definite (PD) similarity. For arbitrary small $\epsilon > 0$ there exists $h_{\rm IPS}(\boldsymbol{x}, \boldsymbol{x}') = \langle \boldsymbol{f}_{\rm NN}(\boldsymbol{x}), \boldsymbol{f}_{\rm NN}(\boldsymbol{x}') \rangle$ such that $|s_{\rm PD}(\boldsymbol{x}, \boldsymbol{x}') - h_{\rm IPS}(\boldsymbol{x}, \boldsymbol{x}')| < \epsilon$ for any $\boldsymbol{x}, \boldsymbol{x}' \in \mathcal{X}$

$$s_{\mathrm{PD}}(\boldsymbol{x}, \boldsymbol{x}') \Leftrightarrow h_{\mathrm{IPS}}(\boldsymbol{x}, \boldsymbol{x}')$$

Proof: assuming s_{PD} is continuous on a compact set X, the theorem follows by combining Mercer's theorem of kernels and Universal approximation theorem of NN (Funahashi, 1989; Cybenko, 1989; Yarotsky, 2017; Telgarsky, 2017)

Positive definite kernel

$$\sum_{i=1}^{n} \sum_{j=1}^{n} c_i c_j s_{\text{PD}}(\boldsymbol{x}_i, \boldsymbol{x}_j) \ge 0$$

for arbitrary $c_1, c_2, \ldots, c_n \in \mathbb{R}, \quad oldsymbol{x}_1, oldsymbol{x}_2, \ldots, oldsymbol{x}_n \in \mathcal{X}$

(example)
$$\boldsymbol{y} = \boldsymbol{f}(\boldsymbol{x}), \ \boldsymbol{y}' = \boldsymbol{f}(\boldsymbol{x}'); \quad s_{\text{PD}}(\boldsymbol{x}, \boldsymbol{x}') = g(\boldsymbol{y}, \boldsymbol{y}')$$

 $g(\boldsymbol{y}, \boldsymbol{y}') = \langle \boldsymbol{y}, \boldsymbol{y}' \rangle$ It is obvious that IPS is PD.
 $g(\boldsymbol{y}, \boldsymbol{y}') = \left\langle \frac{\boldsymbol{y}}{\|\boldsymbol{y}\|}, \frac{\boldsymbol{y}'}{\|\boldsymbol{y}'\|} \right\rangle$

(counter example) $g(\boldsymbol{y}, \boldsymbol{y}') = -\|\boldsymbol{y} - \boldsymbol{y}'\|^2$

A hierarchy of similarity models



Note: we are not attempting

$$s_{\text{general}}(\boldsymbol{x}, \boldsymbol{x}') \Leftrightarrow f_{\text{NN}}(\boldsymbol{x}, \boldsymbol{x}')$$

Adding bias terms to IPS

Okuno, Kim, Shimodaira (2019). Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability. AISTATS2019.

Inner product similarity (IPS)

$$h_{\text{IPS}}(\boldsymbol{x}, \boldsymbol{x}') := \langle f_{\text{NN}}(\boldsymbol{x}), f_{\text{NN}}(\boldsymbol{x}') \rangle$$



Shifted Inner product similarity (SIPS) $h_{\rm SIPS}(\boldsymbol{x}, \boldsymbol{x}') := \langle \boldsymbol{f}_{\rm NN}(\boldsymbol{x}), \boldsymbol{f}_{\rm NN}(\boldsymbol{x}') \rangle + u_{\rm NN}(\boldsymbol{x}) + u_{\rm NN}(\boldsymbol{x}')$

Bias terms have been used in some applications (e.g. recommender systems)

Shifted Inner Product Similarity (SIPS) can approximate any Conditionally Positive Definite (CPD) similarity

Theorem 2 (Okuno, Kim, Shimodaira 2019)

$$\begin{split} s_{\mathrm{CPD}}(\boldsymbol{x},\boldsymbol{x}') \text{ is a conditionally positive definite (CPD) similarity.} \\ & \text{For arbitrary small } \epsilon > 0 \\ & \text{there exists } h_{\mathrm{SIPS}}(\boldsymbol{x},\boldsymbol{x}') \coloneqq \langle \boldsymbol{f}_{\mathrm{NN}}(\boldsymbol{x}), \boldsymbol{f}_{\mathrm{NN}}(\boldsymbol{x}') \rangle + u_{\mathrm{NN}}(\boldsymbol{x}) + u_{\mathrm{NN}}(\boldsymbol{x}') \\ & \text{such that } |s_{\mathrm{CPD}}(\boldsymbol{x},\boldsymbol{x}') - h_{\mathrm{SIPS}}(\boldsymbol{x},\boldsymbol{x}')| < \epsilon \quad \text{for any } \quad \boldsymbol{x}, \boldsymbol{x}' \in \mathcal{X} \end{split}$$

$$s_{\mathrm{CPD}}(\boldsymbol{x}, \boldsymbol{x}') \Leftrightarrow h_{\mathrm{SIPS}}(\boldsymbol{x}, \boldsymbol{x}')$$

Proof: combining Mercer's theorem + the Universal approximation theorem + Berg et al. (1984)

Conditionally Positive Definite (CPD) kernel

$$\sum_{i=1}^{n} \sum_{j=1}^{n} c_i c_j s_{\text{CPD}}(\boldsymbol{x}_i, \boldsymbol{x}_j) \ge 0$$

for arbitrary $c_1, c_2, \dots, c_n \in \mathbb{R}, \quad \boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_n \in \mathcal{X}$
with the condition $\sum_{i=1}^{n} c_i = 0$

$$\begin{array}{ll} \text{(example)} \quad {\bm{y}} = {\bm{f}}({\bm{x}}), \ {\bm{y}}' = {\bm{f}}({\bm{x}}'); & s_{\mathrm{CPD}}({\bm{x}}, {\bm{x}}') = g({\bm{y}}, {\bm{y}}') \\ g({\bm{y}}, {\bm{y}}) = -\mathrm{distance}({\bm{y}}, {\bm{y}}') \\ d_{\mathrm{Euclid}}({\bm{y}}, {\bm{y}}') = \|{\bm{y}} - {\bm{y}}'\|_2^{\alpha}, & 0 < \alpha \leq 2 \\ d_{\mathrm{Poincaré}}({\bm{y}}, {\bm{y}}') = \cosh^{-1} \left(1 + 2 \frac{\|{\bm{y}} - {\bm{y}}'\|^2}{(1 - \|{\bm{y}}\|^2)(1 - \|{\bm{y}}'\|^2)}\right), & \|{\bm{y}}\|, \|{\bm{y}}'\| < 1 \\ d_{\mathrm{Wasserstein-}p}({\bm{y}}, {\bm{y}}') = \inf_{\pi \in \Pi({\bm{y}}, {\bm{y}}')} \left(\int_{\mathcal{Z} \times \mathcal{Z}} d_Z(z, z')^p d\pi(z, z')\right)^{1/p} & \Pr_{\mathtt{p=2 \ if \ Z = \ R}}^{\mathtt{p=1 \ if \ -d_Z \ is \ CPD}} \end{array}$$

The reason why bias terms are necessary for CPD kernel

$$-rac{1}{2}\|m{y}-m{y}'\|^2 = \langlem{y},m{y}'
angle - rac{1}{2}\|m{y}\|^2 - rac{1}{2}\|m{y}'\|^2$$

A tree structure is easily expressed in Poincaré Nickel & Kiela (NIPS 2017)





(b) Embedding of a tree in \mathcal{B}^2



 $\|\mathbf{u} - \mathbf{v}\|^2$

 $d(\mathbf{u}, \mathbf{v})$

Figure 1: (a) Due to the negative curvature of \mathcal{B} , the distance of points increases exponentially (relative to their Euclidean distance) the closer they are to the boundary. (c) Growth of the Poincaré distance d(u, v) relative to the Euclidean distance and the norm of v (for fixed ||u|| = 0.9). (b) Embedding of a regular tree in \mathcal{B}^2 such that all connected nodes are spaced equally far apart (i.e., all black line segments have identical hyperbolic length).





Hyperbolid model and Poincare disk model (from Wikipedia)

Inner-Product Difference Similarity (IPDS)
can approximate any general similarities

$$s_{\text{general}}(\boldsymbol{x}, \boldsymbol{x}') = s_{\text{PD}}^+(\boldsymbol{x}, \boldsymbol{x}') - s_{\text{PD}}^-(\boldsymbol{x}, \boldsymbol{x}')$$

(Indefinite kernel: Ong et al., 2004)

Theorem 3 (Okuno, Kim, Shimodaira 2019)

$$egin{aligned} h_{ ext{IPDS}}(oldsymbol{x},oldsymbol{x}') &:= \langle oldsymbol{f}_{ ext{NN}}^+(oldsymbol{x}),oldsymbol{f}_{ ext{NN}}^-(oldsymbol{x}')
angle - \langle oldsymbol{f}_{ ext{NN}}^-(oldsymbol{x}),oldsymbol{f}_{ ext{NN}}^-(oldsymbol{x}')
angle \ s_{ ext{general}}(oldsymbol{x},oldsymbol{x}') & \Leftrightarrow h_{ ext{IPDS}}(oldsymbol{x},oldsymbol{x}') \end{aligned}$$

Vector representation: $oldsymbol{y} = (oldsymbol{f}_{\mathrm{NN}}^+(oldsymbol{x}), oldsymbol{f}_{\mathrm{NN}}^-(oldsymbol{x})) \in \mathbb{R}^K$

$$oldsymbol{f}_{\mathrm{NN}}^+(oldsymbol{x})\in\mathbb{R}^{K-q},\quadoldsymbol{f}_{\mathrm{NN}}^-(oldsymbol{x})\in\mathbb{R}^q$$

Choosing q from (0,K) is important but it can be difficult in practice

Weighted Inner Product Similarity (WIPS) avoids model selection

G. Kim, A. Okuno, K. Fukui & H. Shimodaira (2019). Representation Learning with Weighted Inner Product for Universal Approximation of General Similarities. Proceedings of 28th International Joint Conference on Artificial Intelligence (IJCAI-19).

Weighted inner product allows negative weight values

$$\langle \boldsymbol{y}, \boldsymbol{y}' \rangle_{\boldsymbol{\lambda}} = \lambda_1 y_1 y_1' + \lambda_2 y_2 y_2' + \dots + \lambda_K y_K y_K'$$

 $\lambda_1,\lambda_2,\ldots,\lambda_K\in\mathbb{R}$ are parameters to be learned

Theorem 4 (Kim, Okuno, Fukui, Shimodaira 2019) $h_{
m WIPS}(m{x},m{x}') := \langle m{f}_{
m NN}(m{x}), m{f}_{
m NN}(m{x}')
angle_{m{\lambda}}$ $s_{
m general}(m{x},m{x}') \Leftrightarrow h_{
m WIPS}(m{x},m{x}')$

WIPS performs very well

Experiments - datasets

3 graph-structured datasets are used.



Each webpage has A. semantic label in {Student, Faculty, Staff, Course, Project}

B. university label in {Cornell, Texas, Washington, Wisconsin}

From presentation slides of Kim, Okuno, Fukui, Shimodaira (IJCAI 2019)





The hypertexts are colored by its semantic labels (upper) for Student (navy) and Course (pink), and also university labels (lower) for Cornell (red), Texas (orange), Washington (green) and Wisconsin (blue).

Both class labels are clearly identified with IPDS and WIPS, whereas they become obscure in the other embeddings.

Hypertext Classification

Each webpage in Hypertext Network has

- A. semantic label ∈ {Student, Faculty, Staff, Course, Project}
- B. university label ∈ {Cornell, Texas, Washington, Wisconsin}

train set (observed) is used to train the embedder and classifiers



test set (invisible) is used for evaluation

			5115	ILD2	wIP5
A 56.	08 46.19 50 30.17	47.22	<u>69.09</u> 03.81	$\frac{71.70}{93.81}$	73.35

Graph Reconstruction

At training, assume that all nodes and links are visible. Use all data to train the model (f_{θ} and g_{λ}).

Embed all nodes \boldsymbol{x}_2 \boldsymbol{x}_1 $oldsymbol{x}_4$ $\boldsymbol{f}_{\boldsymbol{ heta}}: \mathcal{X} \mapsto \mathcal{Y}$ x_6 \boldsymbol{x}_3 x_5 \boldsymbol{y}_4 \boldsymbol{y}_2 $oldsymbol{y}_1$ \boldsymbol{y}_6 $oldsymbol{y}_5$ $oldsymbol{y}_3$ 2 4 $g_{\lambda}: \mathcal{Y}^2 \mapsto \mathbb{R}$ 6 3 Predict (reconstruct) all links ROC-AUC for prediction errors are calculated ►

		Reconstruction		
		10	50	100
t	IPS	91.99	94.23	94.24
ex	Poincaré	94.09	94.13	94.11
er	SIPS	95.11	95.12	95.12
yp	IPDS	95.12	95.12	95.12
H	WIPS	95.11	95.12	95.12
	IPS	85.01	86.02	85.80
ho	Poincaré	86.84	86.69	86.72
III	SIPS	90.01	91.35	91.06
-5	IPDS	90.13	91.68	91.59
U	WIPS	90.50	92.44	92.95
v	IPS	79.95	75.80	74.97
H	Poincaré	91.69	89.10	88.97
OU	SIPS	98.78	99.75	99.77
XC	IPDS	99.65	99.89	99.90
Ē	WIPS	99.64	99.85	99.87

Then, at evaluation,

Link Prediction

At training, assume that some nodes (and its links) are invisible. Use only observed sub-graph to train the model (f_{θ} and g_{λ}).



		Link prediction		
		10	50	100
t	IPS	77.73	77.62	77.16
ex	Poincaré	82.21	79.64	79.48
ert	SIPS	82.01	81.84	81.13
yp	IPDS	82.59	82.75	82.19
H	WIPS	82.38	82.68	82.93
1	IPS	83.83	84.41	84.02
ho	Poincaré	85.82	85.92	85.93
III	SIPS	88.24	88.69	88.67
0-5	IPDS	88.42	88.97	88.85
C	WIPS	88.16	89.43	89.40
v	IPS	67.25	65.71	65.38
H	Poincaré	83.04	79.52	78.97
DUC	SIPS	90.42	92.12	92.09
aXC	IPDS	95.99	96.37	96.41
E	WIPS	95.07	96.36	96.51

Summary of similarity models



$$\begin{split} h_{\text{WIPS}}(\boldsymbol{x}, \boldsymbol{x}') &:= \langle \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}), \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}') \rangle_{\boldsymbol{\lambda}} \\ h_{\text{IPDS}}(\boldsymbol{x}, \boldsymbol{x}') &:= \langle \boldsymbol{f}_{\text{NN}}^+(\boldsymbol{x}), \boldsymbol{f}_{\text{NN}}^+(\boldsymbol{x}') \rangle - \langle \boldsymbol{f}_{\text{NN}}^-(\boldsymbol{x}), \boldsymbol{f}_{\text{NN}}^-(\boldsymbol{x}') \rangle \\ h_{\text{SIPS}}(\boldsymbol{x}, \boldsymbol{x}') &:= \langle \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}), \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}') \rangle + u_{\text{NN}}(\boldsymbol{x}) + u_{\text{NN}}(\boldsymbol{x}') \\ h_{\text{IPS}}(\boldsymbol{x}, \boldsymbol{x}') &:= \langle \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}), \boldsymbol{f}_{\text{NN}}(\boldsymbol{x}') \rangle \end{split}$$

References

- A. Okuno and H. Shimodaira. Hyperlink Regression via Bregman Divergence. arXiv:1908.02573, 2019
- G. Kim, A. Okuno, K. Fukui & H. Shimodaira, Representation Learning with Weighted Inner Product for Universal Approximation of General Similarities, Proceedings of 28th International Joint Conference on Artificial Intelligence (IJCAI-19), Main track, 5031-5038, 2019
- A. Okuno & H. Shimodaira, Robust Graph Embedding with Noisy Link Weights, Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS), PMLR 89, 664-673, 2019
- A. Okuno, G. Kim & H. Shimodaira, Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability, Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS), PMLR 89, 644-653, 2019
- A. Okuno, T. Hada & H. Shimodaira, A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks, Proceedings of the 35th International Conference on Machine Learning (ICML), PMLR 80, 3888-3897, 2018
- K. Fukui, T. Oshikiri & H. Shimodaira, Spectral Graph-Based Method of Multimodal Word Embedding. TextGraphs-11, the Workshop on Graph-based Methods for Natural Language Processing, ACL 2017
- M. Nickel and D. Kiela, Poincare Embeddings for Learning Hierarchical Representations. In Advances in Neural Information Processing Systems (NIPS) 30, 6338–6347, 2017
- K. Fukui, A. Okuno & H. Shimodaira, Image and tag retrieval by leveraging image-group links with multi-domain graph embedding, 2016 IEEE International Conference on Image Processing (ICIP), 221–225, 2016
- T. Oshikiri, K. Fukui & H. Shimodaira, Cross-Lingual Word Representations via Spectral Graph Embeddings, Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL), 493–498, 2016
- H. Shimodaira, Cross-validation of matching correlation analysis by resampling matching weights, Neural Networks, 75, 126-140, 2016
- Harold Hotelling, Relations Between Two Sets of Variates, Biometrika, 28, 321-377, 1936