

Cross-Graph Learning of Multi-Relational Associations

Hanxiao Liu, Yiming Yang
Carnegie Mellon University
`{hanxiaol, yiming}@cs.cmu.edu`

June 22, 2016

Outline

Task Description

New Contributions

Framework

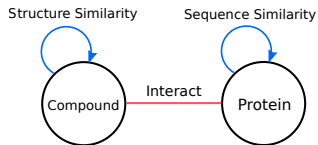
Scalable Inference

Empirical Evaluation

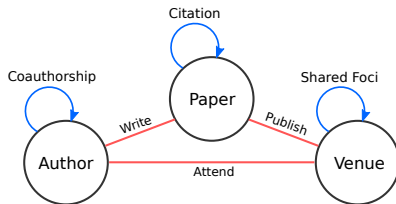
Summary

Task Description

Goal: Predict **associations** among heterogeneous **graphs**.



(a) Drug-Target Interaction



(b) Citation Network

“John publish a reinforcement learning paper at ICML.”
(John, RL_Paper, ICML)

Outline

Task Description

New Contributions

Framework

Scalable Inference

Empirical Evaluation

Summary

New Contributions

- ▶ A unified framework to integrating heterogeneous information in multiple graphs.
- ▶ Transductive learning to leverage both labeled data (sparse) and unlabeled data (massive).
- ▶ A convex approximation for the scalable inference over the combinatorial number of possible tuples.

Outline

Task Description

New Contributions

Framework

Scalable Inference

Empirical Evaluation

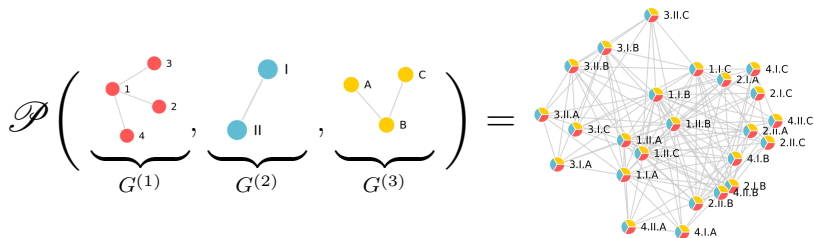
Summary

Notation

- ▶ $G^{(1)}, G^{(2)}, \dots, G^{(J)}$ are individual graphs;
- ▶ n_j is the #nodes in $G^{(j)}$;
- ▶ (i_1, i_2, \dots, i_J) is a tuple (multi-relation);
- ▶ f_{i_1, i_2, \dots, i_J} is the predicted score for the tuple;
- ▶ f is a tensor in $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_J}$.

Framework

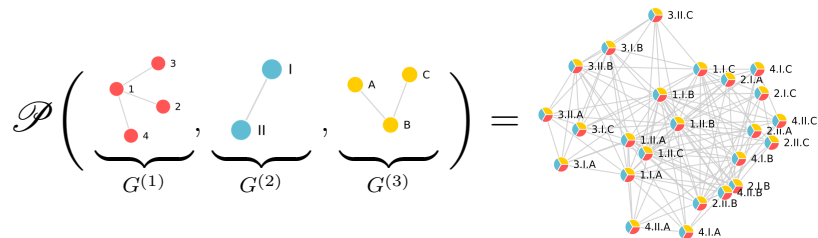
Product Graph (\mathcal{P}) induced from $G^{(1)}, \dots, G^{(J)}$.



Tensor product: $\mathcal{P}(G^{(1)}, G^{(2)}, G^{(3)}) = G^{(1)} \otimes G^{(2)} \otimes G^{(3)}$

Framework

Product Graph (\mathcal{P}) induced from $G^{(1)}, \dots, G^{(J)}$.



Tensor product: $\mathcal{P}(G^{(1)}, G^{(2)}, G^{(3)}) = G^{(1)} \otimes G^{(2)} \otimes G^{(3)}$

Why product graph?

- ▶ Mapping heterogeneous graphs onto a unified graph for label propagation (transductive learning).

Framework

Assuming

$$\text{vec}(f) \sim \mathcal{N}(0, \mathcal{P}) \quad (1)$$

which implies:

$$-\log p(f|\mathcal{P}) \propto \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) := \|f\|_{\mathcal{P}}^2 \quad (2)$$

Optimization problem

$$\min_f \ell_{\mathcal{O}}(f) + \frac{\gamma}{2} \|f\|_{\mathcal{P}}^2 \quad (3)$$

Framework

Assuming

$$\text{vec}(f) \sim \mathcal{N}(0, \mathcal{P}) \quad (1)$$

which implies:

$$-\log p(f|\mathcal{P}) \propto \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) := \|f\|_{\mathcal{P}}^2 \quad (2)$$

Optimization problem

$$\min_f \ell_{\mathcal{O}}(f) + \frac{\gamma}{2} \|f\|_{\mathcal{P}}^2 \quad (3)$$

Framework

Assuming

$$\text{vec}(f) \sim \mathcal{N}(0, \mathcal{P}) \quad (1)$$

which implies:

$$-\log p(f|\mathcal{P}) \propto \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) := \|f\|_{\mathcal{P}}^2 \quad (2)$$

Optimization problem

$$\min_f \ell_{\mathcal{O}}(f) + \frac{\gamma}{2} \|f\|_{\mathcal{P}}^2 \quad (3)$$

For computational tractability, we focus on the spectral graph product family of \mathcal{P} .

Spectral Graph Product (SGP)

The eigensystem of $\mathcal{P}_\kappa(G^{(1)}, \dots, G^{(J)})$ is parametrized by the eigensystems of individual graphs, i.e.,

$$\left\{ \kappa(\lambda_{i_1}, \dots, \lambda_{i_J}), \bigotimes_j v_{i_j} \right\}_{i_1, \dots, i_J} \quad (4)$$

λ_{i_j}/v_{i_j} is the i_j -th eigenvalue/eigenvector of the j -th graph.

Framework

Nice properties of SGP:

Subsuming basic operations

$$\kappa(x, y) = x \times y \implies \mathcal{P}_\kappa(G, H) = G \otimes H \quad \text{Tensor} \quad (5)$$

$$\kappa(x, y) = x + y \implies \mathcal{P}_\kappa(G, H) = G \oplus H \quad \text{Cartesian} \quad (6)$$

Supporting graph diffusions

$$\sigma_{Heat}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \frac{1}{2}\mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{e^\kappa} \quad (7)$$

$$\sigma_{von-Neumann}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{\frac{1}{1-\kappa}} \quad (8)$$

Order-insensitive: If κ is commutative, then SGP is commutative (up to graph isomorphism).

Framework

Nice properties of SGP:

Subsuming basic operations

$$\kappa(x, y) = x \times y \implies \mathcal{P}_\kappa(G, H) = G \otimes H \quad \text{Tensor} \quad (5)$$

$$\kappa(x, y) = x + y \implies \mathcal{P}_\kappa(G, H) = G \oplus H \quad \text{Cartesian} \quad (6)$$

Supporting graph diffusions

$$\sigma_{Heat}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \frac{1}{2}\mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{e^\kappa} \quad (7)$$

$$\sigma_{von-Neumann}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{\frac{1}{1-\kappa}} \quad (8)$$

Order-insensitive: If κ is commutative, then SGP is commutative (up to graph isomorphism).

Framework

Nice properties of SGP:

Subsuming basic operations

$$\kappa(x, y) = x \times y \implies \mathcal{P}_\kappa(G, H) = G \otimes H \quad \text{Tensor} \quad (5)$$

$$\kappa(x, y) = x + y \implies \mathcal{P}_\kappa(G, H) = G \oplus H \quad \text{Cartesian} \quad (6)$$

Supporting graph diffusions

$$\sigma_{Heat}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \frac{1}{2}\mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{e^\kappa} \quad (7)$$

$$\sigma_{von-Neumann}(\mathcal{P}_\kappa) = I + \mathcal{P}_\kappa + \mathcal{P}_\kappa^2 + \dots = \mathcal{P}_{\frac{1}{1-\kappa}} \quad (8)$$

Order-insensitive: If κ is commutative, then SGP is commutative (up to graph isomorphism).

Outline

Task Description

New Contributions

Framework

Scalable Inference

Empirical Evaluation

Summary

Scalable Inference

For general GP, the semi-norm is computed as

$$\|f\|_{\mathcal{P}}^2 = \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) \quad (9)$$

For SGP, \mathcal{P}_κ no longer has to be explicitly computed.

$$\|f\|_{\mathcal{P}_\kappa}^2 = \sum_{i_1, i_2, \dots, i_J}^{n_1, n_2, \dots, n_J} \frac{f(v_{i_1}, \dots, v_{i_J})^2}{\kappa(\lambda_{i_1}, \dots, \lambda_{i_J})} \quad (10)$$

- ▶ $f(v_{i_1}, v_{i_2}, \dots, v_{i_J}) = f \times_1 v_{i_1} \times_2 v_{i_2} \cdots \times_J v_{i_J}$
- ▶ However, even evaluating (10) is expensive.

Scalable Inference

For general GP, the semi-norm is computed as

$$\|f\|_{\mathcal{P}}^2 = \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) \quad (9)$$

For SGP, \mathcal{P}_κ no longer has to be explicitly computed.

$$\|f\|_{\mathcal{P}_\kappa}^2 = \sum_{i_1, i_2, \dots, i_J}^{n_1, n_2, \dots, n_J} \frac{f(v_{i_1}, \dots, v_{i_J})^2}{\kappa(\lambda_{i_1}, \dots, \lambda_{i_J})} \quad (10)$$

- ▶ $f(v_{i_1}, v_{i_2}, \dots, v_{i_J}) = f \times_1 v_{i_1} \times_2 v_{i_2} \cdots \times_J v_{i_J}$
- ▶ However, even evaluating (10) is expensive.

Scalable Inference

For general GP, the semi-norm is computed as

$$\|f\|_{\mathcal{P}}^2 = \text{vec}(f)^\top \mathcal{P}^{-1} \text{vec}(f) \quad (9)$$

For SGP, \mathcal{P}_κ no longer has to be explicitly computed.

$$\|f\|_{\mathcal{P}_\kappa}^2 = \sum_{i_1, i_2, \dots, i_J}^{n_1, n_2, \dots, n_J} \frac{f(v_{i_1}, \dots, v_{i_J})^2}{\kappa(\lambda_{i_1}, \dots, \lambda_{i_J})} \quad (10)$$

- ▶ $f(v_{i_1}, v_{i_2}, \dots, v_{i_J}) = f \times_1 v_{i_1} \times_2 v_{i_2} \cdots \times_J v_{i_J}$
- ▶ However, even evaluating (10) is expensive.

Scalable Inference

Using low-rank SGP

- ▶ f lies in the linear span of the eigenvectors of \mathcal{P} .
- ▶ Eigenvectors of high volatility can be pruned away.

Scalable Inference

Using low-rank SGP

- ▶ f lies in the linear span of the eigenvectors of \mathcal{P} .
- ▶ Eigenvectors of high volatility can be pruned away.

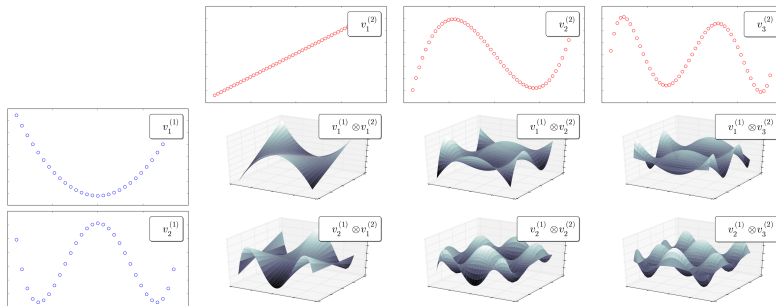


Figure : Eigenvectors of G (blue), H (red) and $\mathcal{P}(G, H)$.

Scalable Inference

Restrict f in the linear span of “smooth” bases of \mathcal{P} .

$$f(\alpha) = \sum_{i_1, i_2, \dots, i_J=1}^{d_1, d_2, \dots, d_J} \alpha_{i_1, i_2, \dots, i_J} \bigotimes_j v_{i_j} \quad (11)$$

where the *core tensor* $\alpha \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_J}$, $d_j \ll n_j$.

The semi-norm becomes

$$\|f(\alpha)\|_{\mathcal{P}_\kappa}^2 = \sum_{i_1, i_2, \dots, i_J=1}^{d_1, d_2, \dots, d_J} \frac{\alpha_{i_1, i_2, \dots, i_J}^2}{\kappa(\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_J})} \quad (12)$$

We then optimize w.r.t. α instead of f . Parameter size:
 $\prod_j n_j \rightarrow \prod_j d_j$.

Scalable Inference

Restrict f in the linear span of “smooth” bases of \mathcal{P} .

$$f(\alpha) = \sum_{i_1, i_2, \dots, i_J=1}^{d_1, d_2, \dots, d_J} \alpha_{i_1, i_2, \dots, i_J} \bigotimes_j v_{i_j} \quad (11)$$

where the *core tensor* $\alpha \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_J}$, $d_j \ll n_j$.

The semi-norm becomes

$$\|f(\alpha)\|_{\mathcal{P}_\kappa}^2 = \sum_{i_1, i_2, \dots, i_J=1}^{d_1, d_2, \dots, d_J} \frac{\alpha_{i_1, i_2, \dots, i_J}^2}{\kappa(\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_J})} \quad (12)$$

We then optimize w.r.t. α instead of f . Parameter size:
 $\prod_j n_j \rightarrow \prod_j d_j$.

Scalable Inference

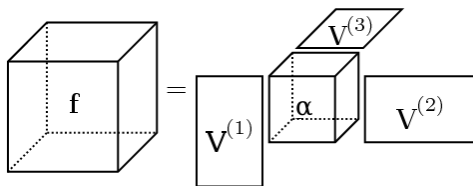


Figure : Tucker Decomposition, where α is the core tensor.

Scalable Inference

Revised optimization objective

$$\min_{\alpha \in \mathbb{R}^{d_1 \times d_2 \cdots \times d_J}} \ell_{\mathcal{O}}(f(\alpha)) + \frac{\gamma}{2} \|f(\alpha)\|_{\mathcal{P}_{\kappa}}^2 \quad (13)$$

Ranking loss function

$$\ell_{\mathcal{O}}(f) = \frac{\sum_{\substack{(i_1, \dots, i_J) \in \mathcal{O} \\ (i'_1, \dots, i'_J) \in \bar{\mathcal{O}}}} (f_{i_1 \dots i_J} - f_{i'_1 \dots i'_J})_+^2}{|\mathcal{O} \times \bar{\mathcal{O}}|} \quad (14)$$

$$\nabla_{\alpha} = \frac{\partial \ell_{\mathcal{O}}}{\partial f} \left(\frac{\partial f_{i_1, \dots, i_J}}{\partial \alpha} - \frac{\partial f_{i'_1, \dots, i'_J}}{\partial \alpha} \right) + \gamma \alpha \oslash \kappa \quad (15)$$

Tensor algebras are carried out on GPU.

Scalable Inference

Revised optimization objective

$$\min_{\alpha \in \mathbb{R}^{d_1 \times d_2 \cdots \times d_J}} \ell_{\mathcal{O}}(f(\alpha)) + \frac{\gamma}{2} \|f(\alpha)\|_{\mathcal{P}_{\kappa}}^2 \quad (13)$$

Ranking loss function

$$\ell_{\mathcal{O}}(f) = \frac{\sum_{\substack{(i_1, \dots, i_J) \in \mathcal{O} \\ (i'_1, \dots, i'_J) \in \bar{\mathcal{O}}}} (f_{i_1 \dots i_J} - f_{i'_1 \dots i'_J})_+^2}{|\mathcal{O} \times \bar{\mathcal{O}}|} \quad (14)$$

$$\nabla_{\alpha} = \frac{\partial \ell_{\mathcal{O}}}{\partial f} \left(\frac{\partial f_{i_1, \dots, i_J}}{\partial \alpha} - \frac{\partial f_{i'_1, \dots, i'_J}}{\partial \alpha} \right) + \gamma \alpha \oslash \kappa \quad (15)$$

Tensor algebras are carried out on GPU.

Scalable Inference

Revised optimization objective

$$\min_{\alpha \in \mathbb{R}^{d_1 \times d_2 \cdots \times d_J}} \ell_{\mathcal{O}}(f(\alpha)) + \frac{\gamma}{2} \|f(\alpha)\|_{\mathcal{P}_{\kappa}}^2 \quad (13)$$

Ranking loss function

$$\ell_{\mathcal{O}}(f) = \frac{\sum_{\substack{(i_1, \dots, i_J) \in \mathcal{O} \\ (i'_1, \dots, i'_J) \in \bar{\mathcal{O}}}} (f_{i_1 \dots i_J} - f_{i'_1 \dots i'_J})_+^2}{|\mathcal{O} \times \bar{\mathcal{O}}|} \quad (14)$$

$$\nabla_{\alpha} = \frac{\partial \ell_{\mathcal{O}}}{\partial f} \left(\frac{\partial f_{i_1, \dots, i_J}}{\partial \alpha} - \frac{\partial f_{i'_1, \dots, i'_J}}{\partial \alpha} \right) + \gamma \alpha \oslash \kappa \quad (15)$$

Tensor algebras are carried out on GPU.

Outline

Task Description

New Contributions

Framework

Scalable Inference

Empirical Evaluation

Summary

Empirical Evaluation

Datasets

Enzyme 445 compounds, 664 proteins.

DBLP 34*K* authors, 11*K* papers, 22 venues.

Representative Baselines

TF/GRTF Tensor Factorization/Graph-Regularized TF

NN One-class Nearest Neighbor

RSVM Ranking SVMs

LTKM Low-Rank Tensor Kernel Machines

Empirical Evaluation

Datasets

Enzyme 445 compounds, 664 proteins.

DBLP 34*K* authors, 11*K* papers, 22 venues.

Representative Baselines

TF/GRTF Tensor Factorization/Graph-Regularized TF

NN One-class Nearest Neighbor

RSVM Ranking SVMs

LTKM Low-Rank Tensor Kernel Machines

Empirical Evaluation

Our method: “TOP” (blue).

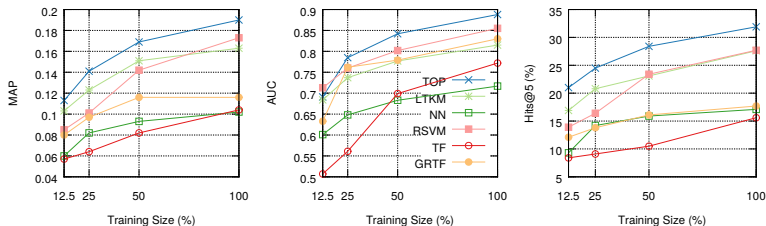
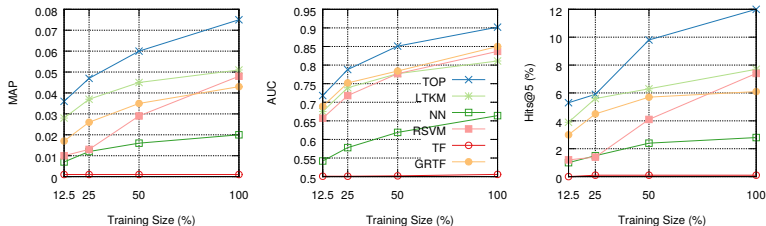


Figure : Performance on Enzyme (above) and DBLP (below).



Outline

Task Description

New Contributions

Framework

Scalable Inference

Empirical Evaluation

Summary

Summary

Contribution

- ▶ A unified framework to integrating heterogeneous information in multiple graphs.
- ▶ Transductive learning to leverage both labeled data (sparse) and unlabeled data (massive).
- ▶ A convex approximation for the scalable inference over the combinatorial number of possible tuples.

Future/On-going Work

- ▶ Learning structured associations.
- ▶ Larger problems: Microsoft Academic Graph (37 GB).

Summary

Contribution

- ▶ A unified framework to integrating heterogeneous information in multiple graphs.
- ▶ Transductive learning to leverage both labeled data (sparse) and unlabeled data (massive).
- ▶ A convex approximation for the scalable inference over the combinatorial number of possible tuples.

Future/On-going Work

- ▶ Learning structured associations.
- ▶ Larger problems: Microsoft Academic Graph (37 GB).

Thank You