

# KDD 2020 paper reading 10/15

Houquan Zhou

October 15, 2020

- **Title:** Data Compression as a Comprehensive Framework for Graph Drawing and Representation Learning
- **Authors:** Claudia Plant, Sonja Biedermann, Christian Böhm
- **Affiliations:** Univerity of Vienna, LMU Munich

## Data Compression as a Comprehensive Framework for Graph Drawing and Representation Learning

Claudia Plant

Faculty of Computer Science, ds:UniVie  
University of Vienna, Vienna, Austria  
claudia.plant@univie.ac.at

Sonja Biedermann

Faculty of Computer Science  
University of Vienna, Vienna, Austria  
sonja.biedermann@univie.ac.at

Christian Böhm

Institute of Informatics, MCML  
LMU Munich, Germany  
boehm@ifi.lmu.de

# Overview

- **Problem:** Representation learning.
- **Idea:** A novel objective function from a **data compression** perspective: **Predictive Entropy (PE)**.
- **Algorithm:** A novel representation learning algorithm GEMPE.
- **Experiments:** Graph drawing & Representation learning.

# Predictive Entropy

**Embeddings quality:** The more effectively the low-dimensional embeddings allow to **compress the adjacency matrix**, the better is the quality of the embedding.

Based on the idea, the authors proposed a novel information-theoretic evaluation function: **Predictive Entropy** (PE).

PE is based on a probabilistic model, which predicts the probability  $\Pr((i, j) \in E)$  from their embeddings  $x_i, x_j$ .

Ideal case:  $\Pr((i, j) \in E)$  close to 1 for  $(i, j) \in E$ , and close to 0 for  $(i, j) \notin E$ .

# Probabilistic model

Edge probability of  $(i, j)$  related to the Euclidean distance of  $x_i$  and  $x_j$ .  
The smaller  $\|x_i - x_j\|$  is, the more possible  $i$  and  $j$  are connected.

$$s_{\mu, \sigma}^+(\delta) := P((i, j) \in E \mid \|x_i - x_j\| = \delta) = \frac{1}{2} - \frac{1}{2} \operatorname{erf}\left(\frac{\delta - \mu}{\sqrt{2}\sigma}\right) \quad (1)$$

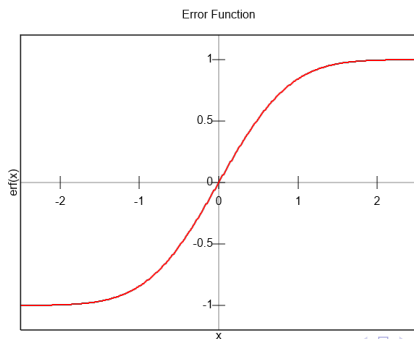
$$s_{\mu, \sigma}^-(\delta) := P((i, j) \notin E \mid \|x_i - x_j\| = \delta) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{\delta - \mu}{\sqrt{2}\sigma}\right) \quad (2)$$

## erf function

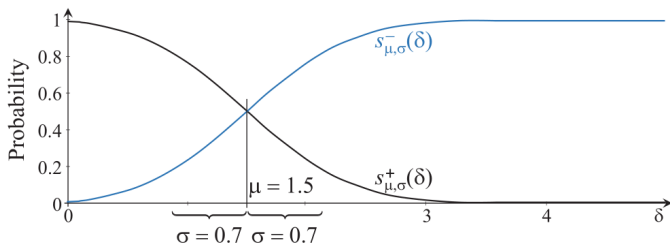
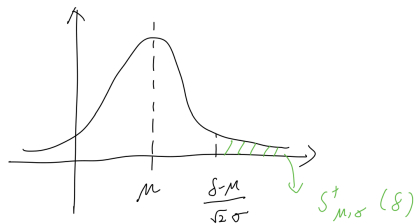
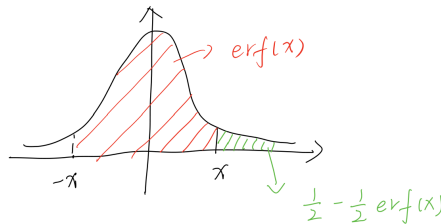
$$\operatorname{erf}(x) = \frac{1}{\sqrt{\pi}} \int_{-x}^x e^{-t^2} dt = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \quad (3)$$

Related to cumulative probability function of Gaussian distribution:

$$\Phi(x) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \quad (4)$$



# Demonstration



# Predictive Entropy

Encoding length of  $(i, j)$  entry of the adjacency matrix:

$$dl_{\mu,\sigma}(i, j) = \begin{cases} dl_{\mu,\sigma}^+(\delta) & \text{if } (i, j) \in E \\ dl_{\mu,\sigma}^-(\delta) & \text{otherwise} \end{cases}$$
$$dl_{\mu,\sigma}^+(\delta) = -\log \left( \frac{1}{2} - \frac{1}{2} \operatorname{erf} \left( \frac{\delta - \mu}{\sqrt{2}\sigma} \right) \right)$$
$$dl_{\mu,\sigma}^-(\delta) = -\log \left( \frac{1}{2} + \frac{1}{2} \operatorname{erf} \left( \frac{\delta - \mu}{\sqrt{2}\sigma} \right) \right)$$

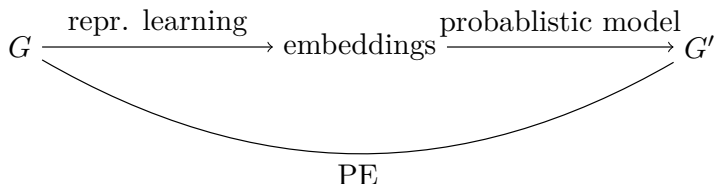
Total objective function  $PE$ :

$$PE_G(x_1, \dots, x_n) = \frac{1}{N} \min_{\mu, \sigma} \sum_{1 \leq i \leq n} \sum_{i < j \leq n} dl_{\mu, \sigma}(i, j), \quad (5)$$

where  $N = \frac{n(n+1)}{2}$ ,  $x_1, \dots, x_n$  are embeddings of nodes.



# Summary: PE



- Learn embedding from  $G$ .
- Use embeddings to predict the edges.
- Encoding cost  $\Rightarrow$  PE.

# Weighted Majorization<sup>1</sup>

## Definition (Weighted Majorization)

- Given  $n$  objects, distance matrix  $d_{i,j}$  and weights  $w_{i,j}$ ;
- Find embeddings  $x_i, i = 1, 2, \dots, n$ ;
- To minimize the following objective:

$$wm(x_1, \dots, x_n) = \sum_{1 \leq i \leq n} \sum_{i < j \leq n} w_{i,j} (\|x_i - x_j\| - d_{i,j})^2 \quad (6)$$

WM aims at finding low-dimensional coordinates  $x_i$  that minimize the squared difference between the distance in the input matrix  $d_{i,j}$  and the Euclidean distance  $\|x_i - x_j\|$  of the generated vectors.

---

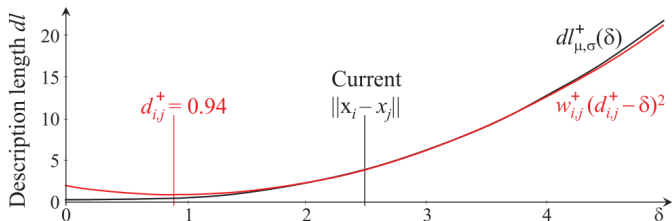
<sup>1</sup>Brandes and Pich 2009; Gansner, Koren, and North 2004.

# Optimizing PE using WM

Focusing on one pair  $(i, j)$ :  $w_{i,j}(\|x_i - x_j\| - d_{i,j})^2$ .

We want the WM objective function to approximate our objective function.  $\Rightarrow$  Two-order approximation.

$$\frac{d}{d\delta} dl_{\mu,\sigma}^+(\delta) = \frac{d}{d\delta} w_{i,j}(\delta - d_{i,j})^2$$
$$\frac{d^2}{d\delta^2} dl_{\mu,\sigma}^+(\delta) = \frac{d^2}{d\delta^2} w_{i,j}(\delta - d_{i,j})^2$$



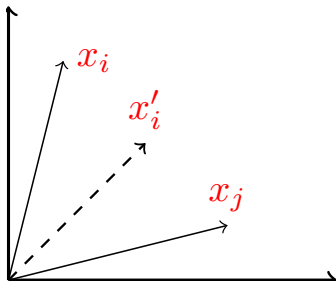
# Optimizing PE using WM

$$w_{i,j}^+ = \frac{\frac{2}{\pi \ln 2} e^{-\frac{(\delta-\mu)^2}{\sigma^2}}}{\sigma^2 (1 - \operatorname{erf}(\frac{\delta-\mu}{\sqrt{2}\sigma}))^2} - \frac{\frac{1}{\sqrt{2\pi} \ln 2} (\delta - \mu) e^{-\frac{(\delta-\mu)^2}{2\sigma^2}}}{\sigma^3 (1 - \operatorname{erf}(\frac{\delta-\mu}{\sqrt{2}\sigma}))}$$
$$d_{i,j}^+ = \delta - \frac{\frac{1}{\sqrt{2\pi} \ln 2} e^{-\frac{(\delta-\mu)^2}{2\sigma^2}}}{w_{i,j} \sigma (1 - \operatorname{erf}(\frac{\delta_{i,j}-\mu}{\sqrt{2}\sigma}))}$$

# Update embeddings

$$x_i := \frac{\sum_{j \neq i} w_{i,j} (x_j + s_{i,j} (x_i - x_j))}{\sum_{j \neq i} w_{i,j}}, \quad (7)$$

where  $s_{i,j} = \begin{cases} d_{i,j} / \|x_i - x_j\| & \text{if } \|x_i - x_j\| \neq 0 \\ 0 & \text{otherwise} \end{cases}$ .



- $s_{i,j} = 0 \Rightarrow d_{i,j} = 0$   
 $x_i \leftarrow x_j$
- $s_{i,j} = 1 \Rightarrow d_{i,j} = \|x_i - x_j\|$   
 $x_i \leftarrow x_i$
- $s_{i,j} = \frac{1}{2} \Rightarrow d_{i,j} = \frac{1}{2} \|x_i - x_j\|$   
 $x_i \leftarrow \frac{x_i + x_j}{2}$

# Finding $\mu$ and $\sigma$

Meaning of  $\mu$ : Boundary of Euclidean distance.

In the experiment,  $\mu$  is fixed to a value between 1 and 2.

The optimal  $\sigma$  is found by a simple bisection search.

# Algorithm

**algorithm** GEMPE ( $V, E, d$ )  $\rightarrow \mathbb{R}^{n \times d}$

$\forall i \in V$ : initialize  $\mathbf{x}_i \in \mathbb{R}^d$  randomly (uniform distribution);

**repeat**

determine  $\sigma$  according to Section 3.2;  $\rightarrow$  每一轮中,找最优的  $\sigma$ .

分子  $\hookrightarrow \forall i \in V: \mathbf{y}_i := (0, \dots, 0)^T \in \mathbb{R}^d$ ; // numerators and...

$\forall i \in V: \mathbf{z}_i := 0 \in \mathbb{R}$ ; // ...denominators of new points

$n$  是固定的

(\*)  $\checkmark$  **foreach**  $(i, j) \in E$  **do**

分母

determine parabola  $(d_{i,j}^+, w_{i,j}^+)$  acc. to Eq. (3) and (4);

weightedMajorizationStep  $(i, j)$ ;

// negative sampling:  $\star$

sample  $k \in V$  such that  $k \neq i$  and  $(k, i) \notin E$ ;

determine parabola  $(d_{k,i}^-, w_{k,i}^-)$  acc. to Eq. (5) and (6);

正负采样

weightedMajorizationStep  $(k, i)$ ;

sample  $k \in V$  such that  $k \neq j$  and  $(k, j) \notin E$ ;

负采样

determine parabola  $(d_{k,j}^-, w_{k,j}^-)$  acc. to Eq. (5) and (6);

weightedMajorizationStep  $(k, j)$ ;

$\forall i \in V: \mathbf{x}_i := \frac{1}{z_i} \cdot \mathbf{y}_i$ ;

更新.

**until convergence**;

**return**  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ;

# Algorithm

```
procedure weightedMajorizationStep( $a, b$ )  
  // one step of weighted majorization for  $\mathbf{x}_a$  and  $\mathbf{x}_b$ :  
   $\mathbf{y}_a := \mathbf{y}_a + \mathbf{w}_{a,b} \cdot \mathbf{x}_b$ ;  
   $\mathbf{y}_b := \mathbf{y}_b + \mathbf{w}_{a,b} \cdot \mathbf{x}_a$ ;  
   $\mathbf{z}_a := \mathbf{z}_a + \mathbf{w}_{a,b}$ ;  
   $\mathbf{z}_b := \mathbf{z}_b + \mathbf{w}_{a,b}$ ;  
  if  $\|\mathbf{x}_a - \mathbf{x}_b\| \neq 0$  then  
     $\mathbf{y}_a := \mathbf{y}_a + \mathbf{w}_{a,b} \cdot \frac{d_{a,b}}{\|\mathbf{x}_a - \mathbf{x}_b\|} \cdot (\mathbf{x}_a - \mathbf{x}_b)$ ;  
     $\mathbf{y}_b := \mathbf{y}_b + \mathbf{w}_{a,b} \cdot \frac{d_{a,b}}{\|\mathbf{x}_a - \mathbf{x}_b\|} \cdot (\mathbf{x}_b - \mathbf{x}_a)$ ;
```



# Experiments: Repr. Learning

Test embedding quality on downstream task: node classification.

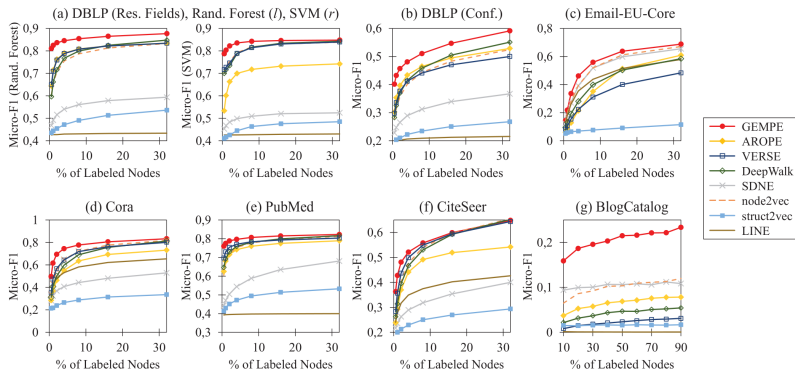


Figure 4: Comparison to Representation Learning Methods: Node Classification.

# Experiments: Clustering

	DBLP $f$ .	DBLP $c$ .	Email	PubMed
GEMPE	<b>0.38</b>	<b>0.27</b>	0.67	<u>0.29</u>
AROPE	0.02	0.09	0.50	<u>0.03</u>
deepWalk	0.35	<u>0.26</u>	<u>0.69</u>	<b>0.30</b>
LINE	0.00	<u>0.00</u>	<u>0.55</u>	0.00
SDNE	0.02	0.04	0.62	0.02
node2vec	0.09	0.21	<b>0.70</b>	0.14
struct2vec	0.00	0.00	0.22	0.00
VERSE	<u>0.36</u>	0.25	0.67	0.28

**Table 1: Comparison of Clustering Results (NMI) in 128D.**

# Experiments: Graph Drawing

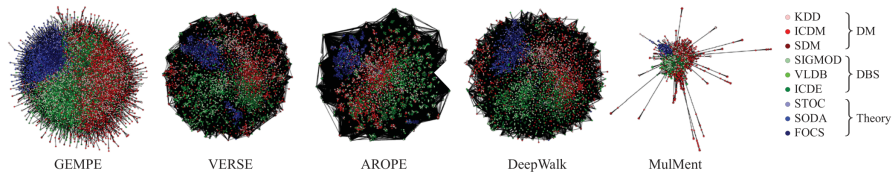


Figure 7: Visualization of DBLP Co-authorship Graph.

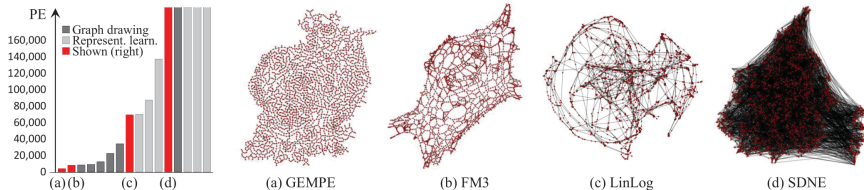


Figure 8: PE as a Quality Measure for Graph Drawing: Clear drawings of the Minnesota road network tend to have a low PE.

# References I



Ulrik Brandes and Christian Pich. “An Experimental Study on Distance-Based Graph Drawing”. In: *Graph Drawing*. Ed. by Ioannis G. Tollis and Maurizio Patrignani. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 218–229. ISBN: 978-3-642-00219-9.



Emden R. Gansner, Yehuda Koren, and Stephen North. “Graph Drawing by Stress Majorization”. In: *Proceedings of the 12th International Conference on Graph Drawing*. GD'04. New York, NY: Springer-Verlag, 2004, pp. 239–250. ISBN: 3540245286. DOI: 10.1007/978-3-540-31843-9\_25. URL: [https://doi.org/10.1007/978-3-540-31843-9\\_25](https://doi.org/10.1007/978-3-540-31843-9_25).