

Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability

(Slides ▷ <http://okuno.co/fimi.pdf>)
to appear in AISTATS2019 (arXiv:1810.03463)
Joint work with H. Shimodaira and G. Kim

Akifumi Okuno

Kyoto University and RIKEN AIP

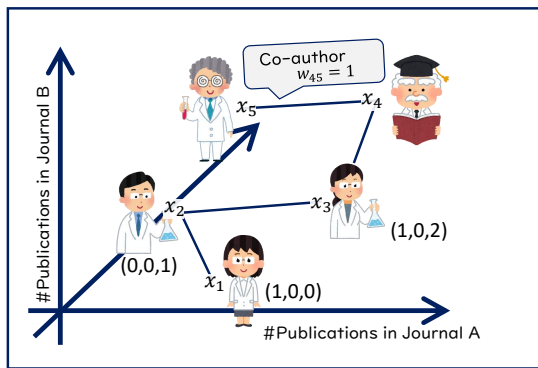
Purpose of this study is clarifying

expressive power

of neural network-based graph embedding (GE),
and propose

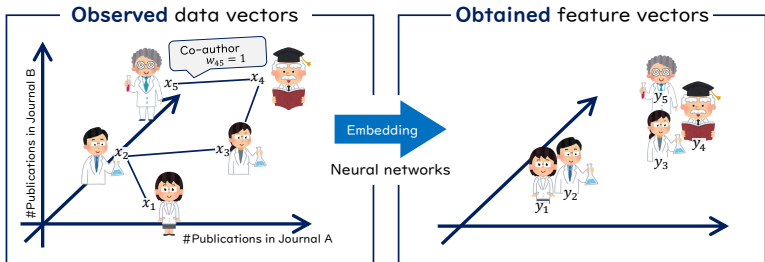
more expressive GE.

Given dataset



- Link weights $w_{ij} \geq 0$: co-authorship relation
- Data vectors \mathbf{x}_i : number of publications in each journal.

Summary of this study



Graph embedding learns a transformation $f : \mathcal{X} \rightarrow \mathcal{Y}$, so that

$$w_{ij} \approx \text{Sigmoid}(\underbrace{g(\mathbf{y}_i, \mathbf{y}_j)}_{\text{"Similarity"}}), \quad \mathbf{y}_i := f(\mathbf{x}_i).$$

g is specified as inner-product (Tang et al., 2015), Poincaré-dist. (Nickel and Kiela, 2017), etc..

Question \triangleright What kind of kernel g is good (in terms of the expressive power?)

Answer \triangleright Proposed **shifted inner-product similarity (SIPS)**.

SIPS can approximate many kernels, e.g., Poincaré-distance!

Today's talk is based on

- Okuno, Kim, and Shimodaira. "Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability," arXiv:1810.03463, to appear in AISTATS2019,

that is an extended version of

- Okuno and Shimodaira. "On representation power of neural network-based graph embedding and beyond," arXiv:1805.12332, ICML2018 workshop on TADGM.

We recently submitted a preprint that proposes a new similarity.

- Kim, Okuno, Fukui, and Shimodaira. "Representation Learning with Weighted Inner Product for Universal Approximation of General Similarities," arXiv:1902.10409 , *submitted*.

1 Introduction: Graph Embedding with Data Vectors

- Formulation: Graph Embedding with Data Vectors
- Existing Graph Embedding Methods
- Examples: Graph with Data Vectors

1 Introduction: Graph Embedding with Data Vectors

- Formulation: Graph Embedding with Data Vectors
- Existing Graph Embedding Methods
- Examples: Graph with Data Vectors

2 PD and CPD similarities (AISTATS2019)

- Poincaré Embedding
- Inner-Product Similarity (IPS)
- Shifted Inner-Product Similarity (SIPS)
- Numerical Experiments

1 Introduction: Graph Embedding with Data Vectors

- Formulation: Graph Embedding with Data Vectors
- Existing Graph Embedding Methods
- Examples: Graph with Data Vectors

2 PD and CPD similarities (AISTATS2019)

- Poincaré Embedding
- Inner-Product Similarity (IPS)
- Shifted Inner-Product Similarity (SIPS)
- Numerical Experiments

3 Recent Progress: General Similarities (arXiv:1902.10409)

- Non-CPD Similarities
- Inner-Product Difference Similarity (IPDS)
- Weighted Inner-Product Similarity (WIPS)
- Numerical Experiments

1 Introduction: Graph Embedding with Data Vectors

- Formulation: Graph Embedding with Data Vectors
- Existing Graph Embedding Methods
- Examples: Graph with Data Vectors

2 PD and CPD similarities (AISTATS2019)

- Poincaré Embedding
- Inner-Product Similarity (IPS)
- Shifted Inner-Product Similarity (SIPS)
- Numerical Experiments

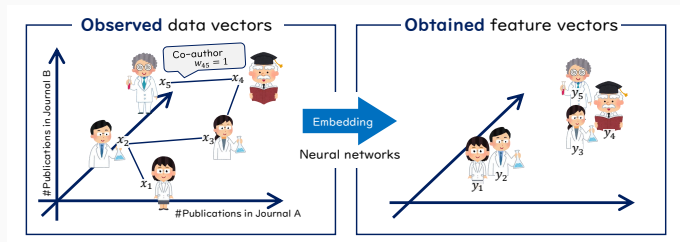
3 Recent Progress: General Similarities (arXiv:1902.10409)

- Non-CPD Similarities
- Inner-Product Difference Similarity (IPDS)
- Weighted Inner-Product Similarity (WIPS)
- Numerical Experiments

4 Conclusion

Introduction: Graph Embedding with Data Vectors

Formulation: Graph Embedding with Data Vectors



Given

- link weights $w_{ij} = w_{ji} \geq 0$, $1 \leq i < j \leq n$,
- data vectors $\mathbf{x}_i \in \mathcal{X} (\subset \mathbb{R}^p)$, $1 \leq i \leq n$,

graph embedding (with data vectors) learns a transformation $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$, so that

$$w_{ij} \approx \exp(\langle \mathbf{y}_i, \mathbf{y}_j \rangle), \quad \mathbf{y}_i := \mathbf{f}(\mathbf{x}_i).$$

$\exp(\cdot)$ can be replaced with Sigmoid, or some other functions.

Existing Graph Embedding Methods

There are (mainly) 3 different types of graph embedding methods.

Spectral

Chung (1997),
Belkin and Niyogi (2003),
He and Niyogi (2004),...

- 😊 eigen-decomposition
- 😞 computationally hard
- 😞 limited to linear

Probabilistic

Mikolov et al. (2013),
Tang et al. (2015),...

- 😊 computationally tractable
- 😊 non-linear
- 😊 inductive

We employ this model.

Convolution

Kipf and Welling (2016),
Kipf and Welling (2017),
Hamilton et al. (2017),...

- 😊 computationally tractable
- 😊 non-linear
- 😊 graph convolution
- 😞 not inductive

Existing Graph Embedding Methods (Spectral)

We **do not consider this model** in this study.

Locality Preserving Projections (He and Niyogi, 2004, LPP) learns a linear-transformation

$$\mathbf{y}_i := \mathbf{A}^\top \mathbf{x}_i.$$

so that the obtained feature vectors $\{\mathbf{y}_i\}_{i=1}^n$ maximize

$$\sum_{i=1}^n \sum_{j=1}^n w_{ij} \langle \mathbf{y}_i, \mathbf{y}_j \rangle,$$

under a quadratic constraint $\sum_{i=1}^n \sum_{j=1}^n w_{ij} \mathbf{y}_i \mathbf{y}_i^\top = \mathbf{I}$. LPP reduces to Spectral Graph Embedding (Chung, 1997) if $\{\mathbf{x}_i\}$ is specified as 1-hot vectors.

😊 can be **solved by eigen-decomposition**.

😞 **computationally hard**.

😞 limited to **linear** setting.

Existing Graph Embedding Methods (Probabilistic)

We employ this model in this study.

Given binary weights $w_{ij} \in \{0, 1\}$, Large-scale Information Network Embedding (Tang et al., 2015, LINE) learns a neural network

$$\mathbf{y}_i := \mathbf{f}(\mathbf{x}_i),$$

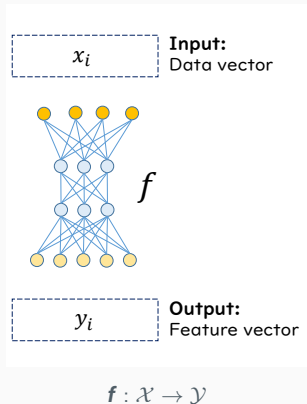
by assuming a probabilistic model¹

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} \text{Bernoulli}(\sigma(\langle \mathbf{y}_i, \mathbf{y}_j \rangle)),$$

where $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is a neural network (NN). **NN can be trained by maximizing the log-likelihood for $(\{w_{ij}\}_{1 \leq i < j \leq n} \mid \{\mathbf{x}_i\}_{i=1}^n)$ without any constraint.**

The log-likelihood can be maximized by stochastic algorithms.

- 😊 **computationally tractable.**
- 😊 **non-linear.**
- 😊 **inductive, i.e.,** unseen data vector \mathbf{x}_{n+1} can be transformed.



¹This model can be replaced, e.g., Poisson distribution (Okuno et al., 2018).

Existing Graph Embedding Methods (Convolution)

We do not consider this model in this study.

For binary weights $w_{ij} \in \{0, 1\}$,
Graph Autoencoder (Kipf and Welling, 2016, GAE) learns a graph convolutional network (Kipf and Welling, 2017, GCN)

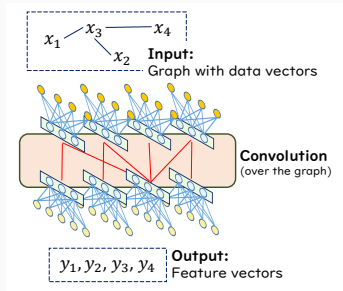
$$(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n) = \text{GCN}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; \mathbf{W}),$$

so that

$$w_{ij} \approx \sigma(\langle \mathbf{y}_i, \mathbf{y}_j \rangle).$$

- 😊 computationally tractable.
- 😊 non-linear.
- 😊 aggregate the information by considering the graph structure.

😞 not inductive, *i.e.*, unseen data vector \mathbf{x}_{n+1} cannot be transformed.²



$$\text{GCN}: \mathcal{X}^n \times \mathcal{W} \rightarrow \mathcal{Y}^n$$

²GraphSAGE Hamilton et al. (2017) and some recent studies try to make GCN inductive, by utilizing associations between the unseen data vector \mathbf{x}_{n+1} and the learned vectors $\{\mathbf{x}_i\}_{i=1}^n$.

Existing Graph Embedding Methods

There are (mainly) 3 different types of graph embedding methods.

Spectral

Chung (1997),
Belkin and Niyogi (2003),
He and Niyogi (2004),...

- 😊 eigen-decomposition
- 😞 computationally hard
- 😞 limited to linear

Probabilistic

Mikolov et al. (2013),
Tang et al. (2015),...

- 😊 computationally tractable
- 😊 non-linear
- 😊 inductive

We employ this model.

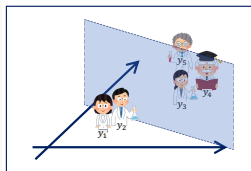
Convolution

Kipf and Welling (2016),
Kipf and Welling (2017),
Hamilton et al. (2017),...

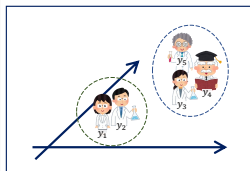
- 😊 computationally tractable
- 😊 non-linear
- 😊 graph convolution
- 😞 not inductive

Further Statistical Analyses

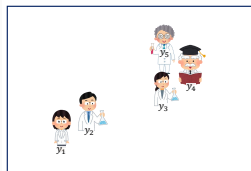
Further statistical analyses can be conducted on the vectors $\{\mathbf{y}_i\}_{i=1}^n \subset \mathcal{Y}$.



(a) Discriminant analysis



(b) Clustering



(c) Visualization

Using obtained feature vectors have **two advantages**

- 😊 **computationally tractable** (usually, $\dim \mathcal{Y} \leq \dim \mathcal{X}$),
- 😊 expected to **improve the performance** of further analyses.

Obtained feature vectors $\{\mathbf{y}_i\}_{i=1}^n \subset \mathcal{Y} (\subset \mathbb{R}^K)$ simultaneously preserves

1. the locality of $\{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{X} (\subset \mathbb{R}^p)$, and
2. the association structure $\{w_{ij}\}_{1 \leq i < j \leq n}$.

Examples: Graph with Data Vectors

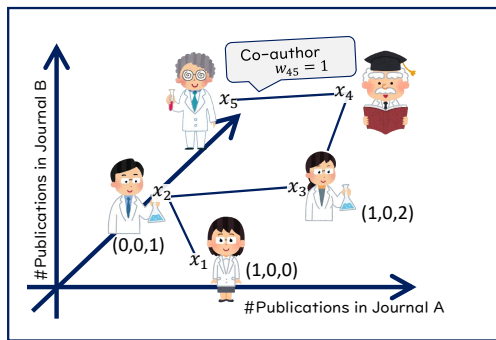


Figure: Co-authorship network

- $w_{ij} \in \{0, 1, 2, \dots\}$: number of co-authored papers
- $\mathbf{x}_i \in \mathbb{N}_0^p$: number of publications in each of p conferences

$$\mathbf{x}_i := \begin{pmatrix} \text{conf. A} & \text{conf. B} & \text{journal A} & \text{journal B} & \text{journal C} \\ 2 & 0 & 1 & 0 & 3, \dots \end{pmatrix} \in \mathbb{N}_0^p$$

Examples: Graph with Data Vectors

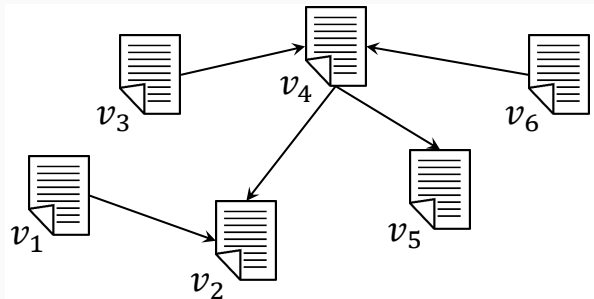


Figure: Citation network

- $w_{ij} \in \{0, 1\}$: whether i -th paper cites j -th paper³
- $\mathbf{x}_i \in \mathbb{N}_0^p$: bag-of-words (i.e., num. of each word) for each paper

$$\mathbf{x}_i = (\text{"the"}, \text{"he"}, \text{"she"}, \text{"animal"}, \text{"fish"}, \text{"dog"}, \dots) \in \mathbb{N}_0^{\#\text{dictionary}}$$

³directed relation can be transformed to undirected relation.

Examples: Graph with Data Vectors

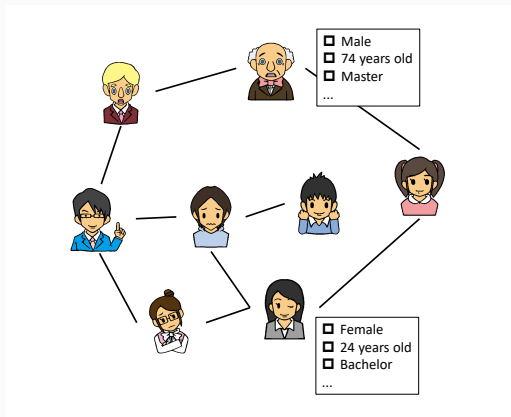


Figure: Friend network

- $w_{ij} \in \{0, 1\}$: friendship
- $\mathbf{x}_i \in \mathbb{R}^P$: property vector (e.g., age, born, degree, etc..)

Examples: Graph with Data Vectors

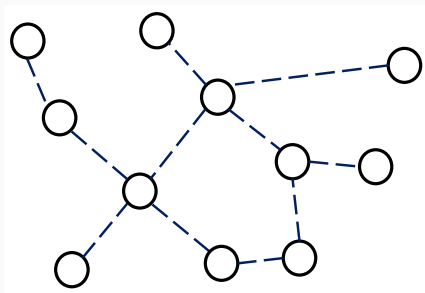


Figure: A **general graph**

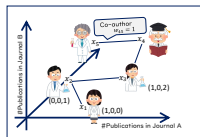
- $w_{ij} \geq 0$: Link weights
- **no data vector**: \triangleright 1-hot vector

$$\mathbf{x}_i := (0, \dots, 0, \underbrace{1}_{i\text{-th entry}}, 0, \dots, 0) \in \{0, 1\}^n$$

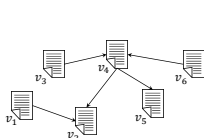
can be used instead. (or, $\mathbf{x}_i = \mathbf{w}_i := (w_{i1}, w_{i2}, \dots, w_{in})$ is used in recommendation systems)

Summary so far

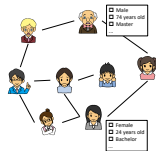
Given link weights $\{w_{ij}\}_{1 \leq i < j \leq n}$ and data vectors $\{\mathbf{x}_i\}_{i=1}^n$, e.g.,



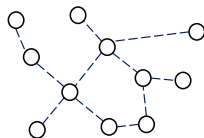
(a) Co-author network



(b) Citation network



(c) Friend network



(d) General graph

graph embedding learns a neural network $\mathbf{f} : \mathcal{X} \ni \mathbf{x}_i \mapsto \mathbf{y}_i \in \mathcal{Y}$, by maximizing the log-likelihood with a probabilistic model

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} \text{Bernoulli}(\sigma(\langle \mathbf{y}_i, \mathbf{y}_j \rangle)).$$

(▷ inner-product similarity has been used,

cf., Spectral: $\sum_{i,j=1}^n w_{ij} \langle \mathbf{y}_i, \mathbf{y}_j \rangle$, and Convolution: $w_{ij} \approx \sigma(\langle \mathbf{y}_i, \mathbf{y}_j \rangle).$)

PD and CPD similarities (AISTATS2019)

Poincaré Embedding

Whereas many of graph embedding trains a neural network $f : \mathcal{X} \ni \mathbf{x}_i \mapsto \mathbf{y}_i \in \mathcal{Y}$, by maximizing the log-likelihood with a probabilistic model

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} \text{Bernoulli}(\sigma(\langle \mathbf{y}_i, \mathbf{y}_j \rangle)),$$

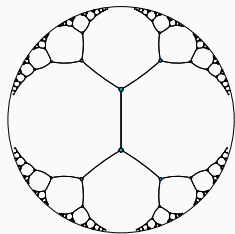
Poincaré embedding (Nickel and Kiela, 2017)⁴ considers

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} \text{Bernoulli}(\sigma(-d_{\text{Poincaré}}(\mathbf{y}_i, \mathbf{y}_j))),$$

where $\mathcal{Y} := \{\mathbf{y} \in \mathbb{R}^p \mid \|\mathbf{y}\|_2 < 1\}$,

$$d_{\text{Poincaré}}(\mathbf{y}_i, \mathbf{y}_j) := \cosh^{-1} \left(1 + 2 \frac{\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{(1 - \|\mathbf{y}_i\|_2^2)(1 - \|\mathbf{y}_j\|_2^2)} \right),$$

$$\cosh^{-1}(z) := \log(z + \sqrt{z^2 - 1}).$$



Tree (=graph) can be efficiently embedded with Poincaré distance. (Nickel and Kiela, 2017, Fig. 1(b))

⁴Lorentz embedding (Nickel and Kiela, 2018) utilizes Lorentz model instead, for efficient computation.

General Graph Embedding Framework

(Probabilistic) graph embedding can be generalized as follows:

General graph embedding

General graph embedding can be obtained by maximizing the log-likelihood of a model

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j \stackrel{\text{indep.}}{\sim} Q(\underbrace{\nu(g(\mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j)))}_{\text{(similarity)}}), \quad \mathbf{y}_i := \mathbf{f}(\mathbf{x}_i),$$

where

- $Q(\nu)$ is a **user-specified** distribution whose expectation is $\nu \in \mathbb{R}$,
 - $\nu : \mathbb{R} \rightarrow \mathbb{R}$ is a **user-specified** function,
 - $g : \mathcal{Y}^2 \rightarrow \mathbb{R}$ is a **user-specified** kernel (i.e., symmetric and continuous function), and
 - $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is a neural network **to be estimated**.
-
- Tang et al. (2015) Q : Bernoulli, ν : Sigmoid, g : inner-product.
 - Nickel and Kiela (2017) Q : Bernoulli, ν : Sigmoid, g : negative Poincaré.
 - Okuno et al. (2018) Q : Poisson, ν : Exponential, g : inner-product. etc...

“Expressive power” of graph embedding

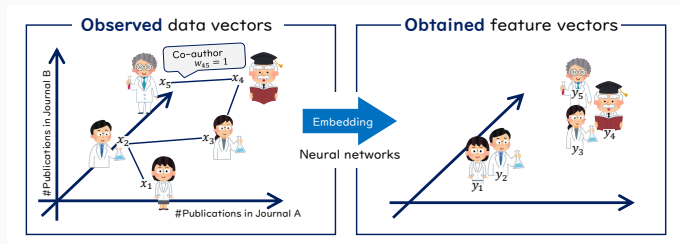


Figure: Graph embedding revisited.

Expressive power of graph embedding

⇔ Expressive power of the **similarity function**⁵

$$h(\mathbf{x}_i, \mathbf{x}_j) := g(\mathbf{y}_i, \mathbf{y}_j) = g(\mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j)),$$

equipped with the neural network \mathbf{f} and the user-specified kernel g .

⁵Similarity between two NNs $g(\mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j))$ is also known as *siamese network* (Bromley et al., 1994).

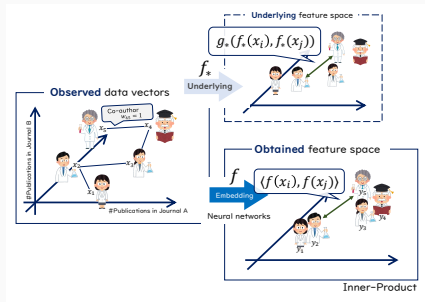
Inner-Product Similarity (IPS)

Okuno et al. (2018) Theorem 5.1, informal

$g_* : \mathcal{Y}_*^2 \rightarrow \mathbb{R}$ is a **PD** kernel ⁶ and $f_* : \mathcal{X} \rightarrow \mathcal{Y}_*$ is a continuous function. Then,

$$g_*(f_*(x_i), f_*(x_j)) \approx \langle f(x_i), f(x_j) \rangle,$$

for some sufficiently large neural network $f : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}^K$ with sufficiently large K .



⁶kernel (i.e., continuous and symmetric function) $g : \mathcal{Y}^2 \rightarrow \mathbb{R}$ satisfying $\sum_{i,j=1}^n c_i c_j g(y_i, y_j) \geq 0$ for all $\{y_i\}_{i=1}^n \subset \mathcal{Y}$, $\{c_i\}_{i=1}^n \subset \mathbb{R}$ is called positive-definite (PD).

Inner-Product Similarity (IPS)

Proof: Mercer's theorem (Minh et al., 2006) indicates that

$$g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j)) \stackrel{\text{Mercer}}{=} \langle \Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j) \rangle$$

for some Ψ . Considering the top- K entries Ψ_K , universal approximation theorem for NN (Cybenko, 1989) proves that $\Psi_K \approx \mathbf{f}$ for some NN \mathbf{f} ; thus

$$\langle \Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j) \rangle \approx \langle \Psi_K(\mathbf{x}_i), \Psi_K(\mathbf{x}_j) \rangle \approx \langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle.$$



The inner-product similarity $\langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle$ with neural network \mathbf{f} is

- 😊 capable of approximating any PD similarity $g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j))$.
- 😊 easy to compute.
- 😞 **limited** to approximate **PD** similarities.

- **Negative Poincaré distance** is not PD.

$$d_{\text{Poincaré}}(\mathbf{y}_i, \mathbf{y}_j) := \cosh^{-1} \left(1 + 2 \frac{\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{(1 - \|\mathbf{y}_i\|_2^2)(1 - \|\mathbf{y}_j\|_2^2)} \right),$$
$$\cosh^{-1}(z) := \log(z + \sqrt{z^2 - 1}).$$

- **Negative squared distance (NSD)** is not PD.

$$d_{\text{NSD}}(\mathbf{y}_i, \mathbf{y}_j) := -\|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

- **Negative Wasserstein distance** is not PD, etc...

- **Negative Poincaré distance** is not PD.

$$d_{\text{Poincaré}}(\mathbf{y}_i, \mathbf{y}_j) := \cosh^{-1} \left(1 + 2 \frac{\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{(1 - \|\mathbf{y}_i\|_2^2)(1 - \|\mathbf{y}_j\|_2^2)} \right),$$
$$\cosh^{-1}(z) := \log(z + \sqrt{z^2 - 1}).$$

- **Negative squared distance (NSD)** is not PD.

$$d_{\text{NSD}}(\mathbf{y}_i, \mathbf{y}_j) := -\|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

- **Negative Wasserstein distance** is not PD, etc...

Okuno et al. (2019) Proposition 4.1

For all $M > 0, p, K \in \mathbb{N}$, and for any continuous functions $\mathbf{f} : \mathbb{R}^p \rightarrow \mathbb{R}^K$,

$$\frac{1}{(2M)^{2p}} \iint_{[-M, M]^{2p}} \underbrace{\left| -\|\mathbf{x} - \mathbf{x}'\|_2^2 - \langle \mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}') \rangle \right|}_{\text{non-PD}} d\mathbf{x}d\mathbf{x}' \geq \frac{2pM^2}{3}.$$

Shifted Inner-Product Similarity (SIPS)

Conditionally PD (CPD) (Berg et al., 1984)

A symmetric function $h : \mathcal{X}^2 \rightarrow \mathbb{R}$ is called **conditionally-PD (CPD)** if $\sum_{i,j=1}^n c_i c_j h(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ for all $\{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{X}$, $\{c_i\}_{i=1}^n \subset \mathbb{R}$ satisfying $\sum_{i=1}^n c_i = 0$.

Conditionally PD

Negative-squared dist.,
Negative-Poincare dist.,
Negative-Earth-mover's,
Negative-sliced-Wasserstein,
...

PD

Linear kernel,
Polynomial kernel,
Gaussian kernel,
Cosine similarity, ...

Shifted Inner-Product Similarity (SIPS)

Shifted Inner-Product Similarity (SIPS) (Okuno et al., 2019)

Proposed SIPS is defined as

$$\langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle + u(\mathbf{x}_i) + u(\mathbf{x}_j),$$

where $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ and $u : \mathcal{X} \rightarrow \mathbb{R}$ are neural networks.

SIPS is capable of approximating any CPD similarities as follows.

Okuno et al. (2019) Theorem 4.1, informal

$g_* : \mathcal{Y}_*^2 \rightarrow \mathbb{R}$ is a **CPD** kernel and $\mathbf{f}_* : \mathcal{X} \rightarrow \mathcal{Y}_*$ is a continuous function. Then,

$$g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j)) \approx \langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle + u(\mathbf{x}_i) + u(\mathbf{x}_j),$$

for some sufficiently large neural networks $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}^K$, $u : \mathcal{X} \rightarrow \mathbb{R}$ with sufficiently large K .

Shifted Inner-Product Similarity (SIPS)

Proof is based on the following Lemma.

Berg et al. (1984) Lemma 2.1

$g_* : \mathcal{Y}^2 \rightarrow \mathbb{R}$ is CPD, if and only if

$$\tilde{g}_*(\mathbf{y}, \mathbf{y}') := g_*(\mathbf{y}, \mathbf{y}') - g_*(\mathbf{y}, \mathbf{y}_0) - g_*(\mathbf{y}_0, \mathbf{y}') + g_*(\mathbf{y}_0, \mathbf{y}_0)$$

with a $\mathbf{y}_0 \in \mathcal{Y}$ is PD.

For any CPD kernel g_* , we have

$$g_*(\mathbf{y}, \mathbf{y}') = \tilde{g}_*(\mathbf{y}, \mathbf{y}') + \underbrace{\left(g_*(\mathbf{y}, \mathbf{y}_0) - \frac{1}{2}g_*(\mathbf{y}_0, \mathbf{y}_0) \right)}_{=:r(\mathbf{y})} + \underbrace{\left(g_*(\mathbf{y}_0, \mathbf{y}') - \frac{1}{2}g_*(\mathbf{y}_0, \mathbf{y}_0) \right)}_{=:r(\mathbf{y}')}. \quad \square$$

Substituting $\mathbf{y}_i := \mathbf{f}_*(\mathbf{x}_i)$ leads to

$$g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j)) = \underbrace{\tilde{g}_*(\mathbf{f}_*(\mathbf{x}), \mathbf{f}_*(\mathbf{x}'))}_{\text{PD}} + r(\mathbf{f}_*(\mathbf{x})) + r(\mathbf{f}_*(\mathbf{x}')).$$

Approximating the terms in right-hand-side by $\langle \mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}') \rangle$, $u(\mathbf{x})$, $u(\mathbf{x}')$, respectively, lead to the assertion. □

Numerical Experiments (Okuno and Shimodaira, 2018)

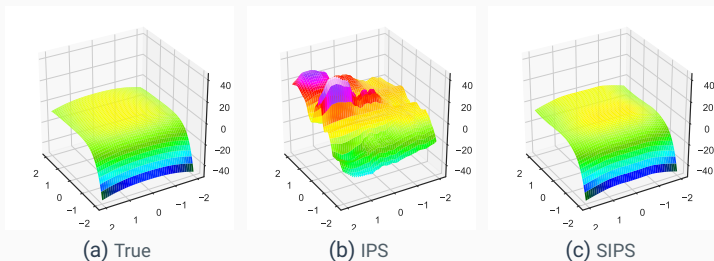


Figure: SIPS can but IPS cannot approximate negative-squared distance (CPD).

For $\mathbf{f}_*(\mathbf{x}) = (x_1, \cos x_2, \exp(-x_3), \sin(x_4 - x_5)) \in \mathbb{R}^4$ w.r.t. $\mathbf{x} \in \mathbb{R}^5$,

(a) $-\|\mathbf{f}_*(s\mathbf{e}_1) - \mathbf{f}_*(t\mathbf{e}_2)\|_2^2$ (CPD),

(b) trained IPS $\langle \mathbf{f}(s\mathbf{e}_1), \mathbf{f}(t\mathbf{e}_2) \rangle$,

(c) trained SIPS $\langle \mathbf{f}(s\mathbf{e}_1), \mathbf{f}(t\mathbf{e}_2) \rangle + u(s\mathbf{e}_1) + u(t\mathbf{e}_2)$,

are plotted on (s, t) -plane along with two orthogonal directions $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{R}^5$.

$\mathbf{f} : \mathbb{R}^5 \rightarrow \mathbb{R}^{10}$ and $u : \mathbb{R}^5 \rightarrow \mathbb{R}$ are two-layer neural networks with 1,000 hidden units and ReLU activation.

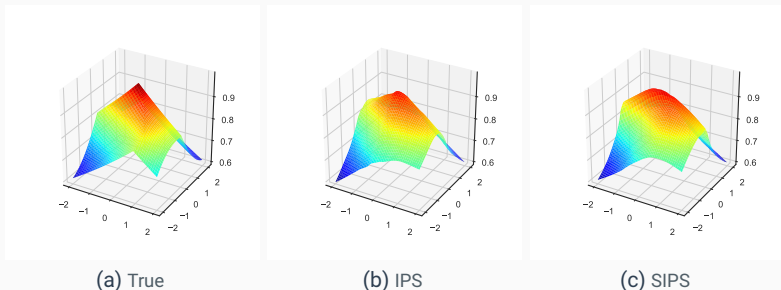


Figure: Both of IPS and SIPS can approximate Laplace kernel (PD).

As well,

(a) Laplace kernel $\exp(-\|\mathbf{f}_*(\mathbf{se}_1) - \mathbf{f}_*(\mathbf{te}_2)\|_1)$ (PD),

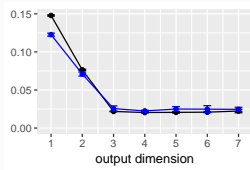
(b) trained IPS $\langle \mathbf{f}(\mathbf{se}_1), \mathbf{f}(\mathbf{te}_2) \rangle$,

(c) trained SIPS $\langle \mathbf{f}(\mathbf{se}_1), \mathbf{f}(\mathbf{te}_2) \rangle + u(\mathbf{se}_1) + u(\mathbf{te}_2)$,

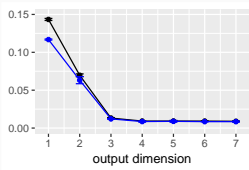
are plotted on (s, t) -plane.

Numerical Experiments (Okuno and Shimodaira, 2018)

Prediction Mean Squared Error (PMSE) for approximating similarities. IPS (Black) and SIPS (Blue).

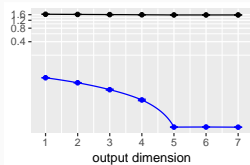


(a) #hidden units=100

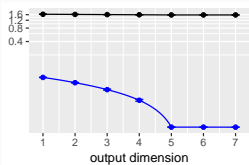


(b) #hidden units=1000

Figure: Cosine similarity (PD)



(a) #hidden units=100



(b) #hidden units=1000

Figure: negative Poincaré distance (CPD)

Approximation Error Rate

Okuno et al. (2019) Theorem 5.2 with $\alpha = 1$, informal

For any

- compact sets $\mathcal{X} \subset \mathbb{R}^p, \mathcal{Y}_* \subset \mathbb{R}^{K_*}$,
- CPD kernel $g_* \in \mathcal{C}^1(\mathcal{Y}_*^2; \mathbb{R})$, and
- $\mathbf{f}_* \in \mathcal{C}^1(\mathcal{X}; \mathcal{Y}_*)$,

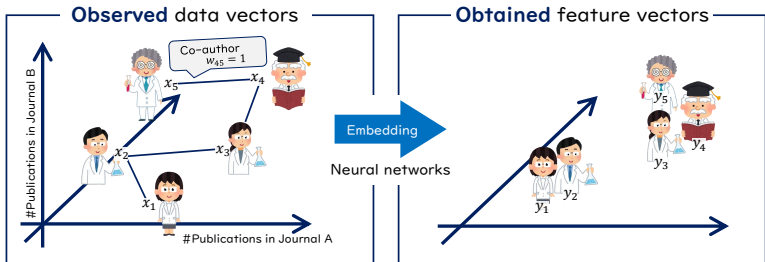
SIPS with ReLU-activated deep MLPs $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^K, u : \mathcal{X} \rightarrow \mathbb{R}$ whose depths are $O(L_f), O(L_u)$ and widths are $O(1)$, uniformly approximates $g_*(\mathbf{f}_*(\mathbf{x}), \mathbf{f}_*(\mathbf{x}'))$ with an error

$$O \left(\underbrace{K^{-1/K_*}}_{(*)} + K^{1/2} L_f^{-2/p} + L_u^{-2/p} \right).$$

Proof is straightforwardly obtained by applying Cobos and Kühn (1990) and Yarotsky (2018). See Okuno et al. (2019) Supplement D.3 with $\alpha = 1$. □

Note: we showed only the comprehensible error rate; the above error rate can be improved by employing other existing studies.

Summary of this study



Graph embedding learns a transformation $f : \mathcal{X} \rightarrow \mathcal{Y}$, so that

$$w_{ij} \approx \text{Sigmoid}(\underbrace{g(\mathbf{y}_i, \mathbf{y}_j)}_{\text{"Similarity"}}), \quad \mathbf{y}_i := f(\mathbf{x}_i).$$

g is specified as inner-product (Tang et al., 2015), Poincaré-dist. (Nickel and Kiela, 2017), etc..

Question \triangleright What kind of kernel g is good (in terms of the expressive power?)

Answer \triangleright Proposed **shifted inner-product similarity (SIPS)**.

SIPS can approximate many kernels, e.g., Poincaré-distance!

Recent Progress: General Similarities (arXiv:1902.10409)

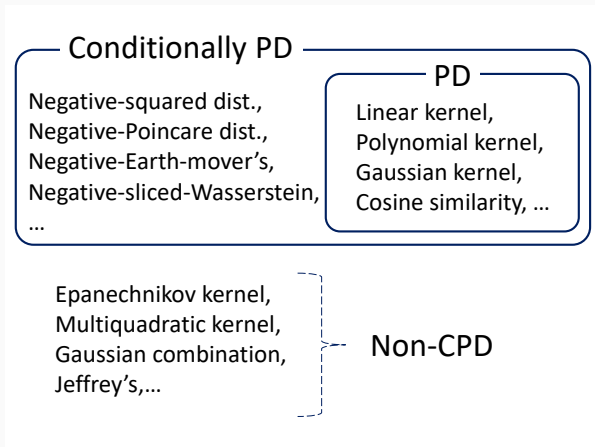


Figure: There remain non-CPD similarities, that cannot be approximated by SIPS.

Inner-Product Difference Similarity (IPDS)

General similarity $g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j))$ can be decomposed as

$$g_*(\mathbf{f}_*(\mathbf{x}_i), \mathbf{f}_*(\mathbf{x}_j)) = \underbrace{h_+(\mathbf{x}_i, \mathbf{x}_j)}_{\text{PD}} - \underbrace{h_-(\mathbf{x}_i, \mathbf{x}_j)}_{\text{PD}},$$

and PD similarities h_+, h_- can be approximated by IPS.

Inner-Product Difference Similarity (IPDS) (Okuno et al., 2019, Supplement E)

IPDS is defined as

$$\langle \mathbf{f}^+(\mathbf{x}_i), \mathbf{f}^+(\mathbf{x}_j) \rangle - \langle \mathbf{f}^-(\mathbf{x}_i), \mathbf{f}^-(\mathbf{x}_j) \rangle,$$

where $\mathbf{f}^+ : \mathcal{X} \rightarrow \mathcal{Y}^+ \subset \mathbb{R}^{K_+}$, $\mathbf{f}^- : \mathcal{X} \rightarrow \mathcal{Y}^- \subset \mathbb{R}^{K_-}$ are neural networks.

IPDS is capable of approximating general similarities⁷.

- 😊 IPDS includes SIPS as a special case⁸.
- 😞 requires to specify the rate K_+/K_- .

⁷ more precisely, similarities dominated by PD similarity. see Ong et al. (2004)

⁸ $\mathbf{f}^+(\mathbf{x}) := (\mathbf{f}(\mathbf{x})^\top, u(\mathbf{x}), 1)^\top$, $\mathbf{f}^-(\mathbf{x}) := u(\mathbf{x}) - 1$.

Weighted Inner-Product Similarity (WIPS)

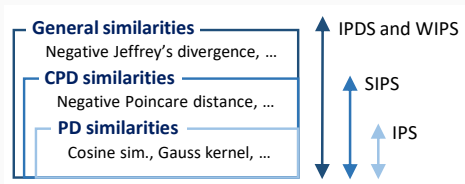
Weighted Inner-Product Similarity (WIPS) (Kim et al., 2019)

WIPS is defined as

$$\langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle_{\lambda},$$

where $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is a neural network, $\langle \mathbf{y}, \mathbf{y}' \rangle_{\lambda} = \sum_{k=1}^K \lambda_k y_k y'_k$, and $\lambda = (\lambda_1, \dots, \lambda_K) \in \mathbb{R}^K$ is a weight vector **to be estimated**.

😊 WIPS using $\lambda := (1, 1, \dots, 1, -1, -1, \dots, -1)$ reduces to IPDS.



😊 no need to specify the rate K_+/K_- .

Numerical Experiments (Kim et al., 2019)

In this numerical experiment, we embed **WebKB hypertext network**⁹ equipped with 1,409-dim. Bag-of-Words data vectors.

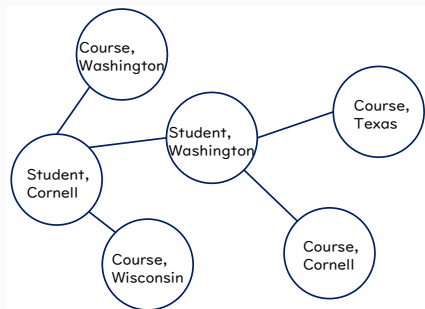


Figure: 877 hypertexts (nodes) are associated by 1,480 hyperlinks (links).

Each node has semantic class labels

- one of {*Student, Faculty, Staff, Course, Project*},
- one of university labels in {*Cornell, Texas, Washington, Wisconsin*}.

⁹<https://linqs.soe.ucsc.edu/data>

Numerical Experiments (Kim et al., 2019)

Obtained 10-dim. feature vectors are mapped to \mathbb{R}^2 with *t*-SNE (Maaten and Hinton, 2008).

Semantic labels are colored as

- Student (navy), Course (pink),
- Cornell (red), Texas (orange), Washington (green) and Wisconsin (blue).

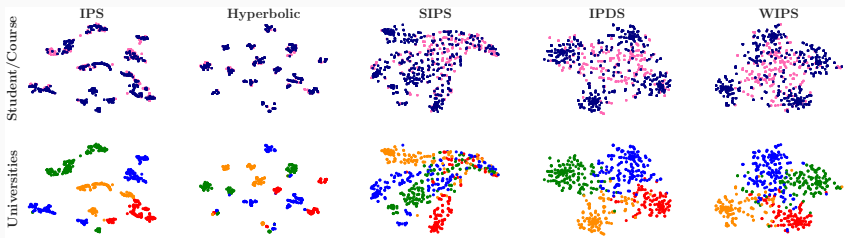


Figure: Selected nodes of embeded WebKB hypertext network using *t*-SNE.

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

Numerical Experiments (Kim et al., 2019)

We conducted following experiments on **WebKB hypertext** network.

- **Reconstruction** of link weights from fully observed data,
- **Link prediction** of unseen nodes.

Table: **ROC-AUC**. Boldface is the best, and underlines are 2nd and 3rd best scores.

		Dimensionality					
		Reconstruction			Link prediction		
		10	50	100	10	50	100
Hypertext	IPS	91.99	94.23	94.24	77.73	77.62	77.16
	Poincaré	94.09	94.13	94.11	<u>82.21</u>	79.64	79.48
	SIPS	<u>95.11</u>	<u>95.12</u>	<u>95.12</u>	82.01	<u>81.84</u>	<u>81.13</u>
	IPDS	95.12	<u>95.12</u>	95.12	82.59	82.75	<u>82.19</u>
	WIPS	<u>95.11</u>	95.12	<u>95.12</u>	<u>82.38</u>	<u>82.68</u>	82.93

Numerical Experiments (Kim et al., 2019)

Same experiments are conducted on

- **DBLP Co-authorship network** (Prado et al., 2013)
 - ▷ 41, 328 nodes, 210, 320 links, 33-dim. data vectors.
- **WordNet Taxonomy Tree** (Nickel and Kiela, 2017)
 - ▷ 37, 623 nodes, 312, 885 links, 300-dim. pre-trained word vectors.

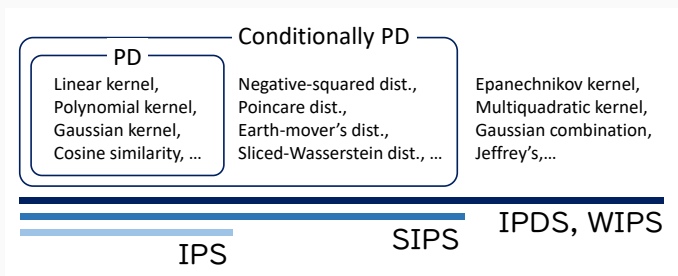
Table: **ROC-AUC**. Boldface is the best, and underlines are 2nd and 3rd best scores.

		Dimensionality					
		Reconstruction			Link prediction		
		10	50	100	10	50	100
Co-author	IPS	85.01	86.02	85.80	83.83	84.41	84.02
	Poincaré	86.84	86.69	86.72	85.82	85.92	85.93
	SIPS	<u>90.01</u>	<u>91.35</u>	<u>91.06</u>	<u>88.24</u>	<u>88.69</u>	<u>88.67</u>
	IPDS	<u>90.13</u>	<u>91.68</u>	<u>91.59</u>	88.42	<u>88.97</u>	<u>88.85</u>
	WIPS	90.50	92.44	92.95	<u>88.16</u>	89.43	89.40
Taxonomy	IPS	79.95	75.80	74.97	67.25	65.71	65.38
	Poincaré	91.69	89.10	88.97	83.04	79.52	78.97
	SIPS	<u>98.78</u>	<u>99.75</u>	<u>99.77</u>	<u>90.42</u>	<u>92.12</u>	<u>92.09</u>
	IPDS	99.65	99.89	99.90	95.99	96.37	<u>96.41</u>
	WIPS	<u>99.64</u>	<u>99.85</u>	<u>99.87</u>	<u>95.07</u>	<u>96.36</u>	96.51

Conclusion

Conclusion

- (1) As the similarity g , we propose **SIPS** and **IPDS** in Okuno et al. (2019), that are highly expressive compared with the simple inner-product similarity.
- (2) We prove that **SIPS approximates any CPD similarities**, and IPDS approximates general similarities.
- (3) We also propose **WIPS** in Kim et al. (2019), that generalizes IPDS.



- (4) Experiments show the high expressive power of SIPS, IPDS, and WIPS.

- Belkin, M. and Niyogi, P. (2003). Laplacian Eigenmaps for Dimensionality Reduction and Data Representation. *Neural Computation*, 15(6):1373–1396.
- Berg, C., Christensen, J. P. R., and Ressel, P. (1984). Harmonic Analysis on Semigroups.
- Bojchevski, A. and Günnemann, S. (2018). Deep gaussian embedding of attributed graphs: Unsupervised inductive learning via ranking. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., and Shah, R. (1994). Signature verification using a “siamese” time delay neural network. In *Advances in Neural Information Processing Systems*, pages 737–744.
- Chung, F. R. (1997). *Spectral Graph Theory*. Number 92. American Mathematical Society.
- Cobos, F. and Kühn, T. (1990). Eigenvalues of Integral Operators with Positive Definite Kernels Satisfying Integrated Hölder Conditions over Metric Compacta. *Journal of Approximation Theory*, 63(1):39–55.
- Cybenko, G. (1989). Approximation by Superpositions of a Sigmoidal Function. *Mathematics of control, signals and systems*, 2(4):303–314.

- Hamilton, W., Ying, Z., and Leskovec, J. (2017). Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034.
- He, X. and Niyogi, P. (2004). Locality Preserving Projections. In *Advances in Neural Information Processing Systems*, pages 153–160.
- Kim, G., Okuno, A., Fukui, K., and Shimodaira, H. (2019). Representation Learning with Weighted Inner Product for Universal Approximation of General Similarities. *arXiv preprint arXiv:1902.10409*. submitted.
- Kipf, T. N. and Welling, M. (2016). Variational Graph Auto-Encoders. *NIPS Workshop*.
- Kipf, T. N. and Welling, M. (2017). Semi-Supervised Classification with Graph Convolutional Networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605.

- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.
- Minh, H. Q., Niyogi, P., and Yao, Y. (2006). Mercer's Theorem, Feature Maps, and Smoothing. In *International Conference on Computational Learning Theory*, pages 154–168. Springer.
- Nickel, M. and Kiela, D. (2017). Poincaré Embeddings for Learning Hierarchical Representations. In *Advances in Neural Information Processing Systems*, pages 6338–6347.
- Nickel, M. and Kiela, D. (2018). Learning Continuous Hierarchies in the Lorentz Model of Hyperbolic Geometry. In *International Conference on Machine Learning*, pages 3776–3785.
- Okuno, A., Hada, T., and Shimodaira, H. (2018). A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3885–3894.

- Okuno, A., Kim, G., and Shimodaira, H. (2019). Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability. In *Proceedings of International Conference on Artificial Intelligence and Statistics*. to appear.
- Okuno, A. and Shimodaira, H. (2018). On representation power of neural-network based graph embedding and beyond. In *ICML2018 workshop on Theoretical Foundations and Applications of Deep Generative Models (TADGM)*.
- Ong, C. S., Mary, X., Canu, S., and Smola, A. J. (2004). Learning with non-positive kernels. In *Proceedings of the International Conference on Machine Learning*, page 81. ACM.
- Prado, A., Plantevit, M., Robardet, C., and Boulicaut, J.-F. (2013). Mining graph topological patterns: Finding covariations among vertex descriptors. *IEEE Transactions on Knowledge and Data Engineering*, 25(9):2090–2104.
- Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J., and Mei, Q. (2015). LINE: Large-scale information network embedding. In *Proceedings of the 24th International Conference on World Wide Web*, pages 1067–1077.

Vilnis, L. and McCallum, A. (2015). Word representations via gaussian embedding. In *Proceedings of the International Conference on Learning Representations (ICLR)*.

Yarotsky, D. (2018). Optimal approximation of continuous functions by very deep ReLU networks. In *Conference On Learning Theory*, pages 639–649.

Shifted Inner-Product Similarity (SIPS)

If $\nu(\cdot) = \exp(\cdot)$, using SIPS incorporates weights into IPS-based model, as

$$\underbrace{\exp(\langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle + u(\mathbf{x}_i) + u(\mathbf{x}_j))}_{\text{SIPS-based model}} = \beta(\mathbf{x}_i) \cdot \beta(\mathbf{x}_j) \cdot \underbrace{\exp(\langle \mathbf{f}(\mathbf{x}_i), \mathbf{f}(\mathbf{x}_j) \rangle)}_{\text{IPS-based model}},$$

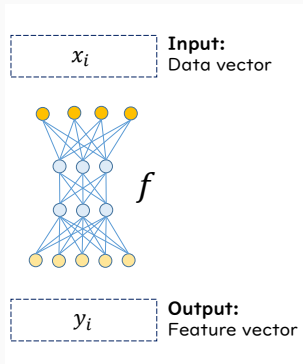
where $\beta(\mathbf{x}) := \exp(u(\mathbf{x}))$ is **weight function**.

In practice, we utilize $\tilde{\mathbf{f}}(\mathbf{x}) := (\mathbf{f}(\mathbf{x})^\top, u(\mathbf{x}))^\top$ for subsequent statistical analyses (e.g. clustering, visualization,...).

- 😊 SIPS is capable of **approximating** any **CPD** similarities.
- 😊 SIPS-based model naturally **incorporates weight function** into IPS-based model.
- 😞 \mathbf{f}, u in SIPS are unidentifiable (**hard to train**).

Proper regularization is required, as a future research.

Neural network architecture is specified as follows.



- We employ **1-hidden layer multi-layer perceptron** with 2,000 hidden Units and ReLU activation, for each neural networks.
- For IPDS, K_+/K_- is grid-searched over $\{0.01, 0.25, 0.5, 0.75, 0.9\}$.
- For WIPS, λ is uniformly randomly initialized in $(0, 1/K)^K$.

For details, see Kim et al. (2019) Section 6.

Gaussian Embedding

Gaussian embedding (Vilnis and McCallum, 2015) can be regarded as learning a neural network¹⁰

$$\mathbf{y}_i = N_K(\underbrace{\mu(\mathbf{x}_i), \Sigma(\mathbf{x}_i)}_{\text{Neural network}}),$$

where N_K represents K -dim. Normal dist., so that

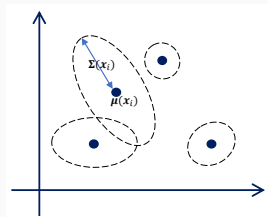
$$w_{ij} \approx \sigma(-D_{\text{KL}}(\mathbf{y}_i, \mathbf{y}_j)).$$

😊 $\Sigma(\mathbf{x}_i)$ represents the **uncertainty** of the embedding (Vilnis and McCallum, 2015).

D_{KL} may be replaced with symmetric Jeffrey's divergence

$$D_{\text{Jeff.}}(\mathbf{y}, \mathbf{y}') := D_{\text{KL}}(\mathbf{y}_i, \mathbf{y}_j) + D_{\text{KL}}(\mathbf{y}_j, \mathbf{y}_i),$$

then **Jeffrey's divergence is not CPD** (Okuno et al., 2019, Supplement E.3).



Gaussian embedding also captures the **uncertainty** $\Sigma(\mathbf{x}_i)$ of the point embedding $\mu(\mathbf{x}_i)$.

¹⁰ Vilnis and McCallum (2015) corresponds to using 1-hot vectors (as no data vectors are available), and Deep Gaussian embedding (Bojchevski and Günnemann, 2018) incorporates neural network into Gaussian embedding.