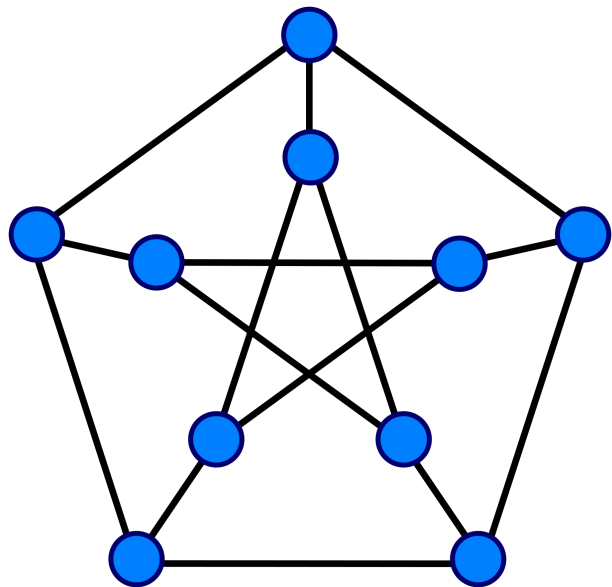# Bayesian approach in applications

Tomasz Kajdanowicz

# Embedding in graphs

# Graph encoding



[0, 1.2, 0.8, 4, ...]

[1, 0.1, 1.2, 3, ...]

[0, 1.2, 0.5, 2, ...]

[1, 0.3, 2.2, 4, ...]

# Node embedding
Applications

- node classification
- user behaviour prediction
- advertisement personalization
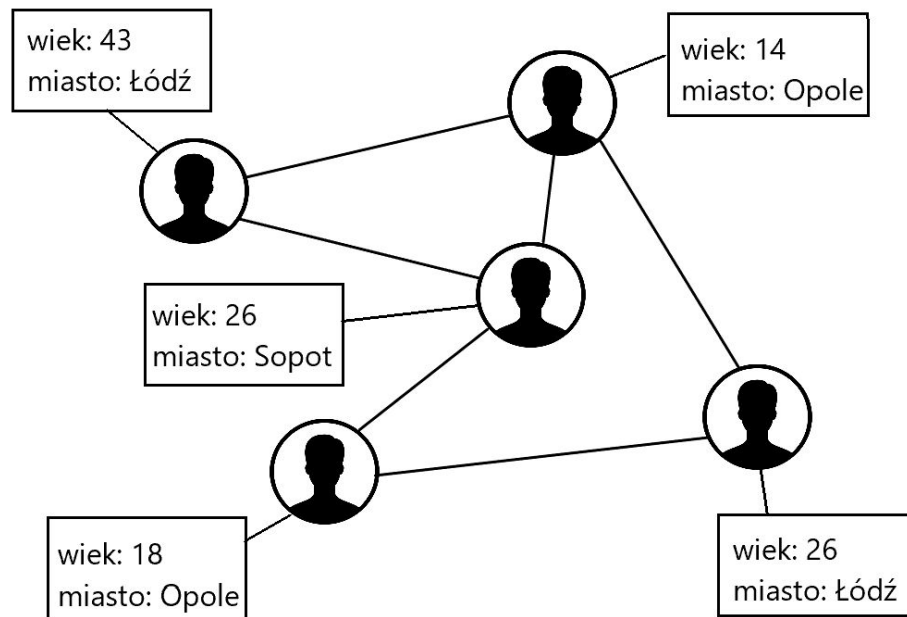- friends recommendation

# Methods for graph encoding
PREPROCESSING

"hard-coded" features

- features included in the dataset
- node labels

network centrality measures:

- node degree
- clustering coefficient
- number of shortest paths passing through a given node

# Methods for graph encoding

REPRESENTATION LEARNING - EMBEDDING

Learning of a mapping from nodes/edges/subgraphs/graphs to a low-dimensional vector space.

# d << |V|



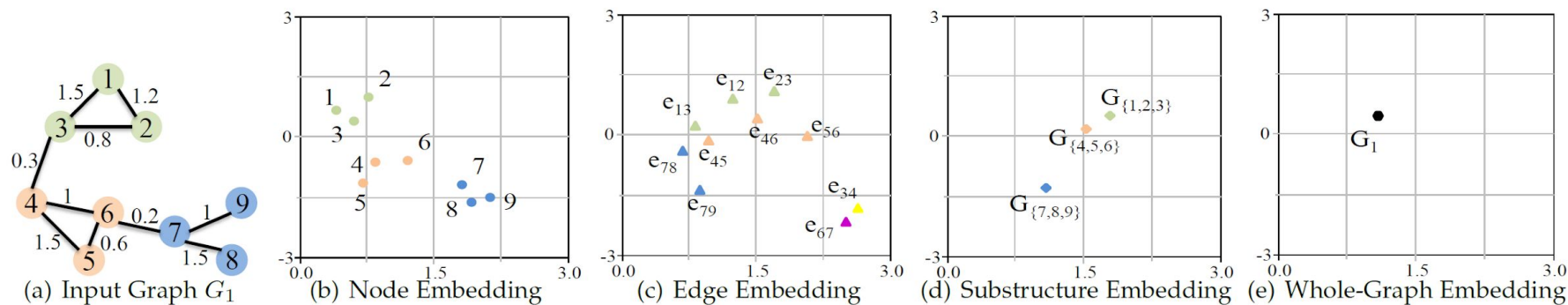(a) Input Graph $G_1$  (b) Node Embedding  (c) Edge Embedding  (d) Substructure Embedding  (e) Whole-Graph Embedding

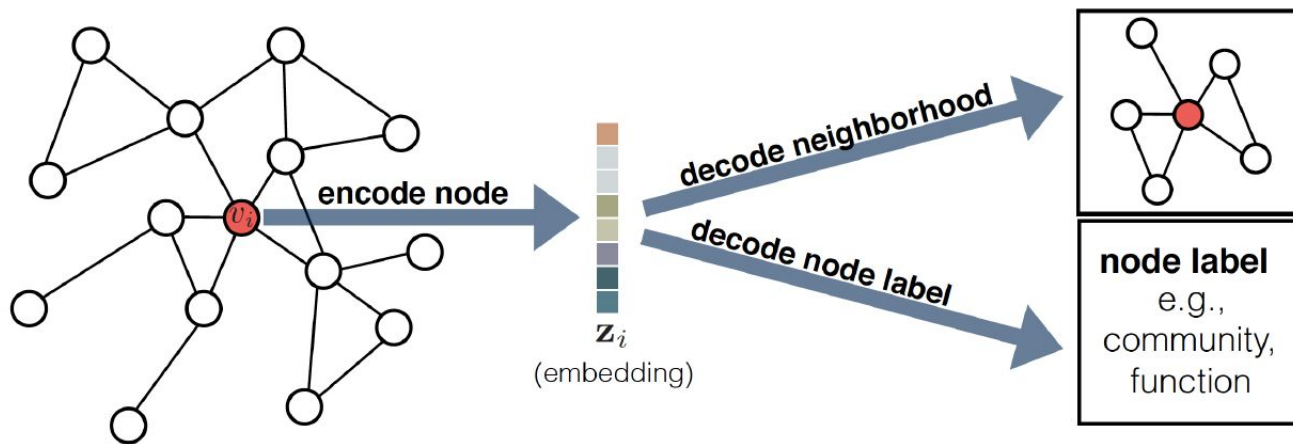Fig. 1. A toy example of embedding a graph into 2D space with different granularities. $G_{\{1,2,3\}}$ denotes the substructure containing node $v_1$, $v_2$, $v_3$.

Source: H. Cai et al., A Comprehensive Survey of Graph Embedding: Problems, Techniques and Applications
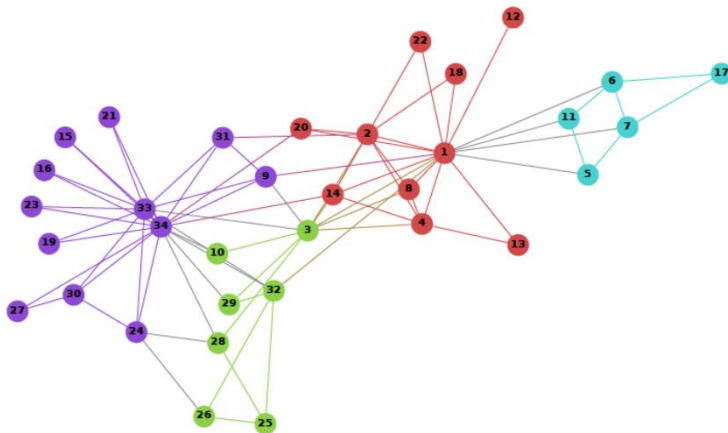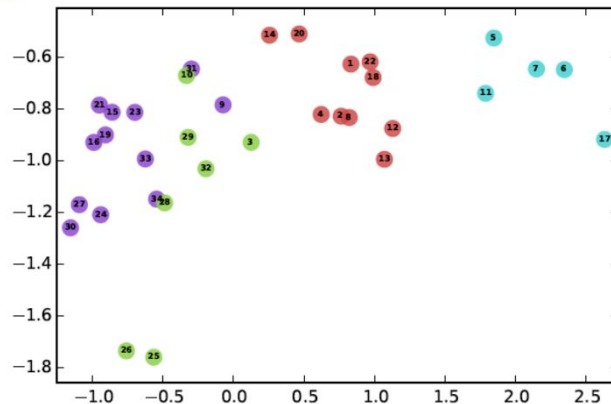
# Node embedding

# Node embedding

- "near" nodes should have a similar vector representation
- node proximity measure:
  - usually 1st or 2nd order node proximity
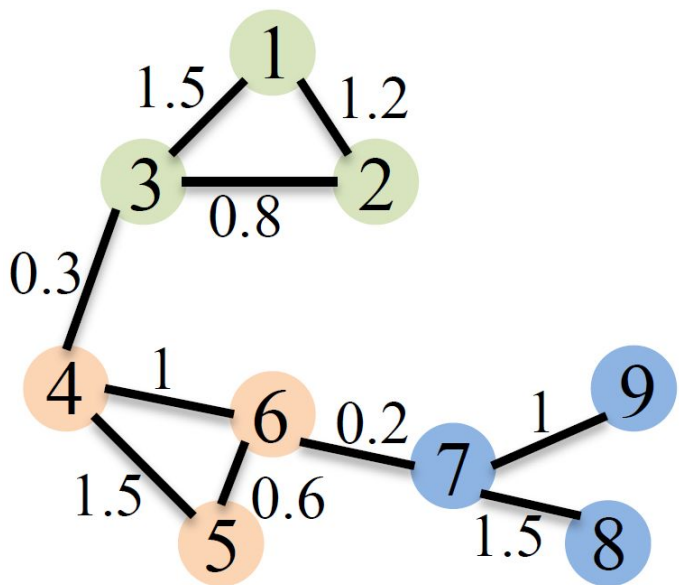


Zachary Karate Club social network

2D visualization of node embeddings

# 1st order node proximity

1st order neighborhood $s_{ij}^{(1)}$ of nodes $v_i$ and $v_j$ is the weight of the edge $e_{ij}$ between those.



$s_i^{(1)}=[s_{i1}^{(1)}, s_{i2}^{(1)},..., s_{i|V|}^{(1)}]$

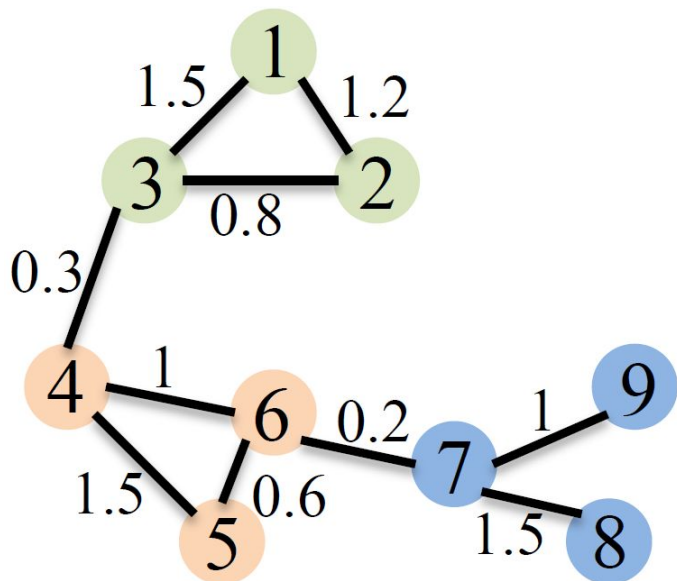$s_1^{(1)}=[s_{11}^{(1)}, s_{12}^{(1)},..., s_{19}^{(1)}]$

$s_1^{(1)}=[0, 1.2, 1.5, 0, 0, 0, 0, 0, 0]$

# 2nd order node proximity

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2} \sqrt{\sum\limits_{i=1}^{n} B_i^2}},$$

2nd order neighborhood $s_{ij}^{(2)}$ of nodes $v_i$ a $v_j$ is the similarity of their 1st order node neighborhoods: $s_i^{(1)}$ for node $v_i$ and $s_j^{(1)}$ for node $v_j$.



$s_1^{(1)}$=[0, 1.2, 1.5, 0, 0, 0, 0, 0, 0]

$s_2^{(1)}$=[1.2, 0, 0.8, 0, 0, 0, 0, 0, 0]

$s_{12}^{(2)}$ =cosine($s_1^{(1)}$, $s_2^{(1)}$)=0.43

$s_{15}^{(2)}$ =cosine($s_1^{(1)}$, $s_5^{(1)}$)=0

0 => no common neighbours

# Edge embedding

| Operator | Symbol | Definition |
|---|---|---|
| Average | $\boxplus$ | $[f(u) \boxplus f(v)]_i = \frac{f_i(u)+f_i(v)}{2}$ |
| Hadamard | $\boxdot$ | $[f(u) \boxdot f(v)]_i = f_i(u) * f_i(v)$ |
| Weighted-L1 | $\|\cdot\|_{\bar{1}}$ | $\|f(u) \cdot f(v)\|_{\bar{1}i} = \|f_i(u) - f_i(v)\|$ |
| Weighted-L2 | $\|\cdot\|_{\bar{2}}$ | $\|f(u) \cdot f(v)\|_{\bar{2}i} = \|f_i(u) - f_i(v)\|^2$ |

u, v - nodes
f(u), f(v) - embedding vectors for nodes u i v

# Edge embedding
LINK PREDICTION

### Prediction of missing edges

- finding such edges in datasets

### Prediction of "probable" edges

- recommender systems: friends, movies
- recommendation of potential scientific research topics

# Edge embedding
KNOWLEDGE GRAPHS

- knowledge database
- edge: <h, r, t>, <head entity, relation, tail entity>
- directed graph
- finding missing entity/relation based on the other two



Fig. 3. A toy example of knowledge graph.

<Alice, isFriendOf, Bob>
<Bob, isSupervisorOf, Chris>

# Subgraph / whole graph embedding

Whole graph embedding

comparing graph structures

Subgraph embedding

comparing communities in graphs

GRAPH DOMAIN

Set of graphs $S$

$$DM = \begin{pmatrix} 0 & d_{1,2} & d_{1,3} & \dots & d_{1,n} \\ d_{2,1} & 0 & d_{2,3} & \dots & d_{2,n} \\ d_{3,1} & d_{3,2} & 0 & \dots & d_{3,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{n,1} & d_{n,2} & d_{n,3} & \dots & 0 \end{pmatrix}$$

VECTOR DOMAIN

n-dimensional Vector Space

# Multi-Modal Bayesian Embeddings for Learning Social Knowledge Graphs

- learn **latent topics** that generate **word embeddings** and **network embeddings** simultaneously
- representation of social network users and knowledge concepts in a shared latent topic space



Male-Female | Verb tense | Country-Capital

Yang, Z., Tang, J., & Cohen, W. (n.d.). *Multi-Modal Bayesian Embeddings for Learning Social Knowledge Graphs.*

# Input

- social network $G^r = (V^r, E^r)$
  - $V^r$ set of users
  - $E^r$ set of edges between the users
- knowledge base $G^k = (V^k, C)$
  - $V^k$ set of knowledge concepts
  - $C$ text associated with or facts between the concepts
- text posted by users of the social network $D$
  - Given a user $u \in V^r$, $d_u \in D$ denotes a document of all text posted by $u$ (each user $u$ has only one document)

# Output

- social knowledge graph $G = (V^r, V^k, P)$
- given a user $u \in V^r$, $P_u$ is a ranked list of top-k knowledge concepts in $V^k$, where order indicates the relatedness to user $u$
- E.g. academic social network
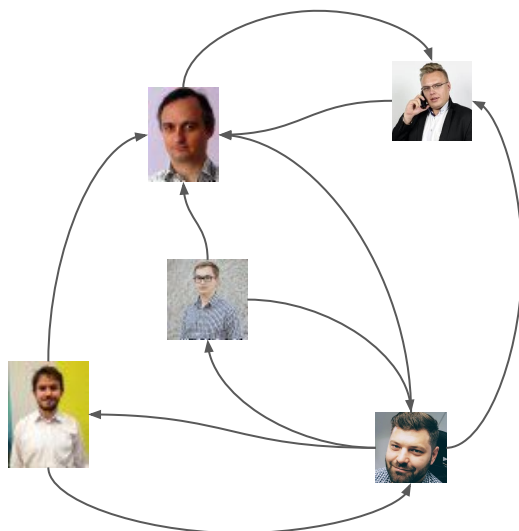  - top-k research interests of each researcher



| 1 | Machine Learning |
| 2 | Embedding |
| 3 | Bayesian inference |
| 4 | Social Networks |
| 5 | Social Media |

# Two modalities

social network of users

knowledge concepts

# More about input

- pretrained knowledge concept embeddings
  - encode information from the knowledge base $G^k$
  - *e.g. skip-gram model [Mikolov et al., 2013]*
- pretrained user embeddings as input
  - social network $G^r$
  - *e.g. DeepWalk [Perozzi et al., 2014]*

# The model

D - documents

Z - topics

M - concepts in a document

$f^r$ - embeddings of users drawn from Normal($\mu$,1/$\lambda$) -> from NormalGamma($\tau^r$)

$f^k_{um}$ - embedding of m-concept in document $d_u$



$\theta_u$ - multinomial topic distribution of document $d_u$ (or user $u$)

$z_{um}$ - topic of the m-th knowledge concept in document $d_u$

$y_u$ - topics of user $u$

# Generative process



1. For each topic $t$, and for each dimension

   (a) Draw $\mu_t^r, \lambda_t^r$ from $\text{NormalGamma}(\tau^r)$
   (b) Draw $\mu_t^k, \lambda_t^k$ from $\text{NormalGamma}(\tau^k)$

2. For each user $u$

   (a) Draw a multinomial distribution $\theta$ from $\text{Dir}(\alpha)$
   (b) For each knowledge concept $w$ in $d_u$
      i. Draw a topic $z$ from $\text{Multi}(\theta)$
      ii. For each dimension of the embedding of $w$, draw $f^k$ from $\mathcal{N}(\mu_z^k, \lambda_z^k)$
   (c) Draw a topic $y$ uniformly from all $z$'s in $d_u$
   (d) For each dimension of the embedding of user $u$, draw $f^r$ from $\mathcal{N}(\mu_y^r, \lambda_y^r)$

# Inference

- collapsed Gibbs sampling [Griffiths, 2002]
- extension of Gaussian LDA [Das et al., 2015] that updates the embeddings during inference

# Experiments

Data for embedding calculations:

- Miner - co-authorships between researchers
  - documents of authors
-  - knowledge base of articles
  - (Wikipedia corpus to learn the knowledge concept embeddings)

Evaluation:

Ranking of concepts from:

- homepage
- linkedin
- crowdsourcing



| 1 | Machine Learning |
|---|---|
| 2 | Embedding |
| 3 | Bayesian inference |
| 4 | Social Networks |
| 5 | Social Media |

| Topic #1 | Topic #2 | Topic #3 |
|---|---|---|
| GenVector | | |
| query expansion | image processing | hepatocellular carcinoma |
| concept mining | face recognition | gastric cancer |
| language modeling | feature extraction | acute lymphoblastic leukemia |
| information extraction | computer vision | renal cell carcinoma |
| knowledge extraction | image segmentation | glioblastoma multiforme |
| entity linking | image analysis | acute myeloid leukemia |
| language models | feature detection | peripheral blood |
| named entity recognition | digital image processing | malignant melanoma |
| document clustering | machine learning algorithms | hepatitis c virus |
| latent semantic indexing | machine vision | squamous cell carcinoma |
| Thorsten Joachims | Anil K. Jain | Keizo Sugimachi |
| Jian Pei | Thomas S. Huang | Setsuo Hirohashi |
| Christopher D. Manning | Peter N. Belhumeur | Masatoshi Makuuchi |
| Raymond J. Mooney | Azriel Rosenfeld | Morito Monden |
| Charu C. Aggarwal | Josef Kittler | Yoshio Yamaoka |
| William W. Cohen | Shuicheng Yan | Kunio Okuda |
| Eugene Charniak | David Zhang | Yasuni Nakanuma |
| Kamal Nigam | Xiaoou Tang | Kendo Kiyosawa |
| Susan T. Dumais | Roberto Cipolla | Masazumi Tsuneyoshi |
| T. K. Landauer | David A. Forsyth | Satoru Todo |

# Embedding Words as Distributions with a Bayesian Skip-gram Model

Idea:

- replace point word embedding with distribution

Havrylov, Serhii, and Ivan Titov. "Embedding Words as Distributions with a Bayesian Skip-gram Model." *Proceedings of the 27th International Conference on Computational Linguistics*. 2018.
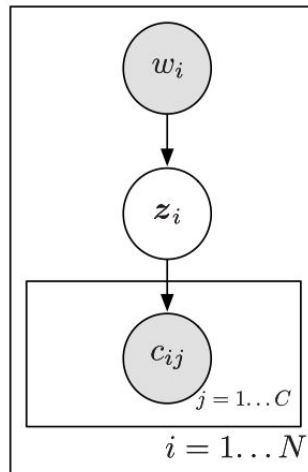
# Models

skip-gram



- For each data point $i = 1...N$
  - ⋆ For each context word $j = 1...C$
    - · Draw a context word $c_{ij} \sim p(c|w_i)$

Bayesian skip-gram



- For each data point $i = 1...N$
  - ⋆ Draw a latent vector $\boldsymbol{z}_i \sim p(\boldsymbol{z}|w_i)$
  - ⋆ For each context word $j = 1\ldots C$
    - · Draw a context word $c_{ij} \sim p(c|\boldsymbol{z}_i, w_i)$

# Examples

| word 1 | word 2 | KL | cosine sim. |
|---|---|---|---|
| dog | cat | 15.47 | 0.71 |
| dog | pet | 18.52 | 0.70 |
| dog | hound | 21.20 | 0.64 |
| dog | animal | 27.69 | 0.52 |
| cappuccino | espresso | 12.59 | 0.76 |
| cappuccino | latte | 13.39 | 0.7 |
| cappuccino | coffee | 22.54 | 0.69 |
| cappuccino | drink | 30.81 | 0.54 |
| microsoft | windows | 24.41 | 0.65 |
| microsoft | google | 24.44 | 0.60 |
| microsoft | corporation | 39.40 | 0.29 |
| microsoft | company | 46.05 | 0.19 |