

Factorized Graph Representations for semi-supervised learning from sparse data

Krishna Kumar, Paul Langton, Wolfgang Gatterbauer

SIGMOD 2020, Thursday, June 18, 2020, R16: 3:00 – 4:30 pm PT

Slides: <https://github.com/northeastern-datalab/factorized-graphs/>

DOI: <https://doi.org/10.1145/3318464.3380577>

Data Lab: <https://db.khoury.northeastern.edu>

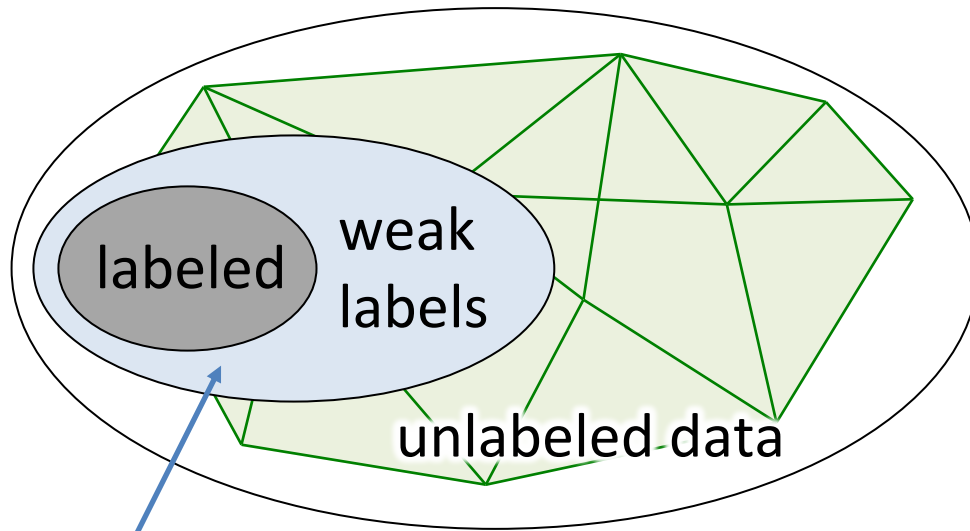


This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 4.0 International License.
See <https://creativecommons.org/licenses/by-nc-sa/4.0/> for details

Learning from few labels with algebraic amplification

Semi-supervised learning

exploit relationships on label distribution (e.g. smoothness in networks)

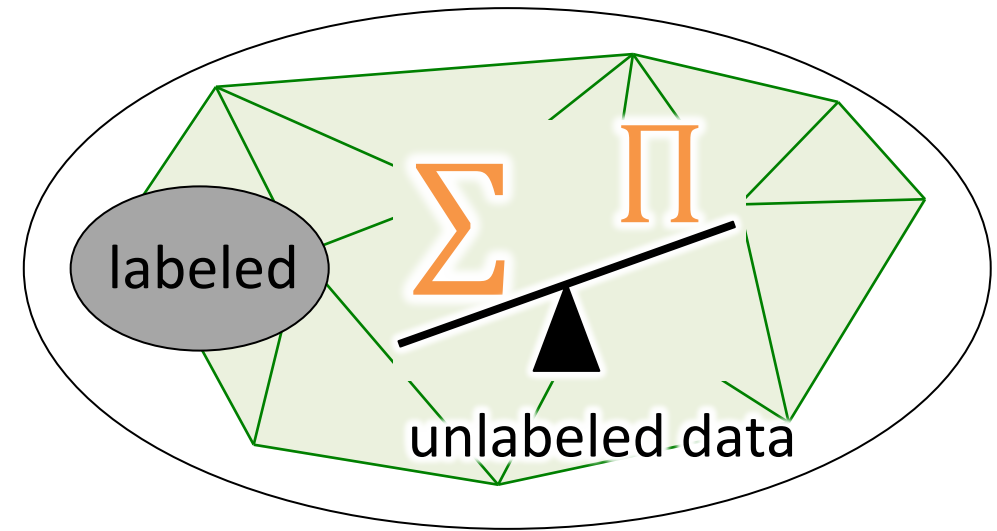


Weak (or distant) supervision

add noisier labels (e.g. heuristics, or external knowledge base)

Algebraic amplification

leverage algebraic properties of the algorithm to amplify signal in sparse data



Algebraic cheating

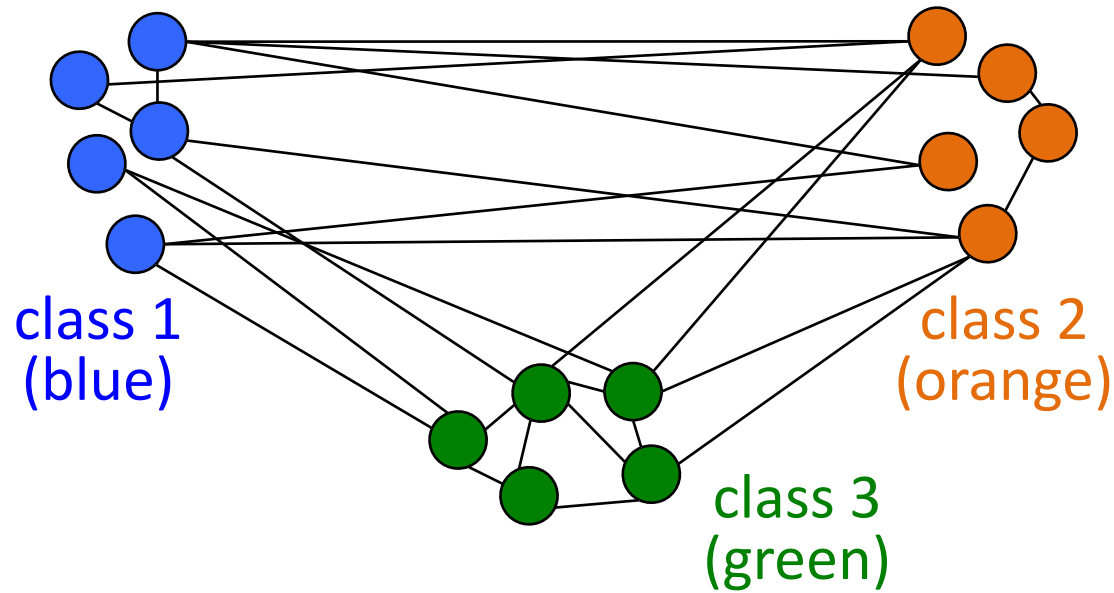
this requires "nice" algebraic properties; we may have to modify the algorithms ☺

Our focus today: Node classification in undirected graphs

Preference among node classes

\Rightarrow

Compatibilities between classes



orange prefers blue (and v.v.)

green prefers green

$\mathbf{H} =$

0.2	0.6	0.2
0.6	0.2	0.2
0.2	0.2	0.6

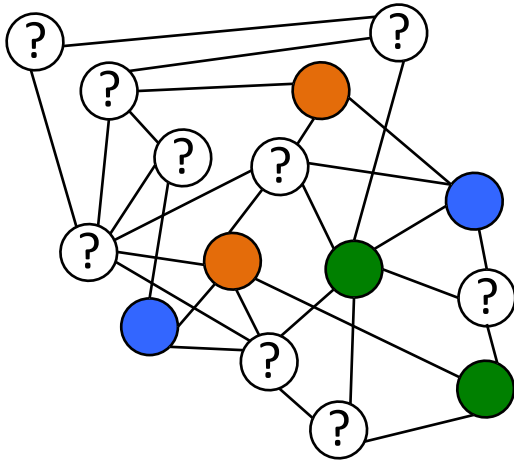
$\Sigma = 1$

Our focus today: Node classification in graphs

**Preference among node classes
most of which are unlabeled**

⇒

**Compatibilities between classes
not known to us ☹**



H=

	0.2	0.6	0.2
	0.6	0.2	0.2
	0.2	0.2	0.6

$\Sigma=1$

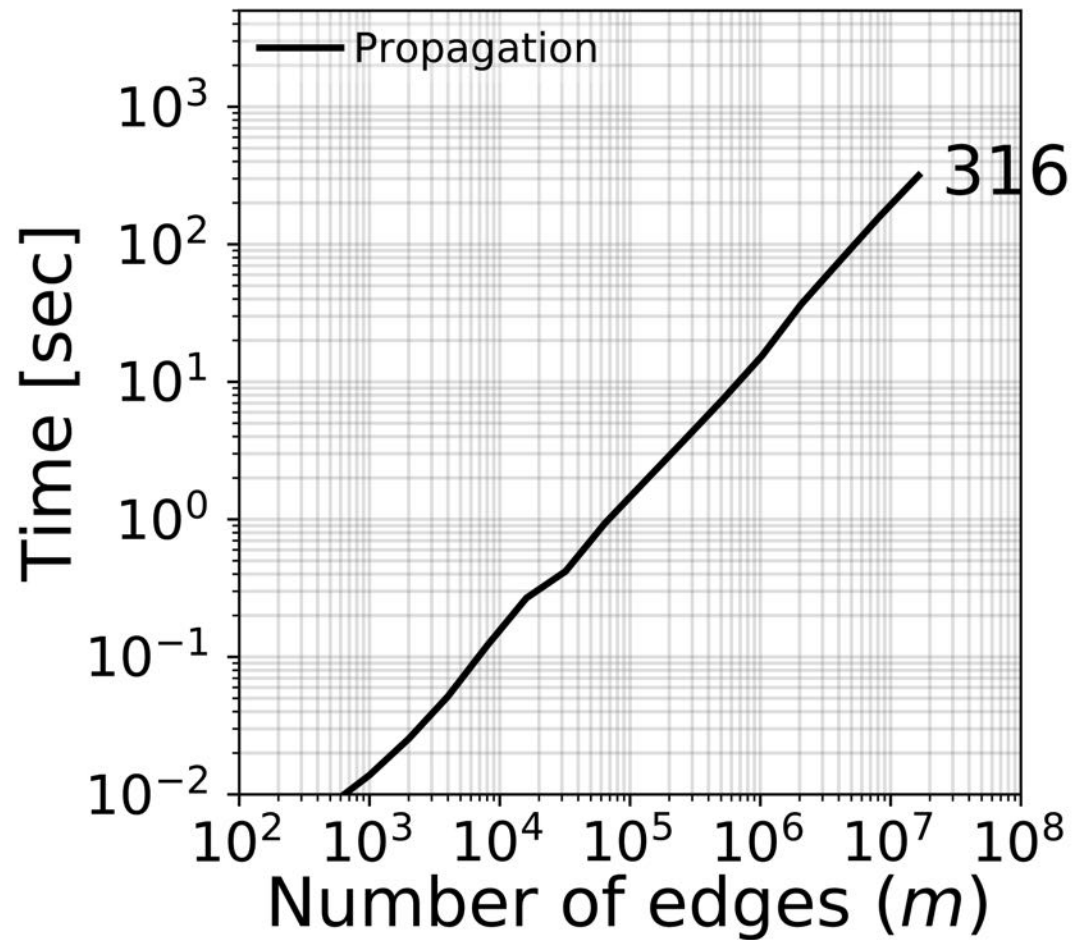
linearized belief propagation,
semi-supervised learning

Goal: Classify the remaining nodes **Estimate** & propagate those compatibilities

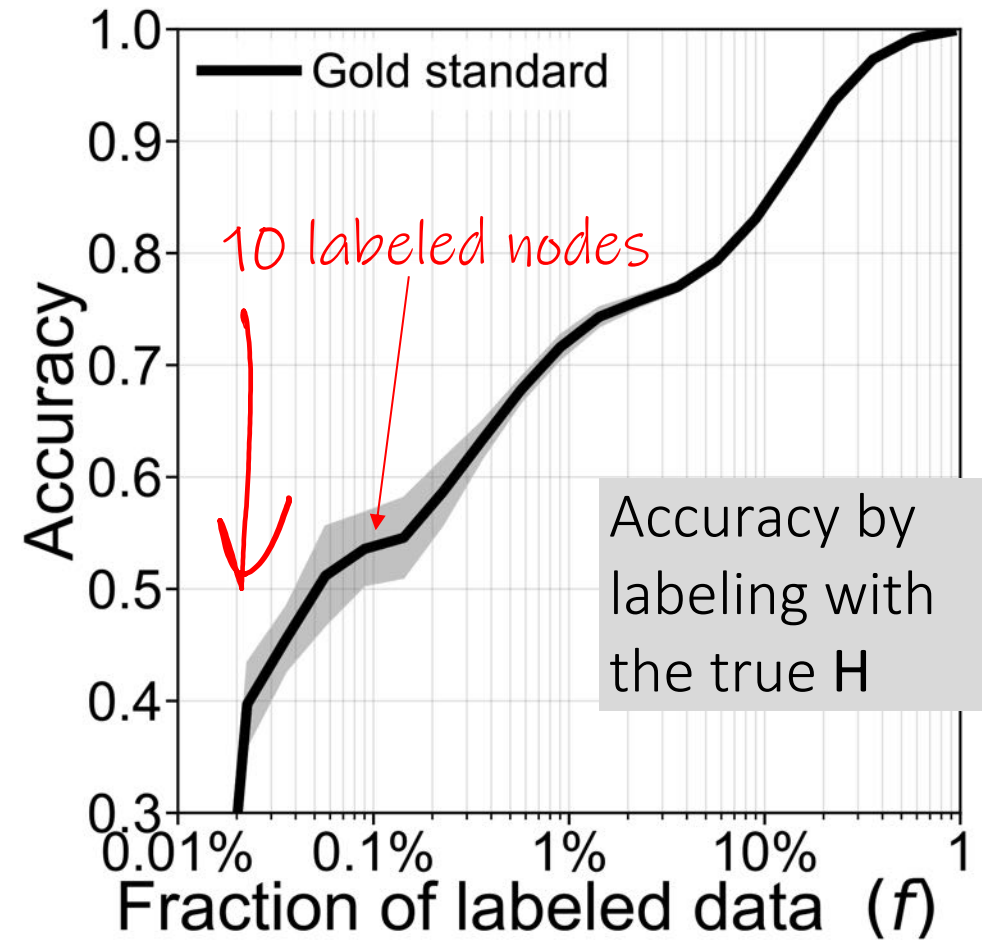
State-of-the-art: Heuristics / domain experts
We will estimate (learn) from sparse data

How well does it work?

Time and Accuracy for label propagation *if we know H*



Label propagation linear in # edges

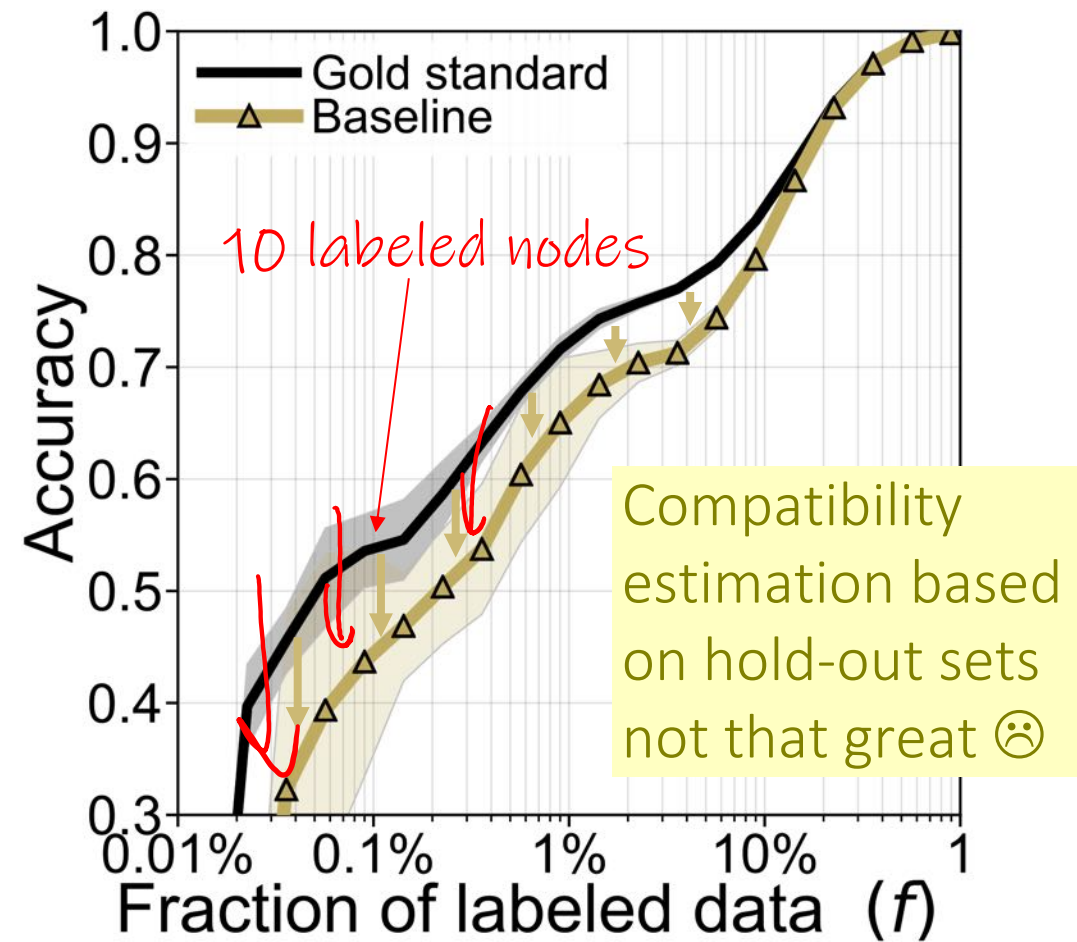
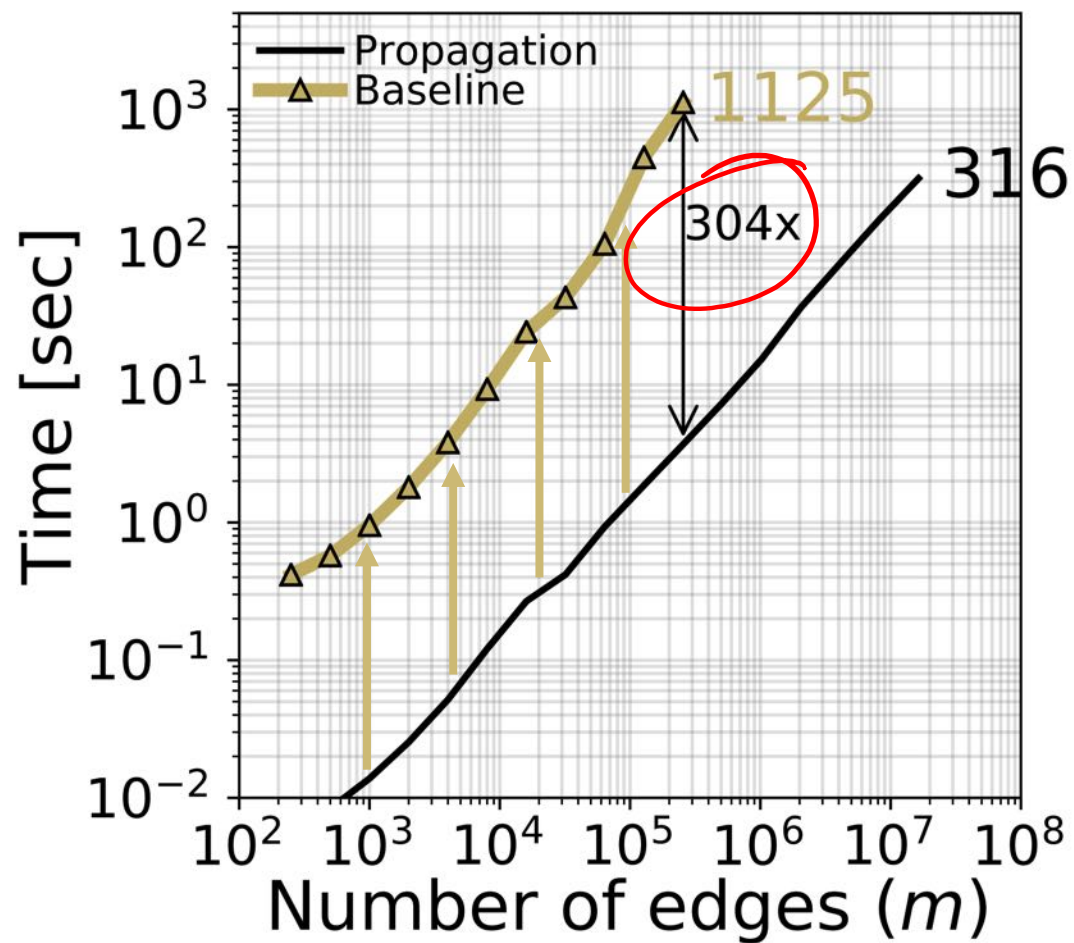


← Fewer labels

Details: 10k nodes, degree $d=25$, $H =$

0.2	0.6	0.2
0.6	0.2	0.2
0.2	0.2	0.6

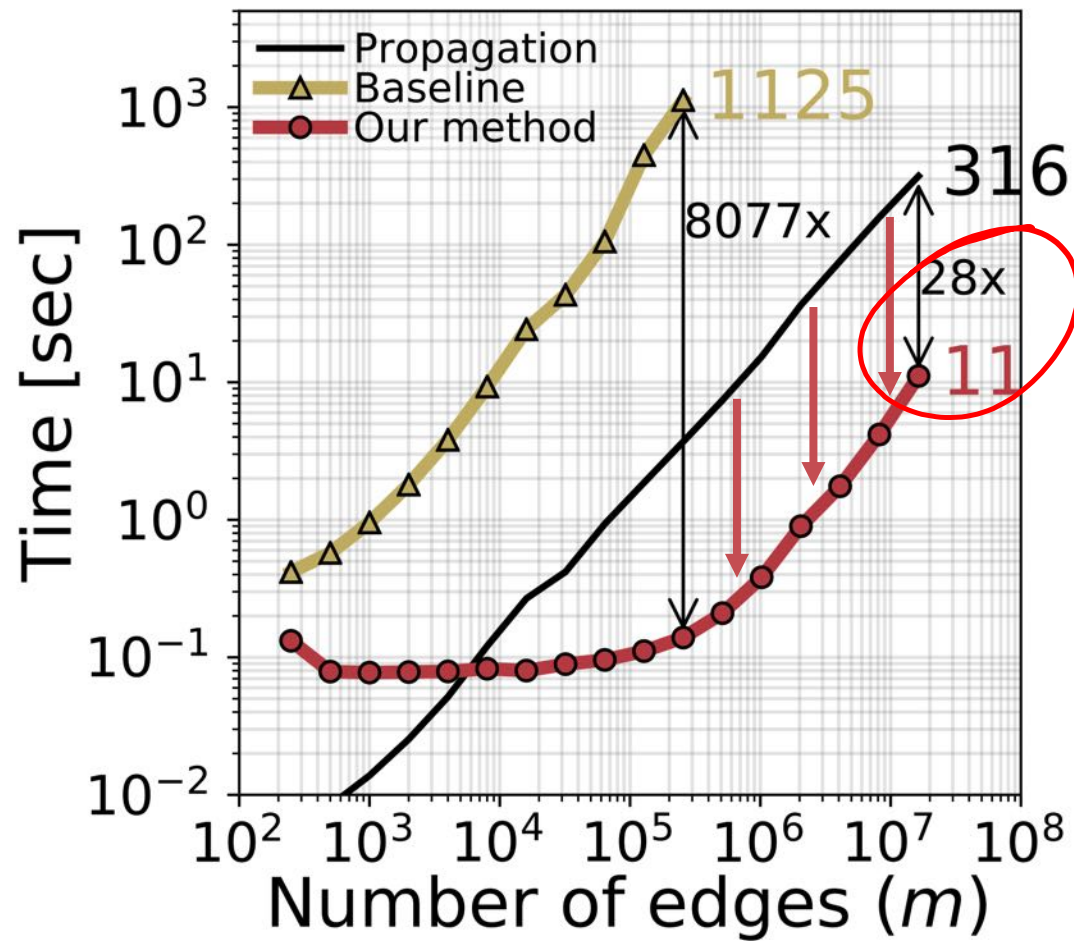
Time and Accuracy if we need to first estimate H ☹️



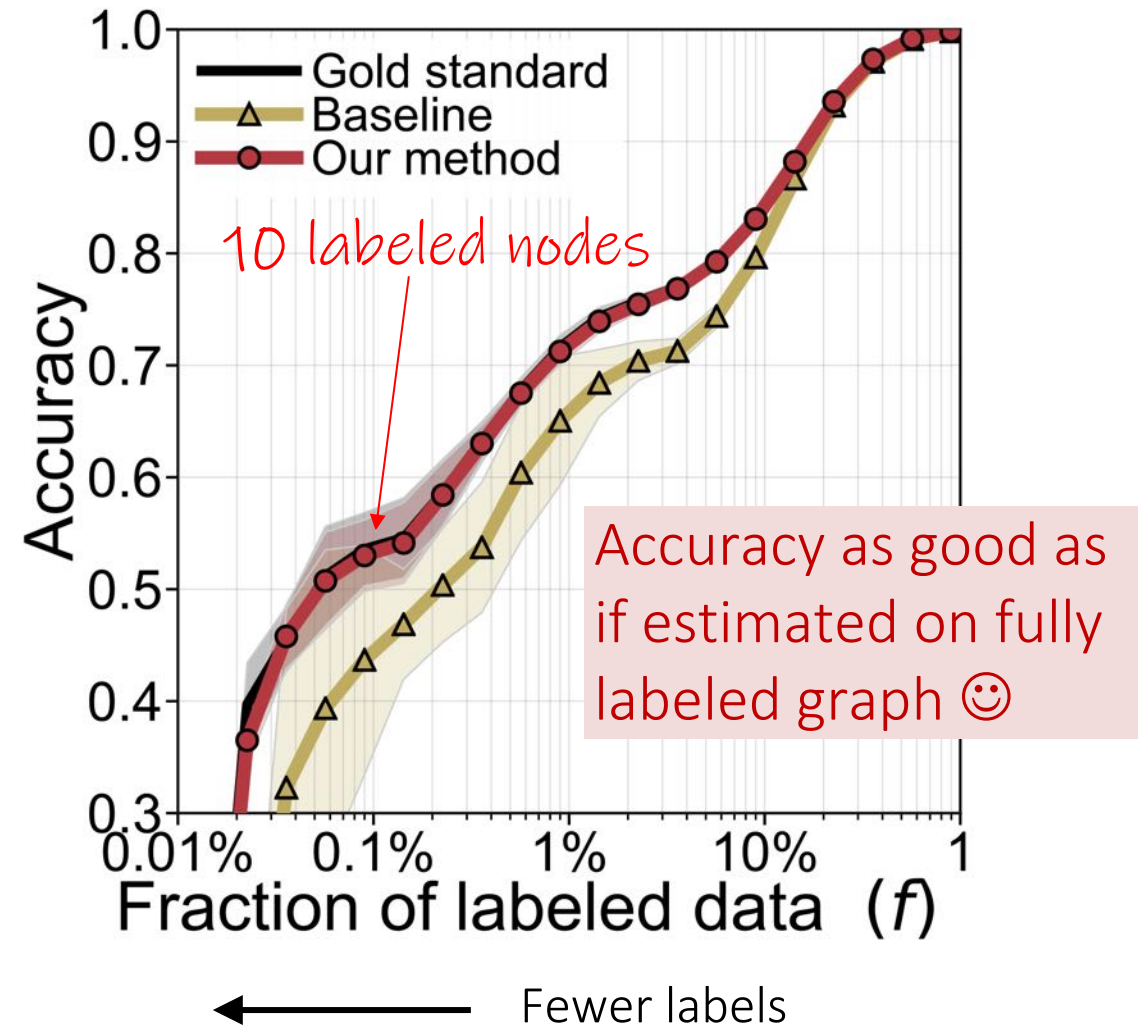
Estimation uses inference as subroutine (thus slower) ☹️

← Fewer labels

Time and Accuracy with our method 😊



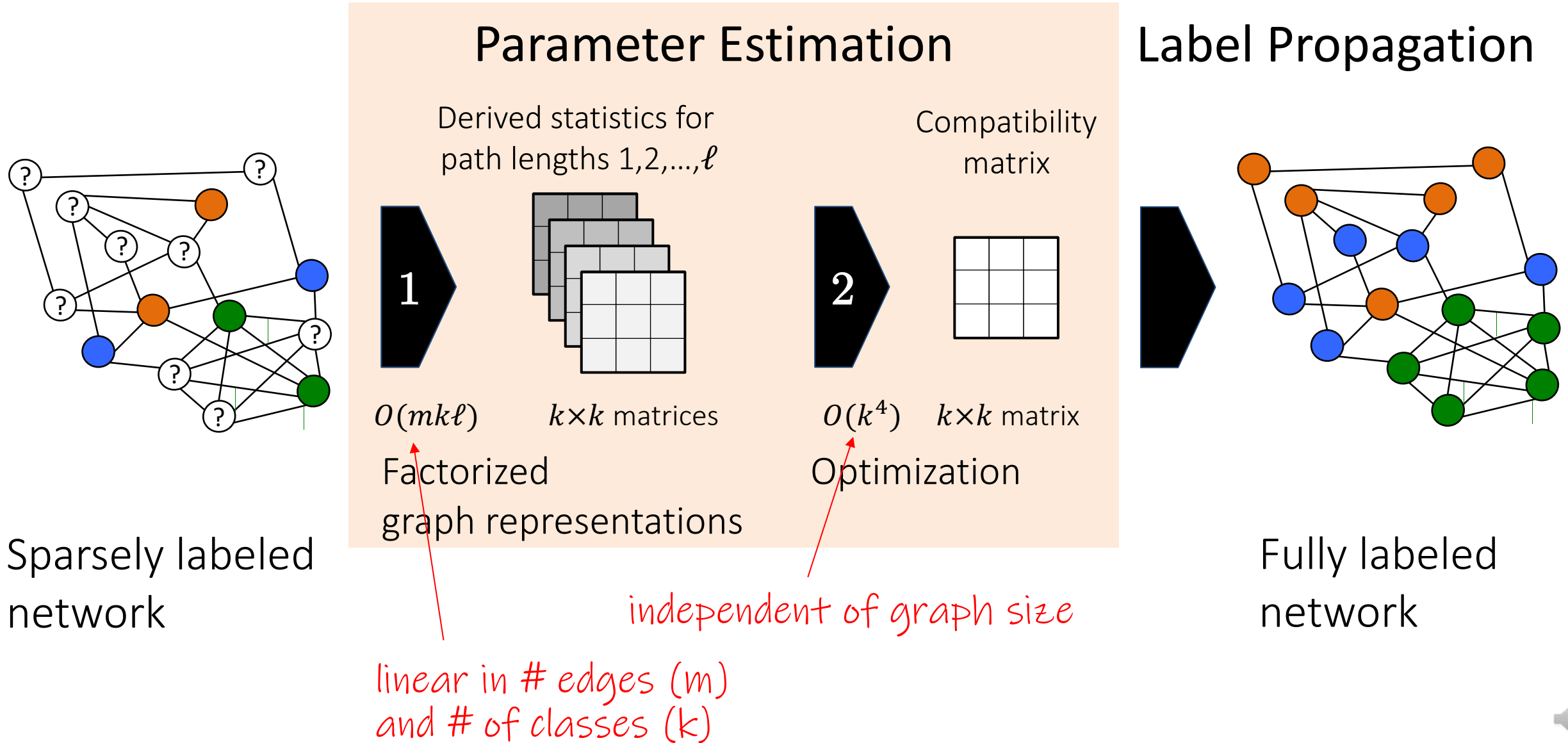
Our method for estimating H needs <5% of the time later needed for labeling 😊



No more need for heuristics or domain experts 😊

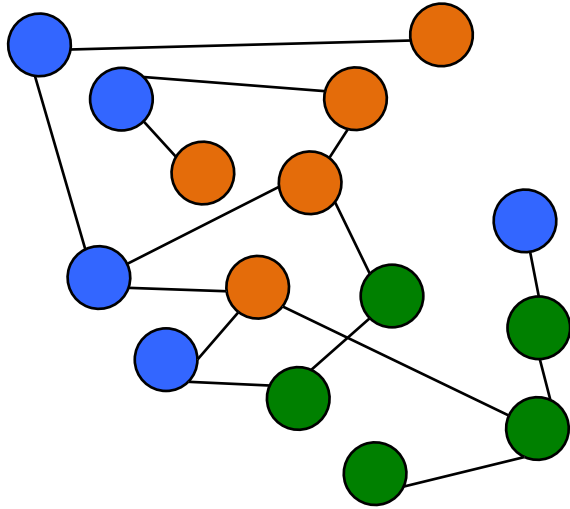
What is the trick?

Splitting parameter estimation into two steps

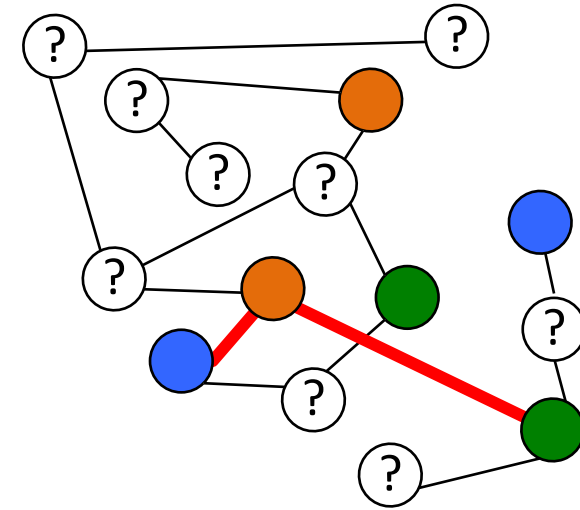


A myopic view: counting relative neighbor frequencies







Fully labeled graph



Sparsely labeled graph






Neighbor count Gold standard compatibilities

$\mathbf{M} =$				\Rightarrow			
	2	6	2		0.2	0.6	0.2
	6	2	2		0.6	0.2	0.2
	2	2	6		0.2	0.2	0.6

normalize

$\Sigma=1$

Labeled neighbor count

$\hat{\mathbf{M}} =$				\Rightarrow	$\hat{\mathbf{H}}$
	0	1	0		
	1	0	1		
	0	1	0		

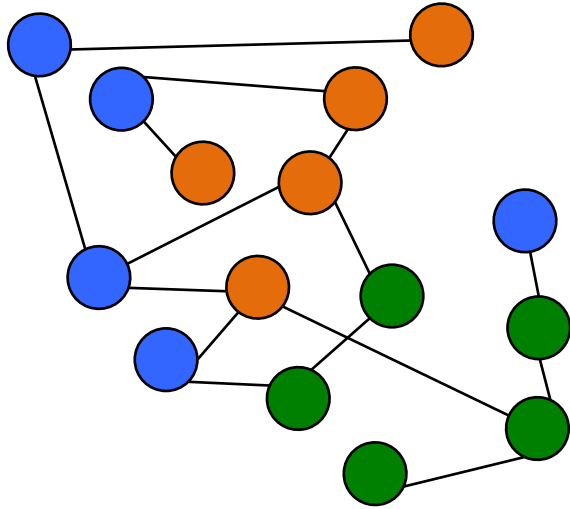
$\Sigma=1$

$\Sigma=2$

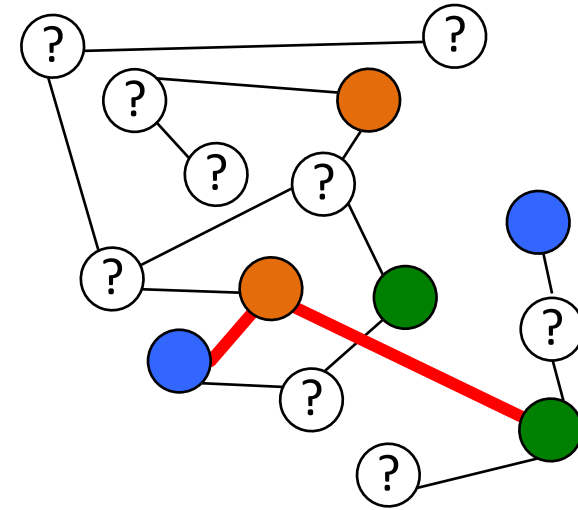
Idea: normalize, then find closest symmetric, doubly-stochastic matrix

A myopic view: counting relative neighbor frequencies

Fully labeled graph



Sparsely labeled graph



Neighbor count

Gold standard compatibilities

$\mathbf{M} =$

	Blue	Orange	Green
Blue	2	6	2
Orange	6	2	2
Green	2	2	6

\Rightarrow

$\mathbf{H} =$

	Blue	Orange	Green
Blue	0.2	0.6	0.2
Orange	0.6	0.2	0.2
Green	0.2	0.2	0.6

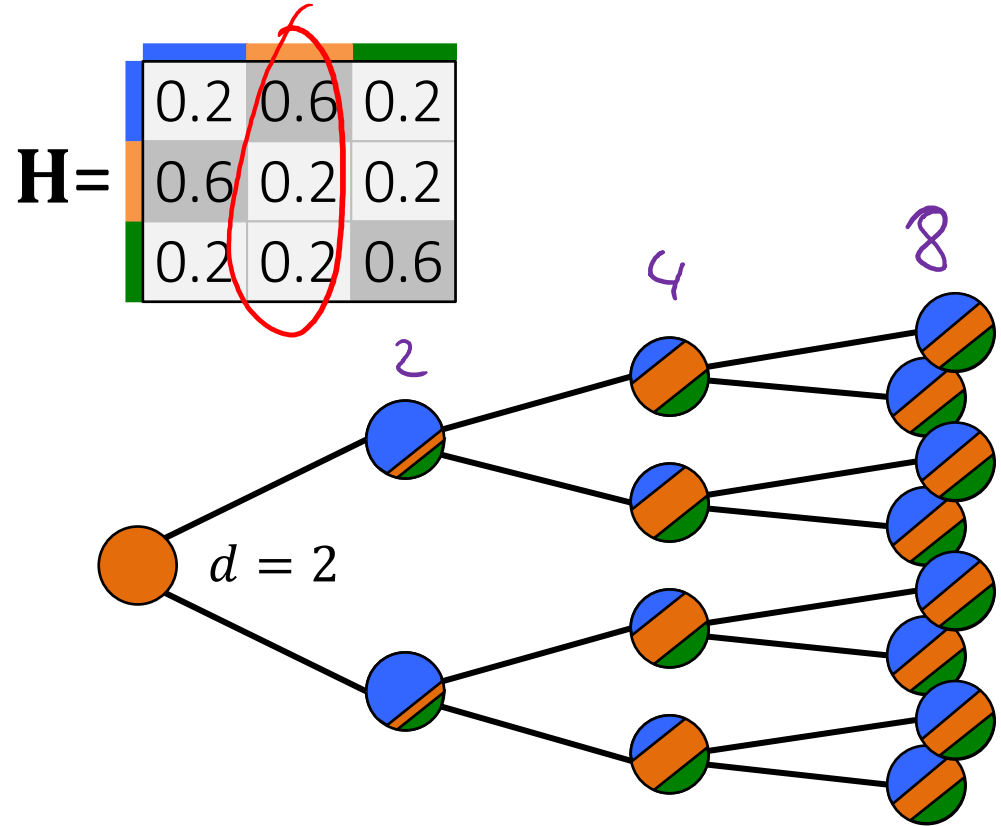
normalize

$\Sigma=1$

Assume $f=10\%$ labeled nodes.
What is the percentage of
edges with labeled end points?

1% 😞 Few nodes \Rightarrow
even fewer edges mf^2

Distant compatibility estimation (DCE)



	$\ell = 1$	$\ell = 2$	$\ell = 3$
0	0.6	0.28	0.38
1	0.2	0.44	0.31
0	0.2	0.28	0.31

Expected signals for neighbors

$\mathbf{H}^2 =$

0.44	0.28	0.28
0.28	0.44	0.28
0.28	0.28	0.44

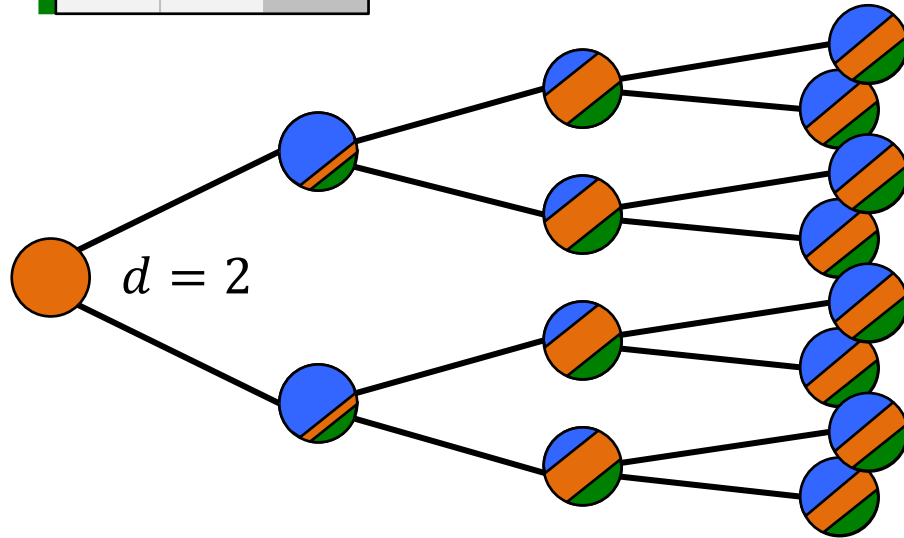
$\mathbf{H}^3 =$

0.31	0.38	0.31
0.38	0.31	0.31
0.31	0.31	0.38

0.6, 0.44, 0.38, 0.35, ...



Distant compatibility estimation (DCE)

$$\mathbf{H} = \begin{array}{|c|c|c|} \hline \text{blue} & \text{orange} & \text{green} \\ \hline 0.2 & 0.6 & 0.2 \\ \hline \text{orange} & 0.6 & 0.2 \\ \hline 0.2 & 0.2 & 0.6 \\ \hline \end{array}$$


	$\ell = 1$	$\ell = 2$	$\ell = 3$	
blue	0	0.6	0.28	0.38
orange	1	0.2	0.44	0.31
green	0	0.2	0.28	0.31

Expected signals for neighbors

graph with:

- m edges
- f fraction labeled nodes
- d node degree

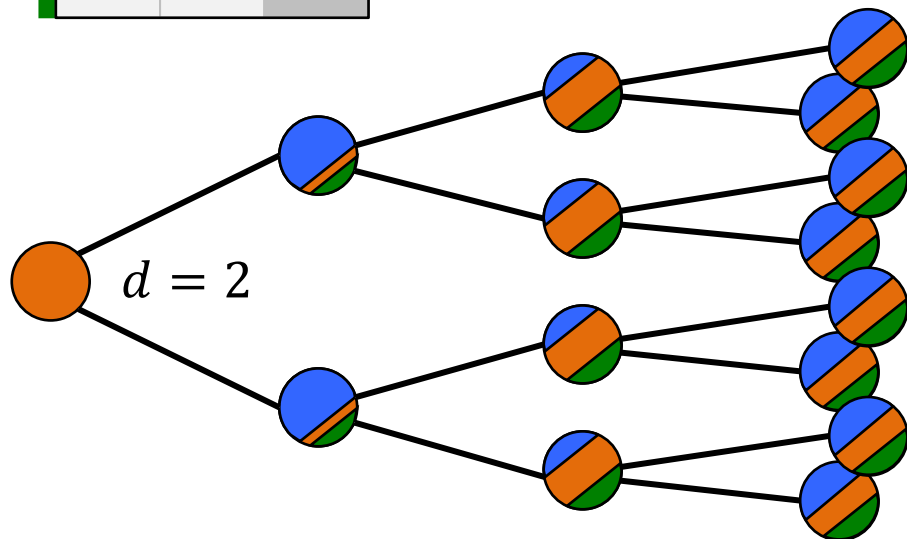
Expected # of labeled neighbors of distance ℓ ?

$d^{\ell-1} m f^2$ expected neighbors of distance ℓ

Idea: amplify the signal from observed length- ℓ paths ☺

Distant compatibility estimation (DCE)

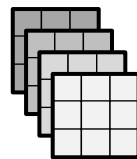
DETAILS

$$\mathbf{H} = \begin{array}{|c|c|c|} \hline \text{blue} & \text{orange} & \text{green} \\ \hline 0.2 & 0.6 & 0.2 \\ \hline 0.6 & 0.2 & 0.2 \\ \hline 0.2 & 0.2 & 0.6 \\ \hline \end{array}$$


	$\ell = 1$	$\ell = 2$	$\ell = 3$	
blue	0	0.6	0.28	0.38
orange	1	0.2	0.44	0.31
green	0	0.2	0.28	0.31

Expected signals for neighbors

distance-smoothed energy function



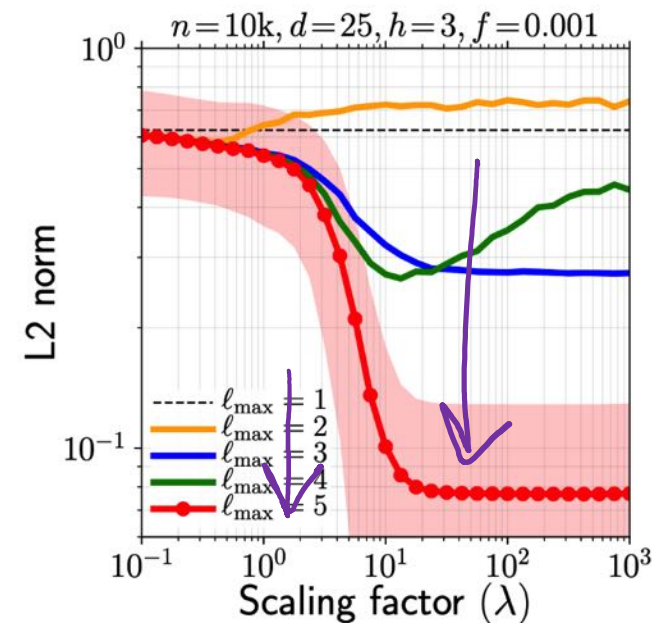
Statistics for path lengths 1, 2, ...

$$E(\mathbf{H}) = \sum_{\ell=1}^{\ell_{\max}} w_{\ell} \|\mathbf{H}^{\ell} - \hat{\mathbf{P}}^{(\ell)}\|^2$$

$$w_{\ell+1} = \lambda w_{\ell} \quad \mathbf{w} = [1, \lambda, \lambda^2, \dots]^T$$

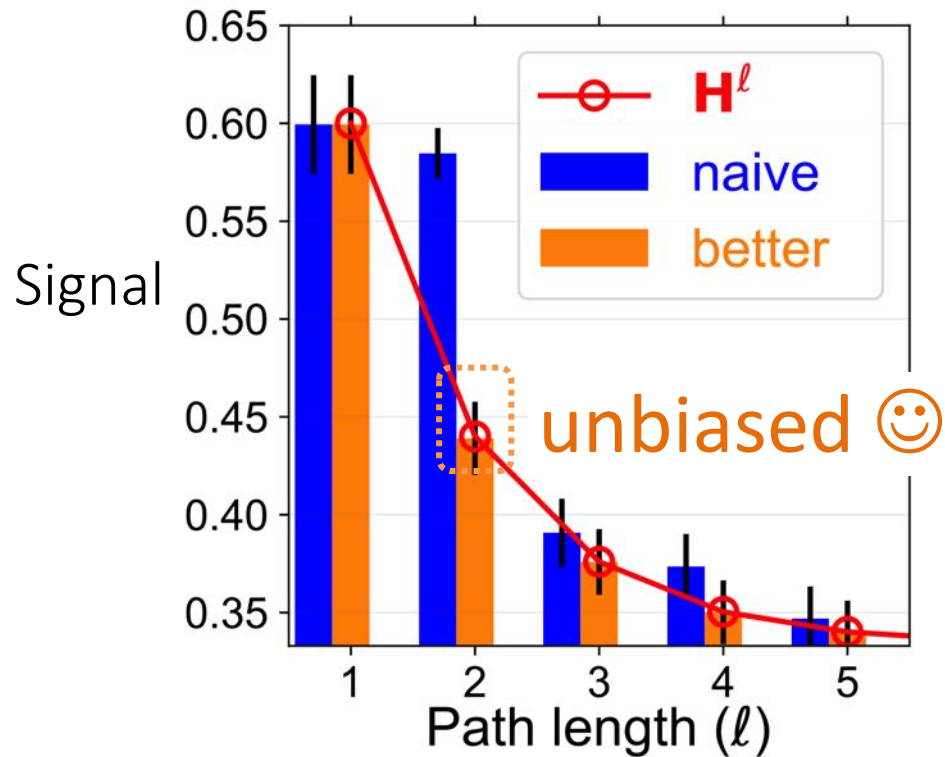
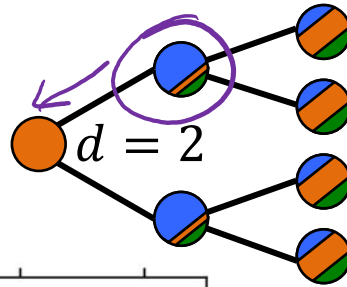
one single hyperparameter ☺

$\|\hat{\mathbf{H}} - \mathbf{H}\|$
(smaller is better)



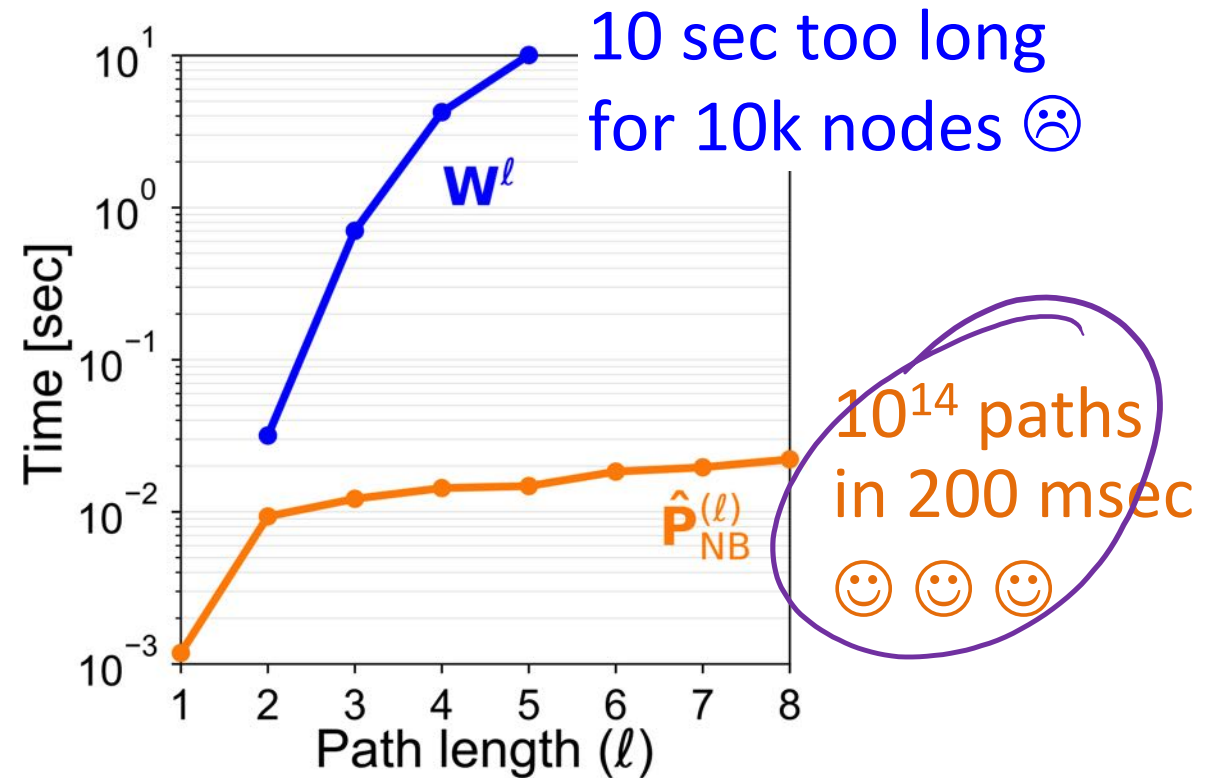
Two technical difficulties

1. Idea from previous page gives **biased estimates** 😞



1. We must **ignore backtracking paths**

2. Calculating longer paths leads to **dense matrix operations** 😞
(\mathbf{W} = sparse adjacency matrix)



2. Requires more careful **re-factorization** of the calculation

→ "factorized graph representations"

Scalable, Factorized Path summation

Details

PROPOSITION 4.2 (NON-BACKTRACKING PATHS). Let $\mathbf{W}_{\text{NB}}^{(\ell)}$ be the matrix with $W_{\text{NB } ij}^{(\ell)}$ being the number of non-backtracking paths of length ℓ from node i to j . Then $\mathbf{W}_{\text{NB}}^{(\ell)}$ for $\ell \geq 3$ can be calculated via following recurrence relation:

$$\mathbf{W}_{\text{NB}}^{(\ell)} = \mathbf{W}\mathbf{W}_{\text{NB}}^{(\ell-1)} - (\mathbf{D} - \mathbf{I})\mathbf{W}_{\text{NB}}^{(\ell-2)} \quad (15)$$

with starting values $\mathbf{W}_{\text{NB}}^{(1)} = \mathbf{W}$ and $\mathbf{W}_{\text{NB}}^{(2)} = \mathbf{W}^2 - \mathbf{D}$. \square

ALGORITHM 4.3 (FACTORIZED PATH SUMMATION). Iteratively calculate the graph summaries $\hat{\mathbf{P}}_{\text{NB}}^{(\ell)}$, for $\ell \in [\ell_{\max}]$ as follows:

- (1) Starting from $\mathbf{N}_{\text{NB}}^{(1)} = \mathbf{W}\mathbf{X}$ and $\mathbf{N}_{\text{NB}}^{(2)} = \mathbf{W}\mathbf{N}_{\text{NB}}^{(1)} - \mathbf{D}\mathbf{X}$, iteratively calculate $\mathbf{N}_{\text{NB}}^{(\ell)} = \mathbf{W}\mathbf{N}_{\text{NB}}^{(\ell-1)} - (\mathbf{D} - \mathbf{I})\mathbf{N}_{\text{NB}}^{(\ell-2)}$.
- (2) Calculate $\mathbf{M}_{\text{NB}}^{(\ell)} = \mathbf{X}^T \mathbf{N}_{\text{NB}}^{(\ell)}$.
- (3) Calculate $\hat{\mathbf{P}}_{\text{NB}}^{(\ell)}$ from normalizing $\mathbf{M}^{(\ell)}$ with Eq. 9.

PROPOSITION 4.4 (FACTORIZED PATH SUMMATION). Algorithm 4.3 calculates all graph statistics $\hat{\mathbf{P}}_{\text{NB}}^{(\ell)}$ for $\ell \in [\ell_{\max}]$ in $\boxed{O(mk\ell_{\max})}$.

Intuition

Relational algebra

$$\pi_{\mathbf{x}}(\mathbf{R}(\mathbf{x}) \bowtie \mathbf{S}(\mathbf{x}, \mathbf{y}))$$

$$\Rightarrow \mathbf{R}(\mathbf{x}) \bowtie \pi_{\mathbf{x}}\mathbf{S}(\mathbf{x}, \mathbf{y})$$

Linear algebra (\mathbf{x} = thin label matrix)

$$(\mathbf{W} \cdot \mathbf{W}) \cdot \mathbf{X}$$

$$\Rightarrow \mathbf{W} \cdot (\mathbf{W} \cdot \mathbf{X})$$

Scalable factorized path summation

Similar ideas of factorized calculation:

- Generalized distributive law
[Aji-McEliece IEEE TIT '00]
- Algebraic path problems
[Mohri JALC'02]
- Valuation algebras
[Kohlas-Wilson AI'08]
- Factorized databases
[Olteanu-Schleich Sigmod-Rec'16]
- FAQ (Functional Aggregate Queries)
[AboKhamis-Ngo-Rudra PODS'16]
- Associative arrays
[Kepner, Jonathan MIT-press'18]
- Optimal ranked enumeration
[Tziavelis+ VLDB'20]

Intuition

Relational algebra

$$\begin{aligned} & \pi_x(R(x) \bowtie S(x, y)) \\ \Rightarrow & R(x) \bowtie \pi_x S(x, y) \end{aligned}$$

Linear algebra (x = thin label matrix)

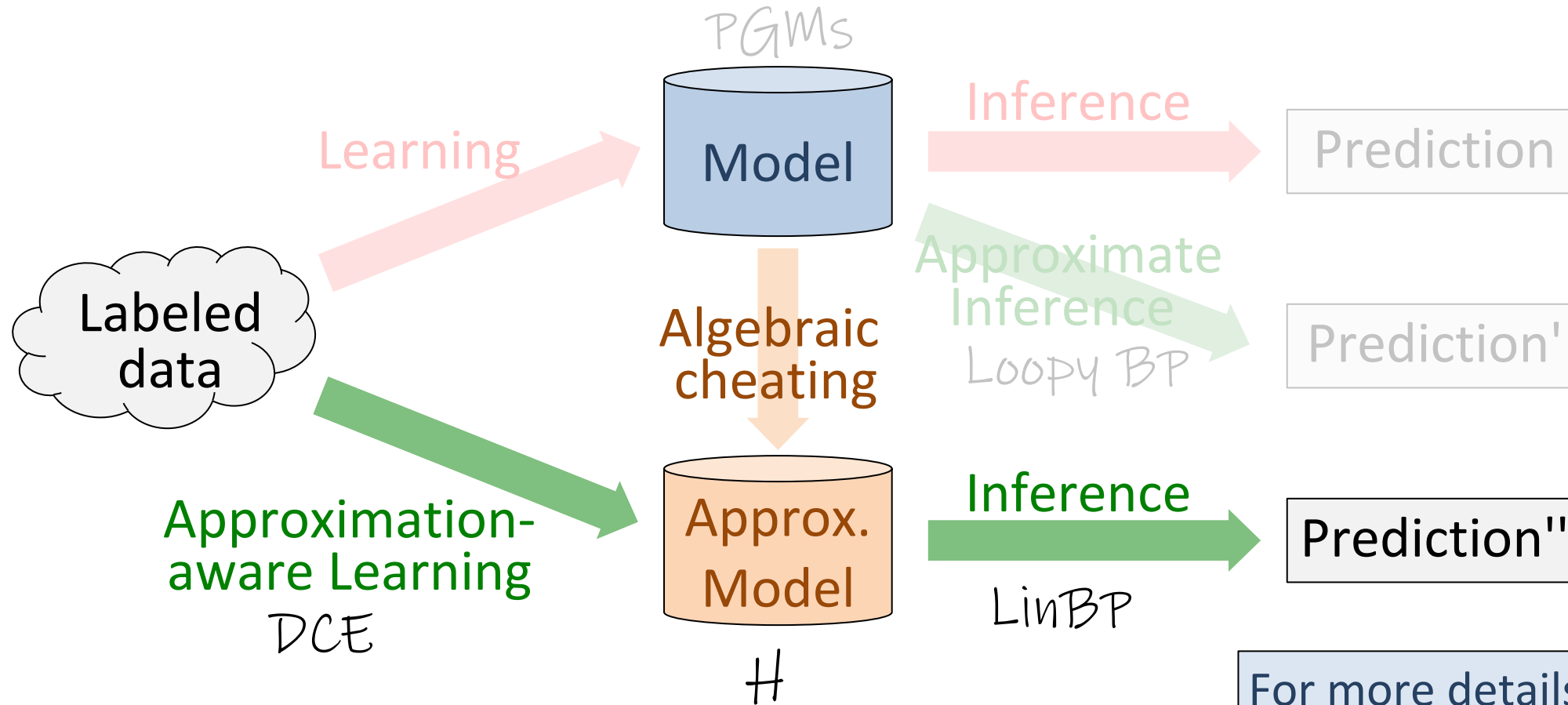
$$\begin{aligned} & (W \cdot W) \cdot X \\ \Rightarrow & W \cdot (W \cdot X) \end{aligned}$$

More details (super happy to discuss further in 1-on-1's)

1. Constrained optimization \rightarrow unconstrained opt. in free parameters
2. Closed form for gradient: gradient-based optimization even faster
3. Random restarts for optimization: but for an optimization on graph sketches, thus independent of n , yet $O(k^4)$
4. Energy-minimization based explanation of LinBP
5. Originally proposed "centering" for LinBP not necessary
6. Proof of unbiased estimator for equal label distribution
7. Non-backtracking paths in factorized calculation that does not require larger $(2m \times 2m)$ "Hashimoto matrix"
8. Lots of experiments on real graphs
9. Even works on graphs without any labeled neighbors 😊

Back to the big picture

"Algebraic cheating" for approximation-aware learning



[Arxiv 2014] Semi-supervised learning with heterophily

[VLDB 2015] Linearized and Single-pass belief propagation

[AAAI 2017] The linearization of pairwise Markov random fields

[VLDBJ 2017] Dissociation and propagation for approximate lifted inference

[UAI 2018] Dissociation-based oblivious bounds for weighted model counting

[SIGMOD 2019] Anytime approximation in probabilistic databases via scaled dissociations

[SIGMOD 2020] Factorized graph representations for semi-supervised learning from sparse data

Supported by [NSF IIS-1762268-CAREER](https://www.nsf.gov/awardsearch/showAward.do?awardNumber=IIS-1762268): Scaling approximate inference and approximation-aware learning

For more details please visit

DATA LAB
@Northeastern

<https://db.khoury.northeastern.edu/>

Thank you 😊