

# Factorized Graph Neural Networks (and the power of algebraic cheating)

*changing the rule of the game*  
("Do I really need to teach PAXOS?")

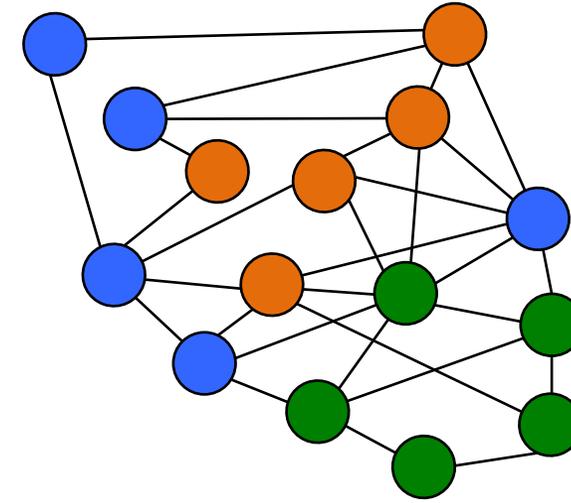
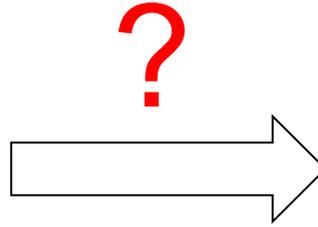
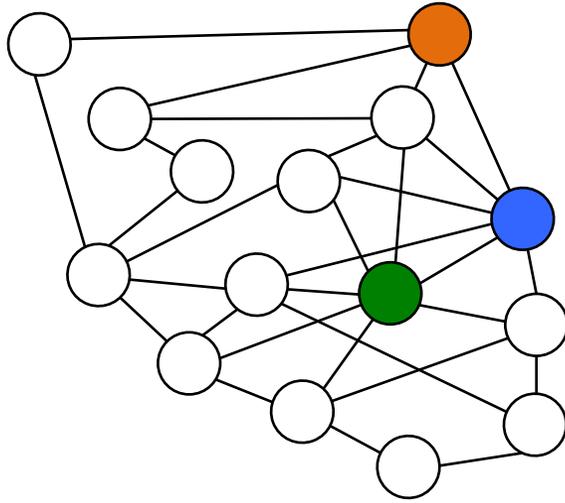
**Wolfgang Gatterbauer**

Based on work with **Krishna Kumar** and **Paul Langton**,  
and earlier work with **Christos Faloutsos**, **Stephan Günnemann**, and **Danai Koutra**

Nov 15, 2023

# Semi-supervised Node classification

*neighbour frequencies*

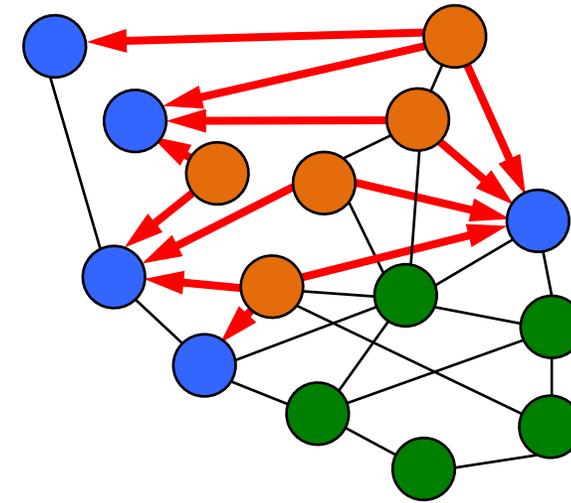
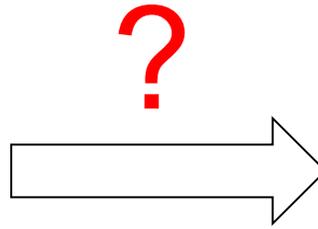
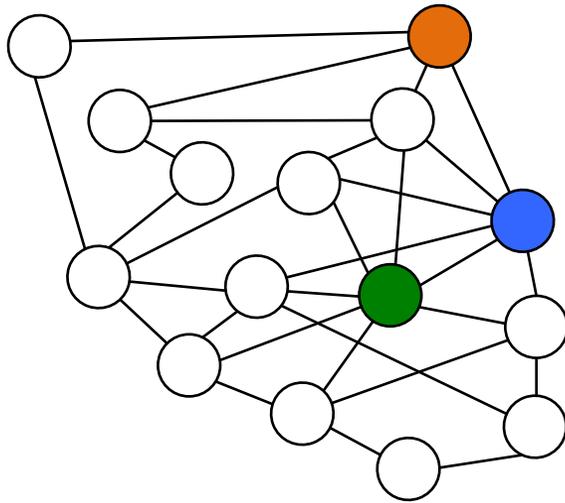


**M**=

4	12	4
12	4	4
4	4	12

PROBLEM: Given a network with labels on some nodes, **what labels should we assign to all other nodes?**

# Semi-supervised Node classification



neighbour frequencies

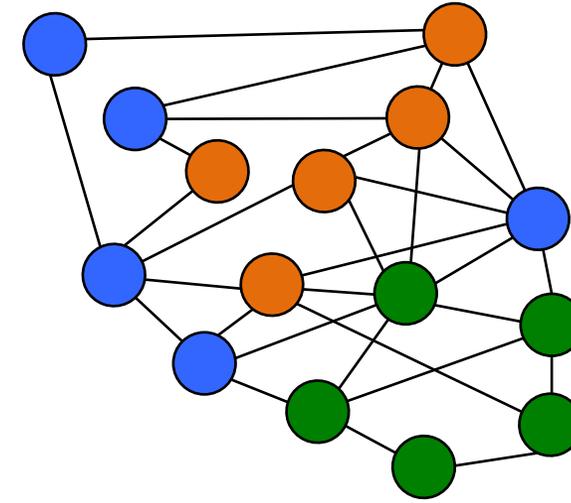
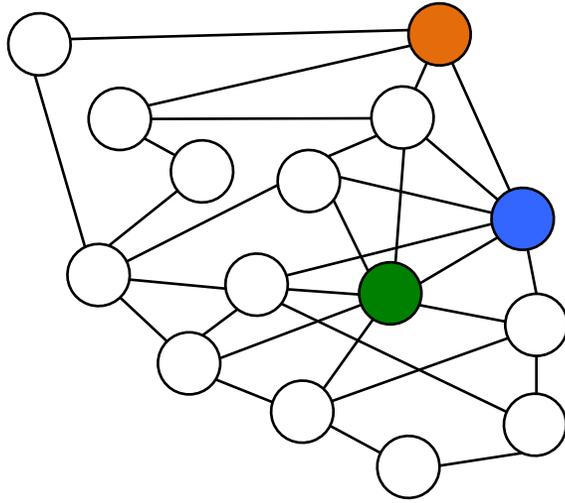
**M** =

4	12	4
12	4	4
4	4	12

PROBLEM: Given a network with labels on some nodes, **what labels should we assign to all other nodes?**

# Semi-supervised Node classification

*neighbour frequencies*



**M**=

4	12	4
12	4	4
4	4	12

**H**=

0.2	0.6	0.2
0.6	0.2	0.2
0.2	0.2	0.6

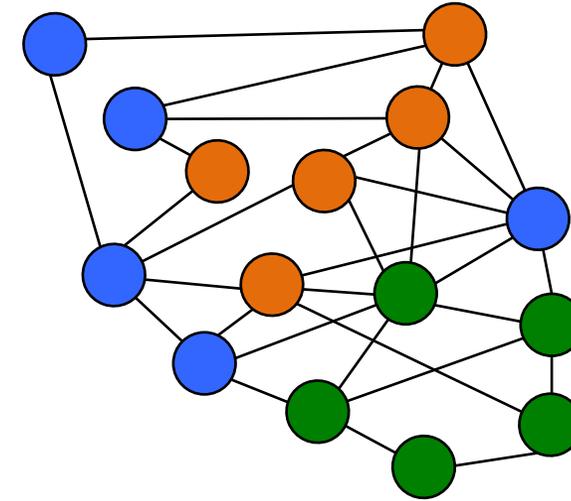
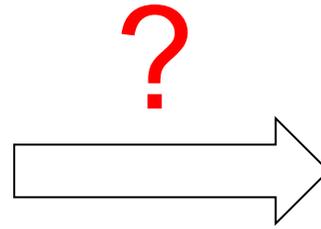
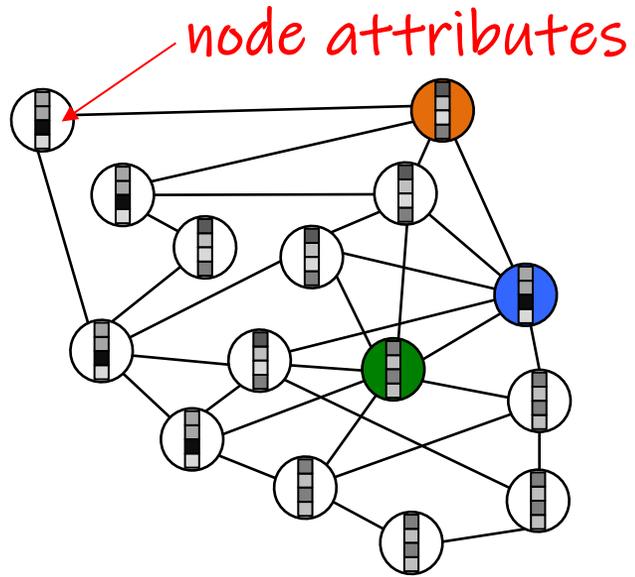
$\Sigma=1$

*"compatibilities"*

PROBLEM: Given a network with labels on some nodes, *what labels should we assign to all other nodes?*

# Semi-supervised Node classification

neighbour frequencies



**M**=

4	12	4
12	4	4
4	4	12

**H**=

0.2	0.6	0.2
0.6	0.2	0.2
0.2	0.2	0.6

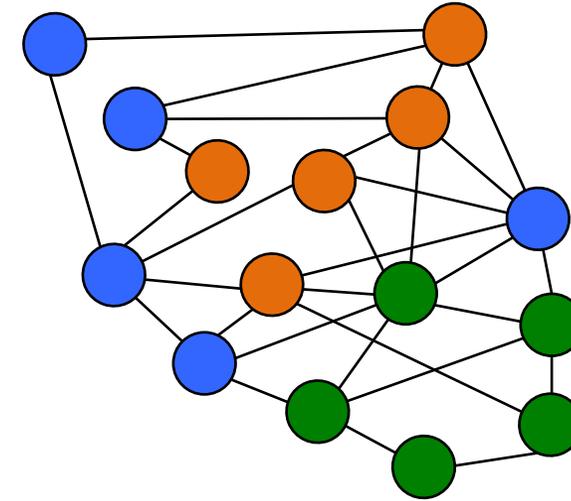
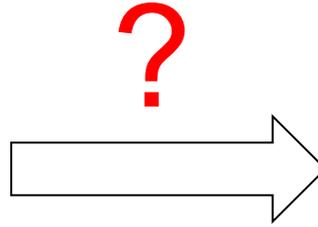
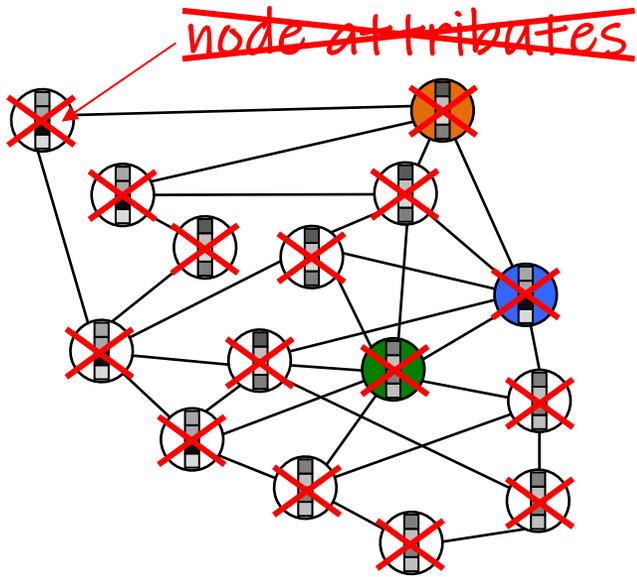
$\Sigma=1$

"compatibilities"

PROBLEM: Given a network with labels on some nodes, what labels should we assign to all other nodes?

# Semi-supervised Node classification

neighbour frequencies


$$\mathbf{M} = \begin{array}{|c|c|c|} \hline 4 & 12 & 4 \\ \hline 12 & 4 & 4 \\ \hline 4 & 4 & 12 \\ \hline \end{array}$$
$$\mathbf{H} = \begin{array}{|c|c|c|} \hline 0.2 & 0.6 & 0.2 \\ \hline 0.6 & 0.2 & 0.2 \\ \hline 0.2 & 0.2 & 0.6 \\ \hline \end{array} \quad \Sigma = 1$$

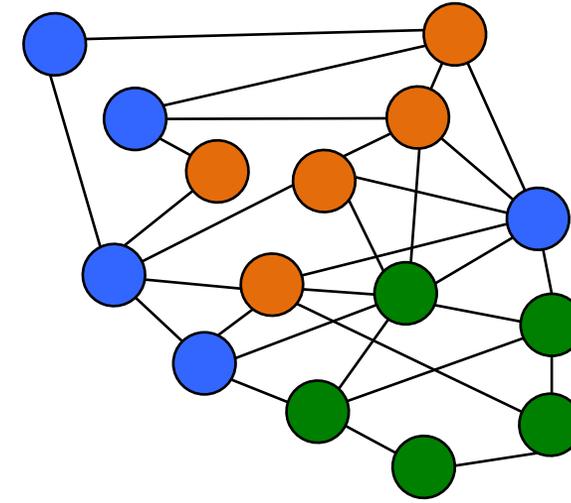
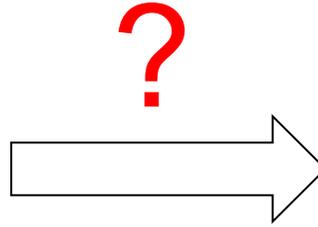
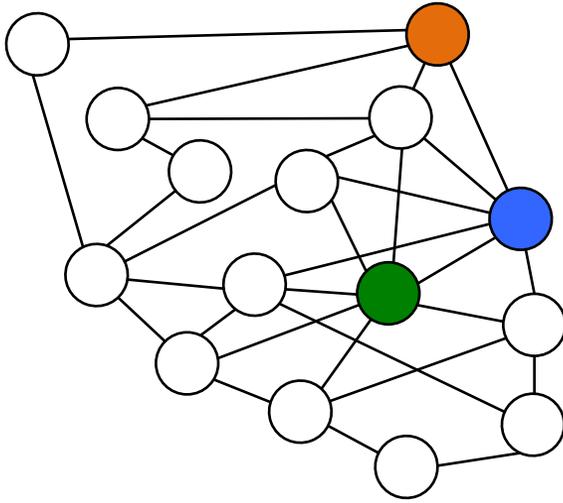
"compatibilities"

PROBLEM: Given a network with labels on some nodes, what labels should we assign to all other nodes?

**NO ATTRIBUTES:** We only use relational information (the graph structure)

# Semi-supervised Node classification

*neighbour frequencies*

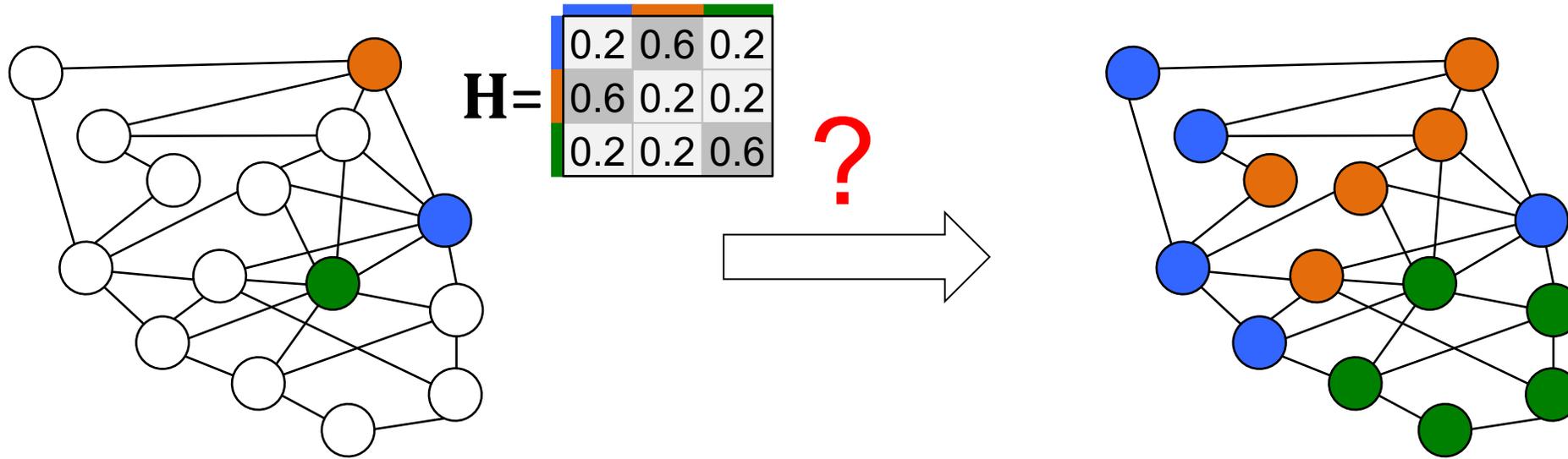

$$\mathbf{M} = \begin{array}{|c|c|c|} \hline 4 & 12 & 4 \\ \hline 12 & 4 & 4 \\ \hline 4 & 4 & 12 \\ \hline \end{array}$$
$$\mathbf{H} = \begin{array}{|c|c|c|} \hline 0.2 & 0.6 & 0.2 \\ \hline 0.6 & 0.2 & 0.2 \\ \hline 0.2 & 0.2 & 0.6 \\ \hline \end{array} \quad \Sigma = 1$$

*"compatibilities"*

PROBLEM: Given a network with labels on some nodes, **what labels should we assign to all other nodes?**

**NO ATTRIBUTES:** We only use relational information (the graph structure)

# Semi-supervised Node classification

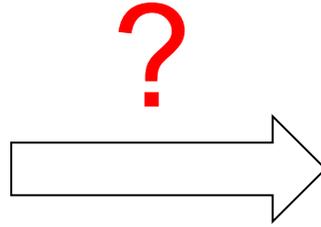
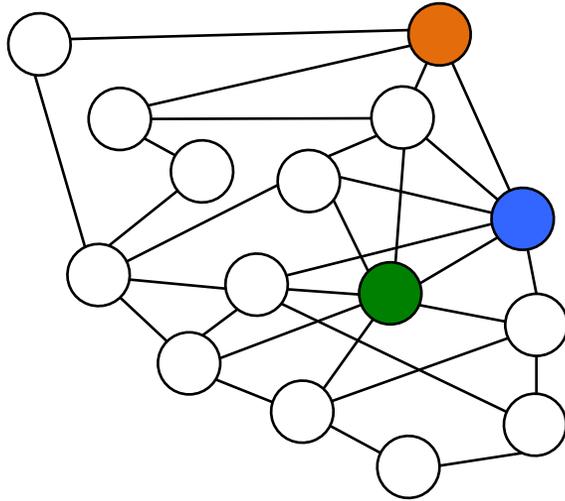


PART 1 (INFERENCE):

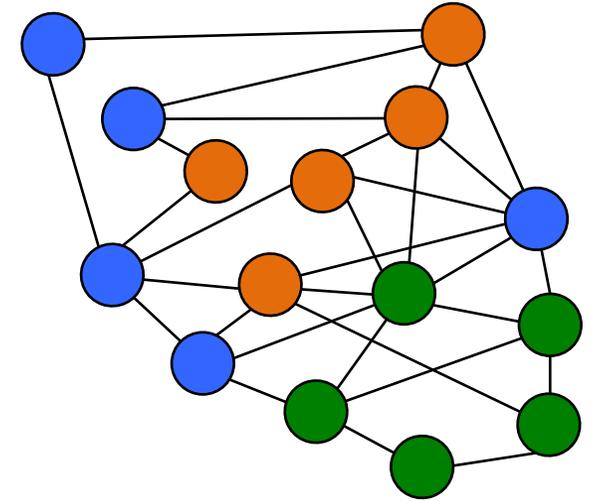
- Given a network  $\mathbf{W}$ ,
- labels on some nodes, and
- compatibilities  $\mathbf{H}$ ,

What labels should we assign to all other nodes?

# Semi-supervised Node classification



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



## PART 1 (INFERENCE):

- Given a network  $\mathbf{W}$ ,
- labels on some nodes, and
- compatibilities  $\mathbf{H}$ ,

What labels should we assign to all other nodes?

## PART 2 (LEARNING):

- Given a network  $\mathbf{W}$ ,
- labels on some nodes, and
- ~~compatibilities  $\mathbf{H}$ ,~~

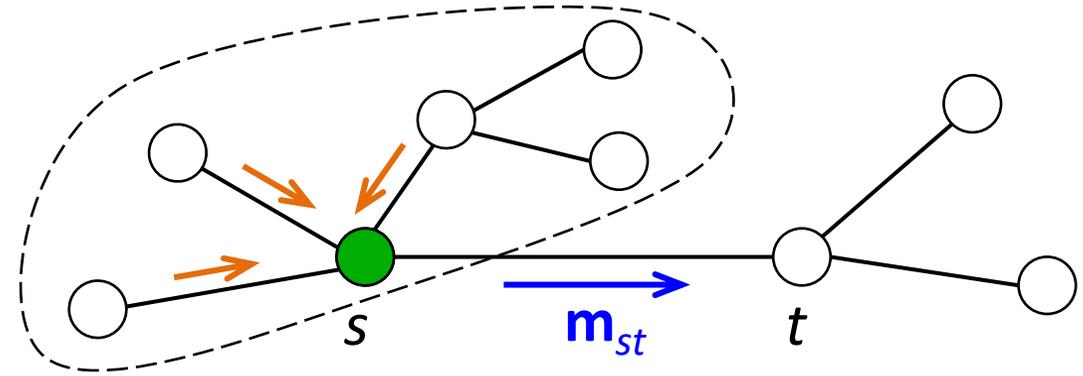
what labels should we assign to all other nodes?

# Approximate Agenda

- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation [VLDB'15]
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations [SIGMOD'20]
- Discussion

# Belief Propagation (BP)

BP is a **Dynamic Programming** (DP) approach to answer conditional probability queries in a **tree-based graphical model**



1) Initialize all message entries to 1

2) Iteratively: calculate messages for each edge and class

$$m_{st}(i) \propto \sum_j H(j, i) x_s(j) \prod_{u \in N(s) \setminus t} m_{us}(j)$$

label-label compatibilities  $\rightarrow$   $H(j, i)$ 
prior beliefs  $\rightarrow$   $x_s(j)$

		$i$		
		0.2	0.6	0.2
$j$	0.2	0.2	0.6	0.2
	0.6	0.6	0.2	0.2
	0.2	0.2	0.2	0.6

label-label compatibilities

$H(j, i)$ : approximately the probability of a node being in state  $i$  given that it has a neighbour in state  $j$

3) After messages converge: calculate final beliefs

$$f_s(i) \propto x_s(i) \prod_{u \in N(s)} m_{us}(j)$$

final beliefs  $\rightarrow$   $f_s(i)$

# Problems with BP (when applied to real graphs with loops)

BP applied as a heuristics to graphs with cycles ("Loopy BP") is difficult to work with 😞

## Collective Classification in Network Data

*Prithviraj Sen, Galileo Namata, Mustafa Bilgic,  
Lise Getoor, Brian Gallagher,  
and Tina Eliassi-Rad*

AI magazine 2008

Cited > 3600 times (11/2023)

cant and useful. However, the LBP accuracy had a sudden drop when the graph became very dense. The reason behind this result is the well known fact that LBP has convergence issues when there are many short loops in the graph.

gorithms. First, although MF and LBP performance is in some cases a bit better than that of ICA and GS, MF and LBP were also the most difficult to work with in both learning and inference. Choosing the initial weights so that the weights will converge during training is nontrivial. Most of the time, we had to initialize the weights with the weights we got from ICA in order to get the algorithms to converge. Thus, the results reported from

# Problems with BP (when applied to real graphs with loops)

BP applied as a heuristics to graphs with cycles ("Loopy BP") is difficult to work with 😞

Our solution for part 1:

1. Linearize and thereby simplify Belief Propagation (it becomes "**algebraically convenient**")
2. Turns out to generalize semi-supervised learning from smoothness (incl. PageRank) to **heterophily**
3. In more modern language: an infinitely deep graph neural network with tied parameters and removed non-linearities, and no "oversmoothing"



# Key Ideas: 1) Centering + 2) Linearizing BP

Original Value = Center point + Residual ①

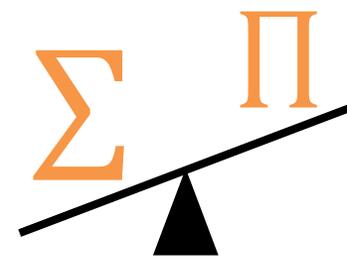
label-label compatibilities  $\mathbf{H}$   $\begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.6 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}$  =  $\begin{bmatrix} 1 \\ k \end{bmatrix}_{k \times k}$   $\begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$  +  $\hat{\mathbf{H}}$   $\begin{bmatrix} -0.13 & 0.26 & -0.13 \\ 0.26 & -0.13 & -0.13 \\ -0.13 & -0.13 & 0.26 \end{bmatrix}$

beliefs  $\mathbf{x}, \mathbf{f}$   $\begin{bmatrix} 0.2 \\ 0.6 \\ 0.2 \end{bmatrix}$  =  $\frac{1}{k} \mathbf{1}$   $\begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$  +  $\hat{\mathbf{x}}, \hat{\mathbf{f}}$   $\begin{bmatrix} -0.13 \\ 0.26 \\ -0.13 \end{bmatrix}$

messages  $\mathbf{m}$   $\begin{bmatrix} 1.1 \\ 0.8 \\ 1.1 \end{bmatrix}$  =  $\mathbf{1}$   $\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$  +  $\hat{\mathbf{m}}$   $\begin{bmatrix} 0.1 \\ -0.2 \\ 0.1 \end{bmatrix}$

②	Expression	Maclaurin series	Approximation
Logarithm	$\ln(1 + \epsilon)$	$= \epsilon - \frac{\epsilon^2}{2} + \frac{\epsilon^3}{3} - \dots$	$\approx \epsilon$
Division	$\frac{\frac{1}{k} + \epsilon_1}{1 + \epsilon_2}$	$= \left(\frac{1}{k} + \epsilon_1\right)(1 - \epsilon_2 + \epsilon_2^2 - \dots)$	$\approx \frac{1}{k} + \epsilon_1 - \frac{\epsilon_2}{k}$

# Intuition behind Centering and Linearizing



Center

component-wise multiplication

$$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \odot \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.25 \end{bmatrix} \propto \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

$\Sigma=0.50$   $\Sigma=1$

Value

$$\begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix} \odot \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.36 \\ 0.16 \end{bmatrix} \propto \begin{bmatrix} 0.69 \\ 0.31 \end{bmatrix}$$

$\Sigma=0.52$   $\Sigma=1$

Residual

$$\begin{bmatrix} 0.1 \\ -0.1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.1 \end{bmatrix} = \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} \Rightarrow \begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix}$$

$\Sigma=0$   $\Sigma=1$

Approx.

Summation instead of multiplication!  
No more normalization necessary 😊

# Matrix formulation of LinBP

## Update equation

$(n \times k)$ , here  $k=3$  labels

$$\mathbf{F} \leftarrow \mathbf{X} + \mathbf{W} \cdot \mathbf{F} \cdot \mathbf{H}$$

final beliefs      explicit beliefs      graph      compatibility matrix

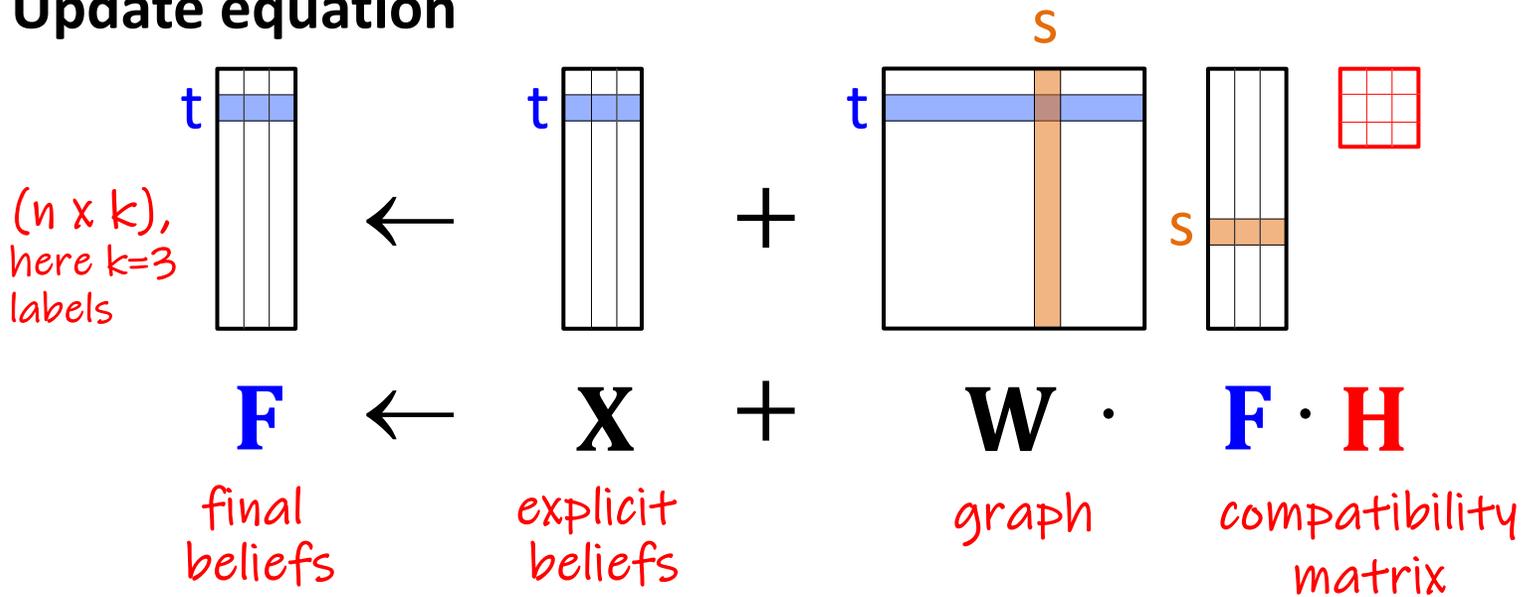
## Compare to Personalized PageRank

$$\mathbf{f} \leftarrow \bar{\alpha} \cdot \mathbf{x} + \alpha \cdot \mathbf{W}^{\text{col}} \cdot \mathbf{f}$$

Basically a generalization of Katz centrality!

# Matrix formulation of LinBP

## Update equation



## Closed form

$$\text{vec}(\mathbf{F}) = (\mathbf{I} - \mathbf{H} \otimes \mathbf{W})^{-1} \text{vec}(\mathbf{X})$$

Annotations:

- $\otimes$ : Kronecker product
- $\text{vec}(\mathbf{X})$ : vectorization of matrix: stacks columns on top of each other

## Convergence

Spectral radius of (...)  $< 1$

Scale with "appropriate"  $\varepsilon$ :  $\mathbf{H}' \leftarrow \varepsilon \mathbf{H}$

Scaling factor:  $\varepsilon = s \cdot \varepsilon^*$  ( $\varepsilon^*$  convergence boundary)

# Geometric sums (intuition for closed-form)

Recall Javier's talk 😊

$$S = 1 + x + x^2 + \dots = \frac{1}{1-x} = (1-x)^{-1}, \quad |x| < 1$$

$$x = \frac{1}{2} \quad S = 1 + \frac{1}{2} + \frac{1}{4} + \dots = \frac{1}{1 - \frac{1}{2}} = 2$$

$$x = -\frac{1}{2} \quad S = 1 - \frac{1}{2} + \frac{1}{4} - \dots = \frac{1}{1 + \frac{1}{2}} = \frac{2}{3}$$

$$x = 2 \quad S = 1 + 2 + 4 + \dots \neq \frac{1}{1-2} = -1$$

# LinBP leads to very concise code (Python)

## BP (Belief Propagation)

```
487 # --- Actual loop: each loop calculates (a) the new messages (with damping) and (b) the new beliefs
488 converged = False
489 actualNumIt = -1 # iterations start with 0th iteration
490 while actualNumIt < numMaxIt and not converged:
491     actualNumIt += 1
492
493 # --- (a) calculate messages
494 if actualNumIt == 0:
495     # --- first pass (counts as 0th iteration): create message dictionaries and initialize messages with ones
496     dict_messages_along_1 = {} # dictionary: messages for each edge (i->j) in direction i->j
497     dict_messages_against_1 = {} # dictionary: messages for each edge (i<-j) in direction i->j
498     default = np.ones(k) # first message vector: all 1s
499     for (i,j) in zip(row, col):
500         dict_messages_along_1[(i,j)] = default
501         dict_messages_against_1[(j,i)] = default
502 else:
503     # --- other iterations: calculate "messages_new" using message-passing with division (from F and messages)
504     dict_messages_along_2 = {} # new dictionary: messages for each edge (i->j) in direction i->j
505     dict_messages_against_2 = {} # new dictionary: messages for each edge (i<-j) in direction i->j
506     for (i,j) in dict_messages_along_1.keys(): # also includes following case: "for (j,i) in dict_messages_against_1.keys()"
507         if dim_pot == 3: # need to reference the correct potential in case dim_pot == 3
508             Pot = P[dict_edges_pot[(i,j)]-1, :, :]
509             dict_messages_along_2[(i,j)] = (F2[i] / dict_messages_against_1[(j,i)].dot(Pot)) # entry-wise division
510             dict_messages_against_2[(j,i)] = (F2[j] / dict_messages_along_1[(i,j)].dot(Pot.transpose()))
511         # TODO above two lines can contain errors
512
513 # --- assign new to old message dictionaries, and optionally damp messages
514 if damping == 1:
515     dict_messages_along_1 = dict_messages_along_2.copy() # requires shallow copy because of later division
516     dict_messages_against_1 = dict_messages_against_2.copy()
517 else:
518     for (i,j) in dict_messages_along_1.keys():
519         dict_messages_along_1[(i,j)] = damping*dict_messages_along_2[(i,j)] + \
520             (1-damping)*dict_messages_along_1[(i,j)]
521     for (i,j) in dict_messages_against_1.keys():
522         dict_messages_against_1[(i,j)] = damping*dict_messages_against_2[(i,j)] + \
523             (1-damping)*dict_messages_against_1[(i,j)]
524
525 # --- (b) create new beliefs by multiplying prior beliefs with all incoming messages (pointing in both directions)
526 for (i, f) in enumerate(F2):
527     if not clamping or implicitVector[i] == 0: # only update beliefs if those are not explicit and clamped
528         F2[i] = X0[i] # need to start multiplying from explicit beliefs, referencing the row with separate variable did not work out
529         for j in dict_edges_out[i]: # edges pointing away
530             F2[i] *= dict_messages_against_1[(j,i)]
531         for j in dict_edges_in[i]: # edges pointing inwards
532             F2[i] *= dict_messages_along_1[(j,i)]
533         # TODO line can contain errors
534
535 # --- normalize beliefs [TODO: perhaps remove later to optimize except in last round]
536 F2 = row_normalize_matrix(F2, norm='l1')
```

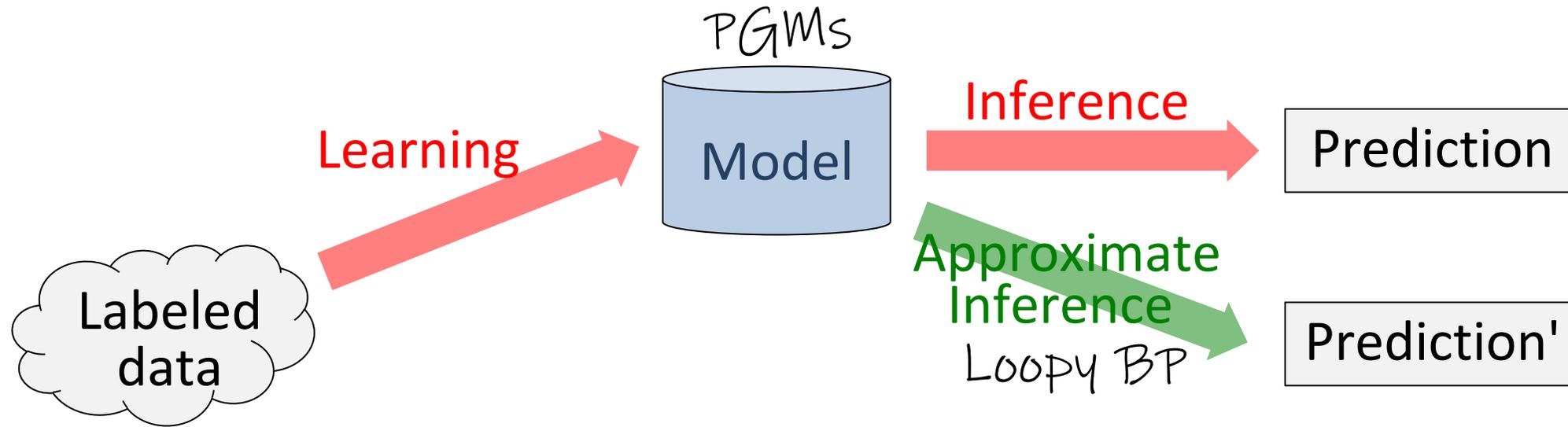
## LinBP

```
155 while actualNumIt < numMaxIt and not converged:
156     actualNumIt += 1
157
158 # --- Calculate new beliefs
159 if not converged:
160     F = X + W.dot(F).dot(H)
161 else:
162     if not compensation:
163         F = X + W.dot(F).dot(H) - D.dot(F).dot(H2) # W.dot(F) is short form for: sparse.csr_matrix.dot(W, F)
164     else:
165         F = X + W.dot(F).dot(H_star) - D.dot(F).dot(H_star2)
```

$$F = X + W \cdot \text{dot}(F) \cdot \text{dot}(H)$$

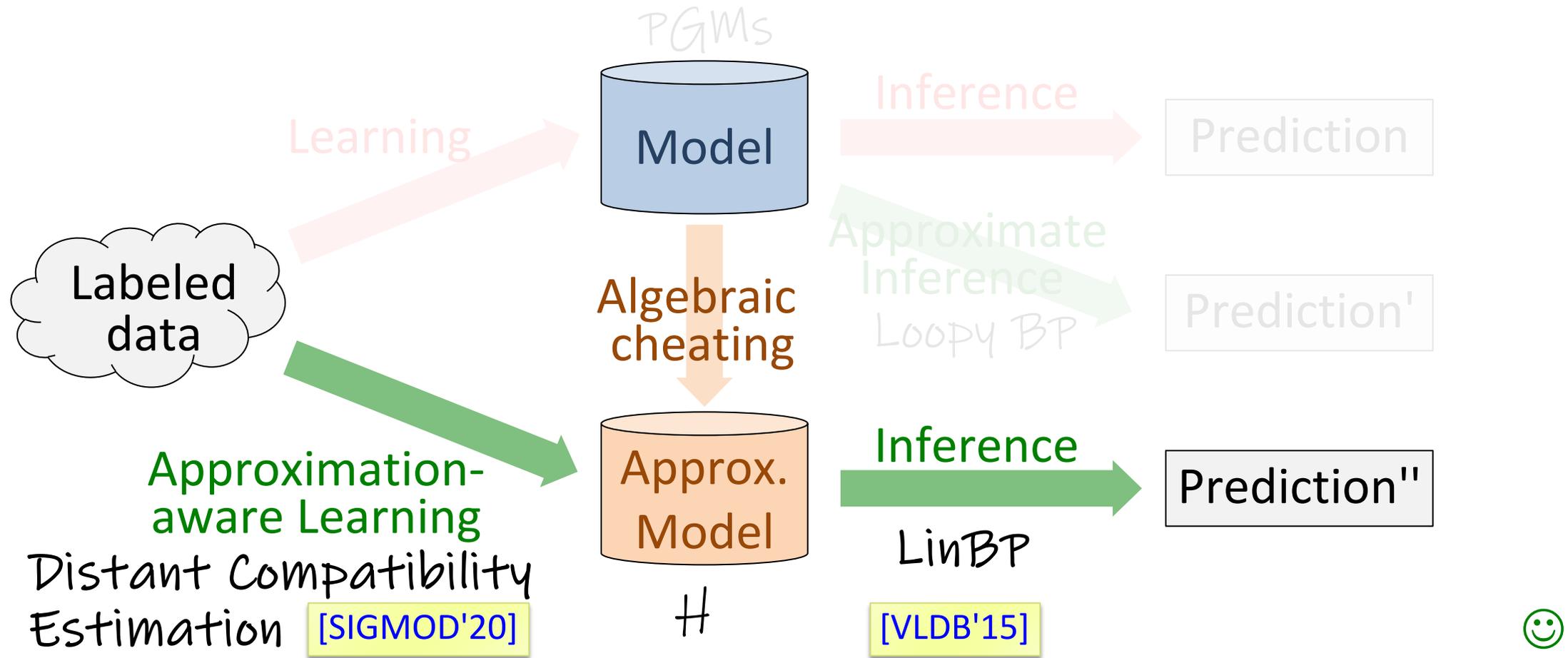
# "Algebraic cheating" for approximation-aware learning

That goes against all the ideas from efficient knowledge compilation ☹️



# "Algebraic cheating" for approximation-aware learning

That goes against all the ideas from efficient knowledge compilation ☹️



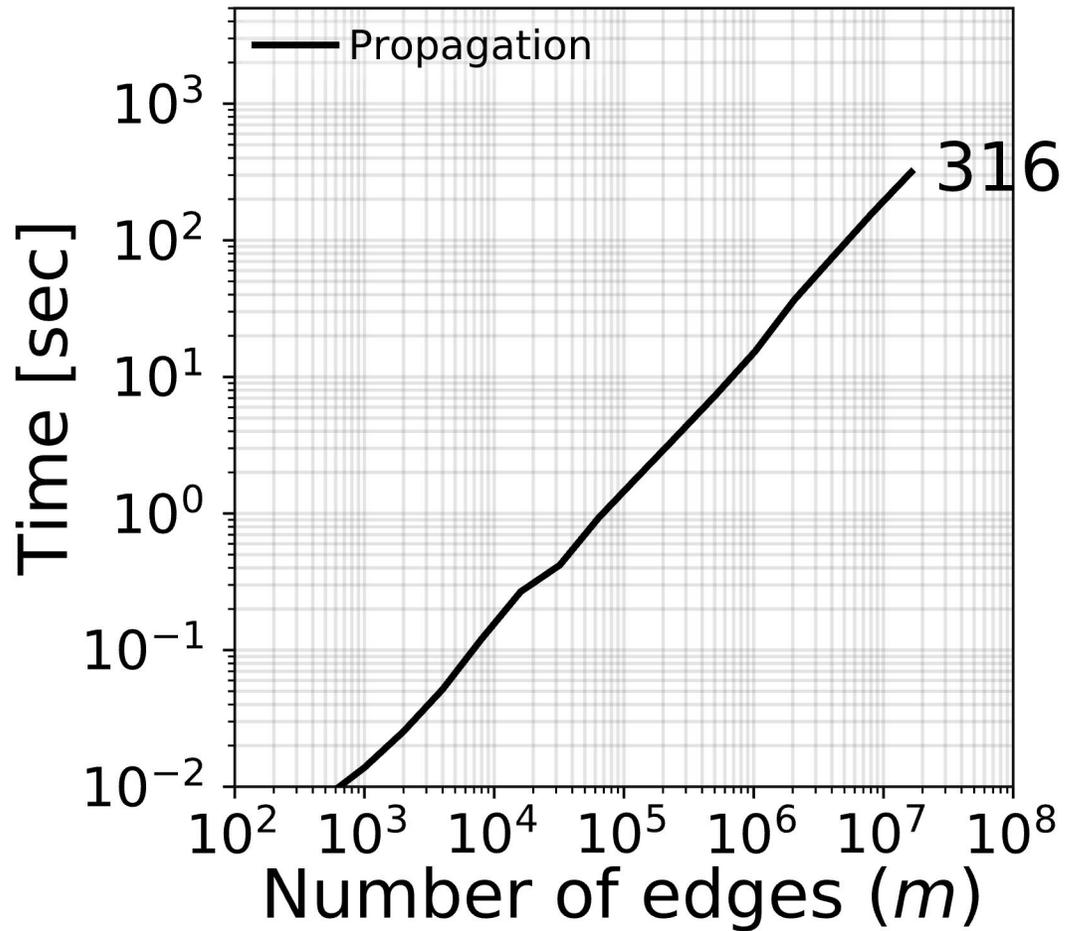
[SIGMOD'20]: "Factorized Graph Representations for Semi-Supervised Learning from Sparse Data", Kumar, Langton, Gatterbauer. SIGMOD'20. <https://doi.org/10.1145/3318464.3380577>

[VLDB'15]: "Linearized and single-pass belief propagation", Gatterbauer, Günemann, Koutra, Faloutsos. VLDB'15. <https://doi.org/10.14778/2735479.2735490>

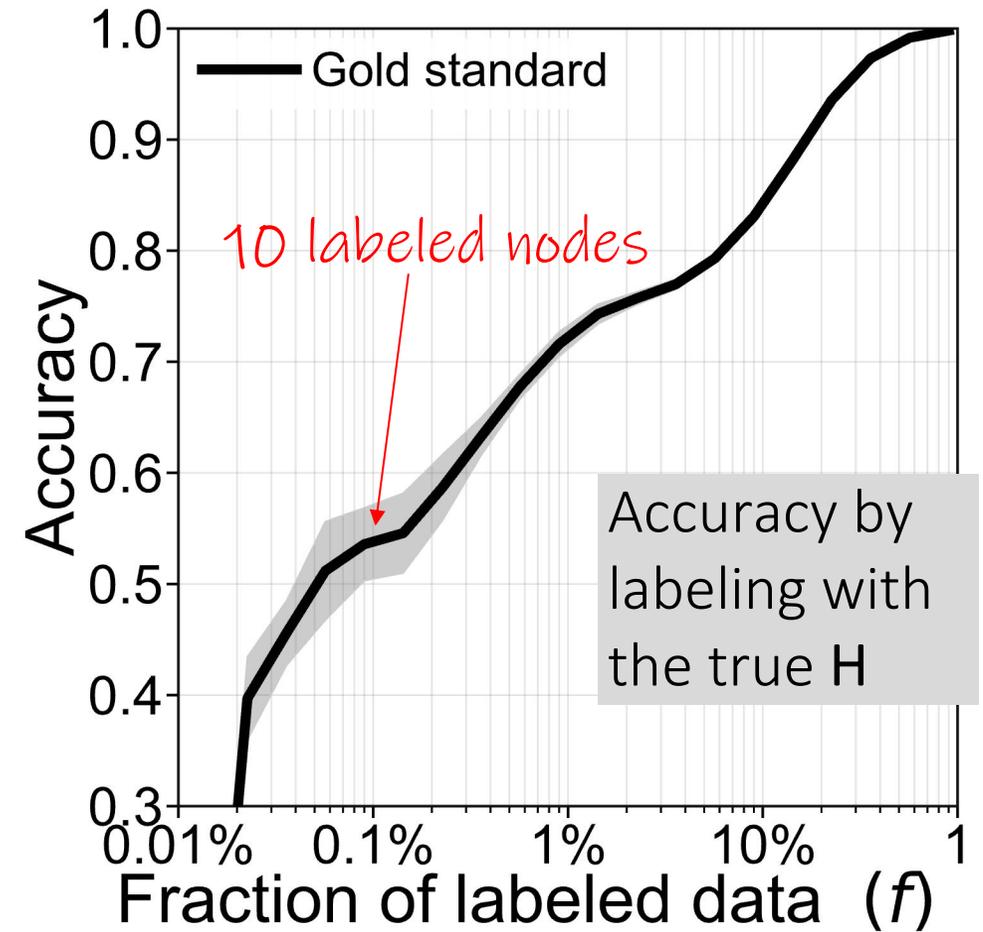
# Approximate Agenda

- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation [VLDB'15]
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations [SIGMOD'20]
  - How well does it work?
  - What is the magic sauce?
  - What we would like to do (but it does not work)
  - What we actually do (Distant Compatibility Estimation)
- Discussion

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



Label propagation linear in # edges

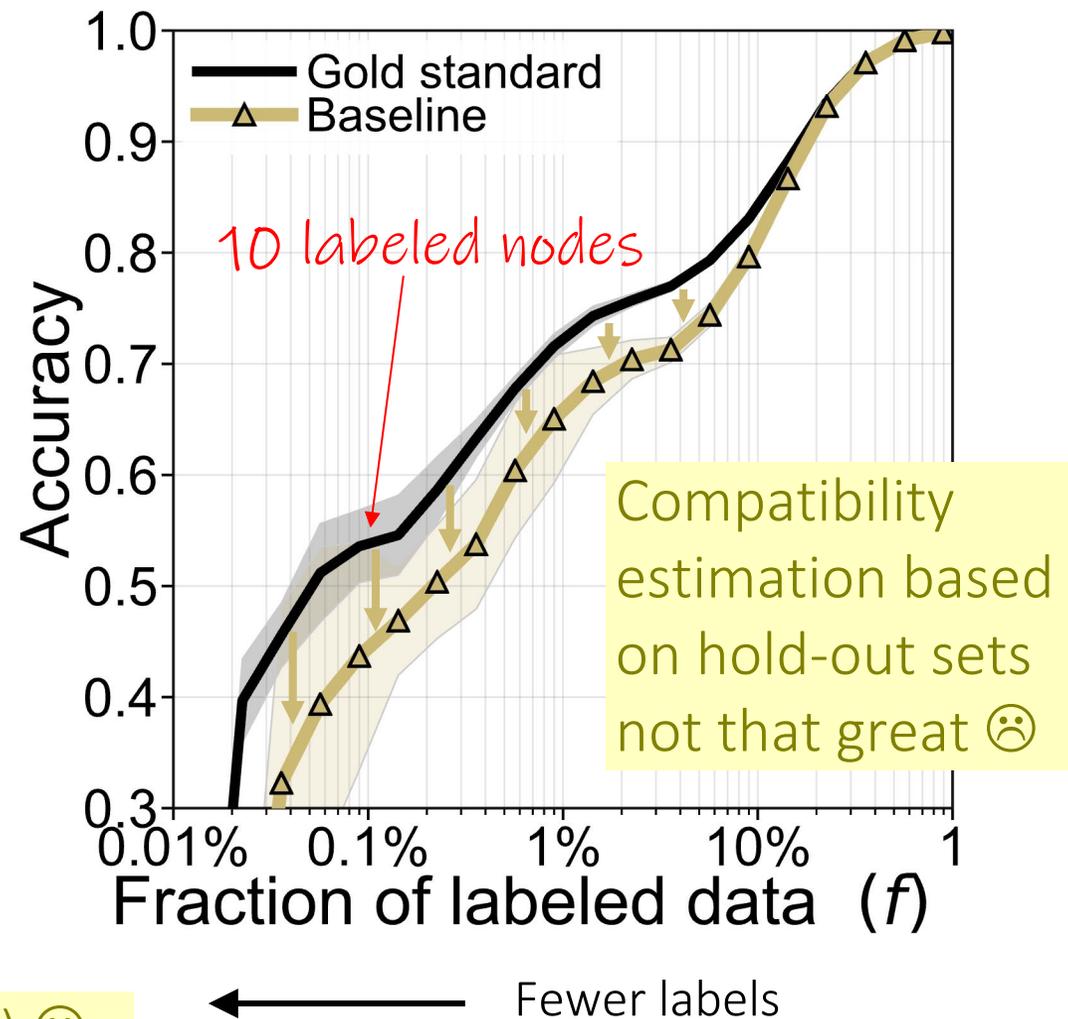
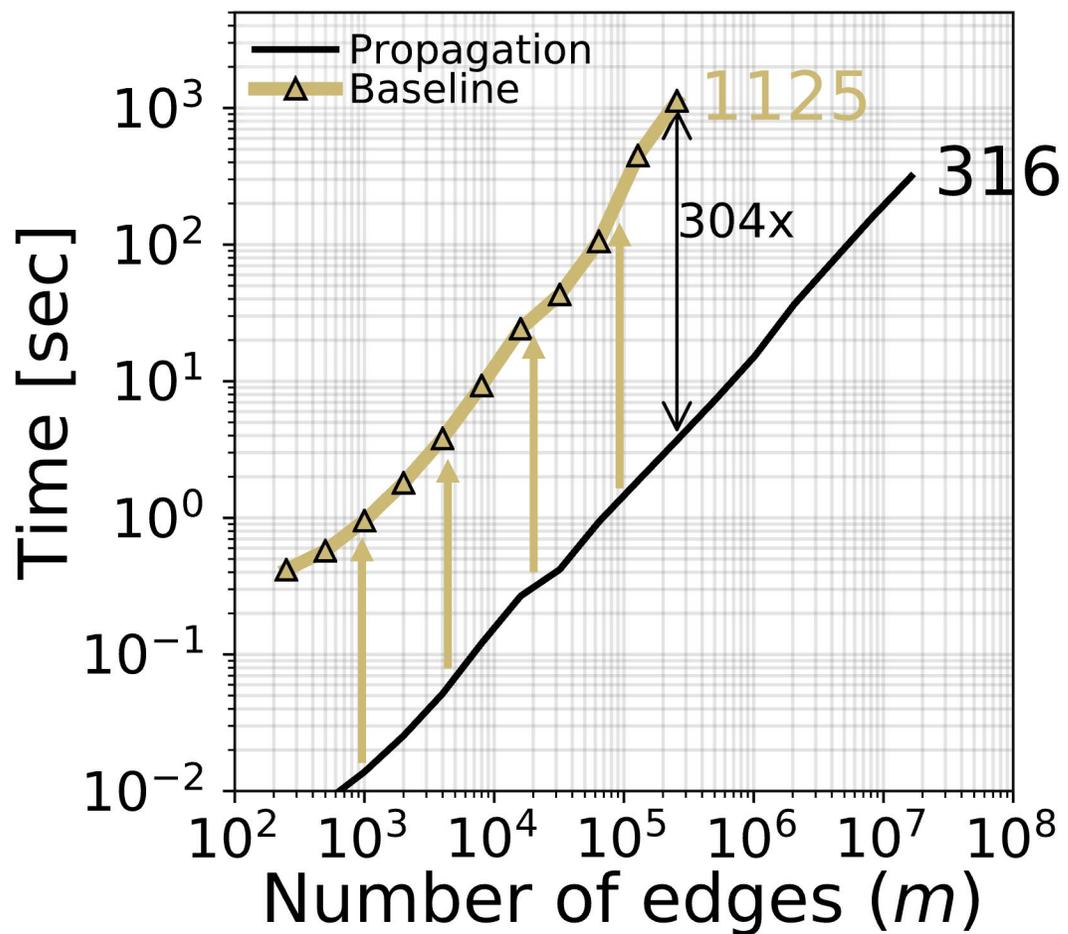


← Fewer labels

Details: 10,000 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 ☹️

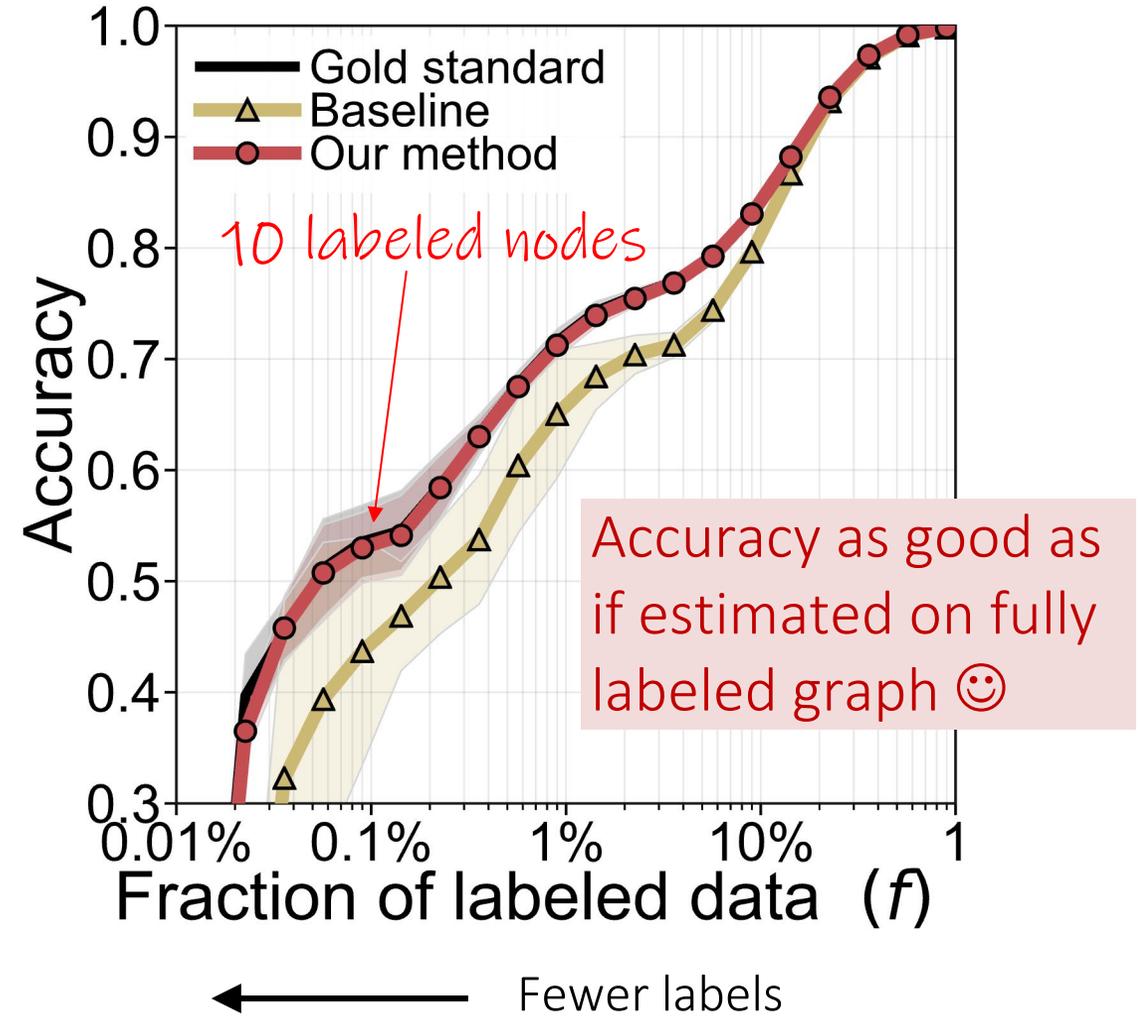
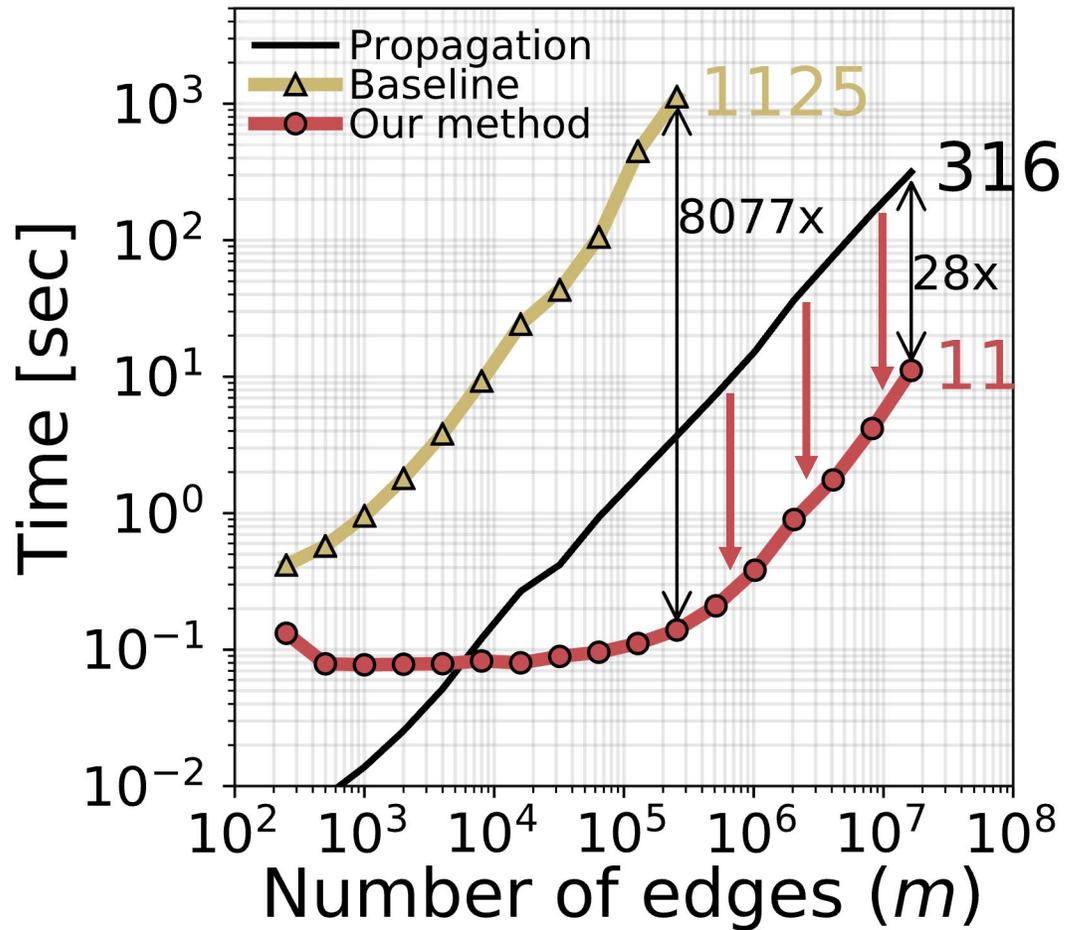


Estimation uses inference as subroutine (thus slower) ☹️

# Time and Accuracy with our method 😊



ACM SIGMOD 2021 Reproducibility Award for papers from SIGMOD 2020



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

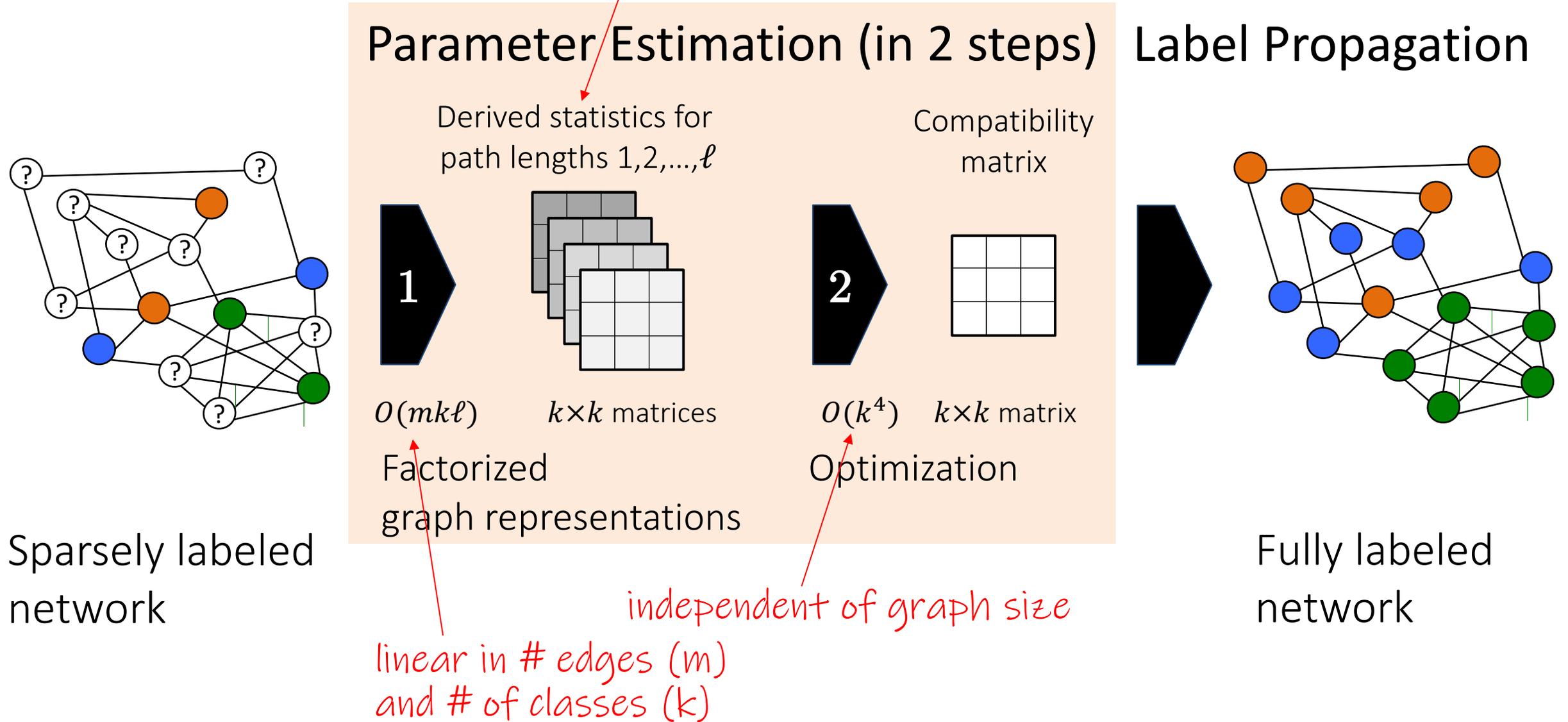
No more need for heuristics or domain experts 😊

# Approximate Agenda

- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation [VLDB'15]
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations [SIGMOD'20]
  - How well does it work?
  - **What is the magic sauce?**
  - What we would like to do (but it does not work)
  - What we actually do (Distant Compatibility Estimation)
- Discussion

# Splitting parameter estimation into two steps

*frequencies of labels across paths of different lengths*

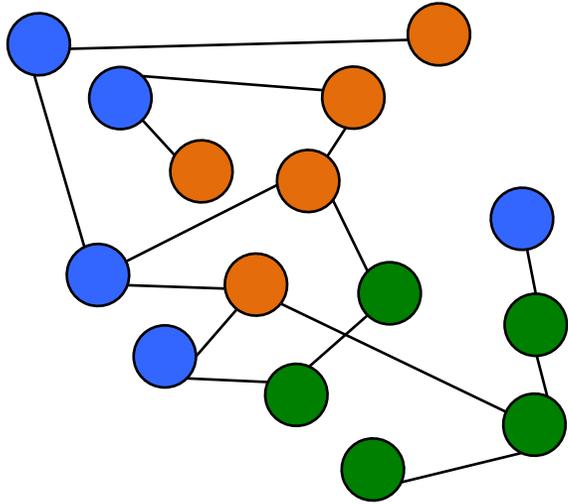


# Approximate Agenda

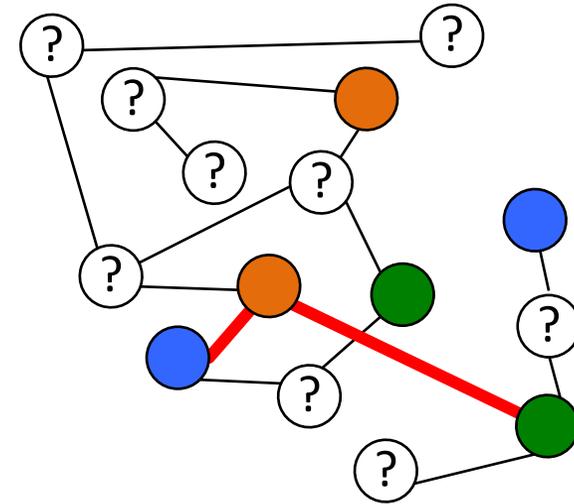
- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation [VLDB'15]
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations [SIGMOD'20]
  - How well does it work?
  - What is the magic sauce?
  - What we would like to do (but it does not work)
  - What we actually do (Distant Compatibility Estimation)
- Discussion

# A myopic view: counting relative neighbor frequencies

## Fully labeled graph



## Sparsely labeled graph



### Neighbor count

$$\mathbf{M} = \begin{array}{|c|c|c|} \hline \text{blue} & \text{orange} & \text{green} \\ \hline 2 & 6 & 2 \\ \hline 6 & 2 & 2 \\ \hline 2 & 2 & 6 \\ \hline \end{array}$$

### Gold standard compatibilities

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

normalize

$\Sigma=1$

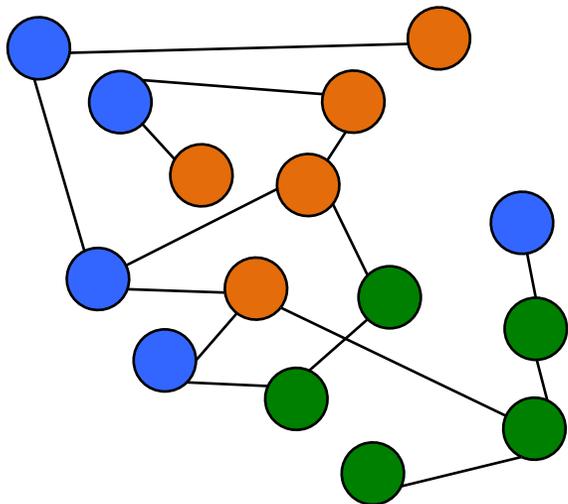
### Labeled neighbor count

$$\hat{\mathbf{M}} = \begin{array}{|c|c|c|} \hline \text{blue} & \text{orange} & \text{green} \\ \hline 0 & 1 & 0 \\ \hline 1 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline \end{array} \begin{array}{l} \Sigma=1 \\ \Sigma=2 \end{array} \Rightarrow \hat{\mathbf{H}}$$

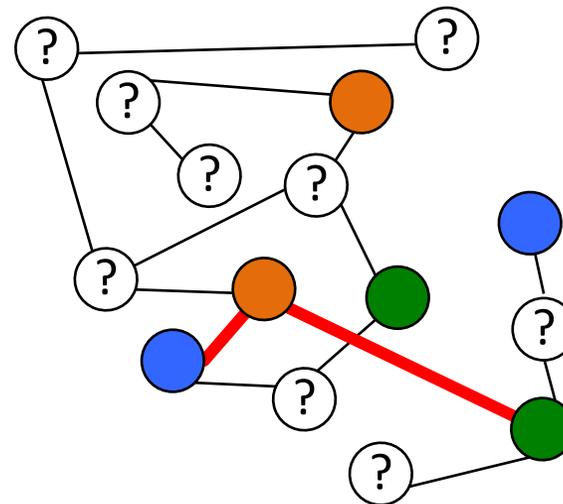
Idea: normalize, then find closest symmetric, doubly-stochastic matrix

# A myopic view: counting relative neighbor frequencies

## Fully labeled graph



## Sparsely labeled graph



### Neighbor count

### Gold standard compatibilities

$$\mathbf{M} = \begin{array}{|c|c|c|} \hline \text{blue} & \text{orange} & \text{green} \\ \hline \text{blue} & 2 & 6 & 2 \\ \hline \text{orange} & 6 & 2 & 2 \\ \hline \text{green} & 2 & 2 & 6 \\ \hline \end{array}$$

⇒

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

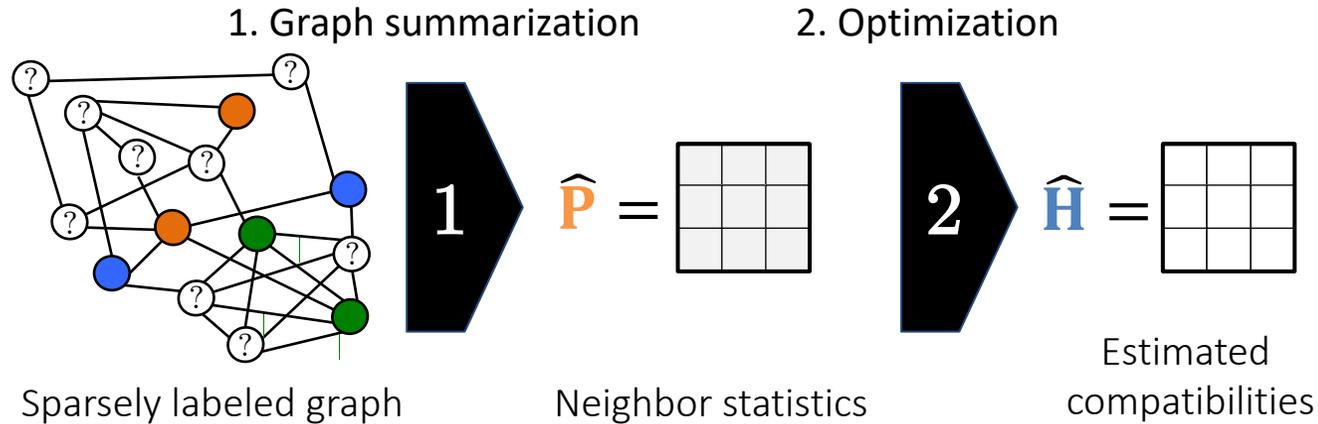
normalize

$\Sigma=1$

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

1% ☹️ Few nodes ⇒  
even fewer edges  $mf^2$

# Myopic compatibility estimation (MCE): from M to H DETAILS



Observed labeled neighbor counts  $\mathbf{M}$

$$\mathbf{M} = \mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{X}$$

*graph labels*

Observed row-stochastic compatibility matrix  $\hat{\mathbf{P}}$

$$\hat{\mathbf{P}} = \mathbf{M}^{\text{row}} \triangleq \text{diag}(\mathbf{M}\mathbf{1})^{-1} \cdot \mathbf{M}$$

Closest doubly stochastic symmetric matrix  $\hat{\mathbf{H}}$

$$E(\mathbf{H}) = \|\mathbf{H} - \hat{\mathbf{P}}\|^2$$

$$\hat{\mathbf{H}} = \min_{\mathbf{H}} E(\mathbf{H}) \quad \text{s.t.} \quad \begin{cases} \mathbf{H}\mathbf{1} = \mathbf{1} \\ \mathbf{H}^T = \mathbf{H} \end{cases}$$

*symmetric doubly stochastic constraints*

Example values

<del>18</del>	135	27	$\Sigma=180$
135	<del>36</del>	189	$\Sigma=360$
27	189	<del>324</del>	$\Sigma=540$

0.1	0.75	0.15	$\Sigma=1$
0.375	0.1	0.525	
0.05	0.35	0.6	

$\Sigma=0.525 \quad \Sigma=1.2 \quad \Sigma=1.275$

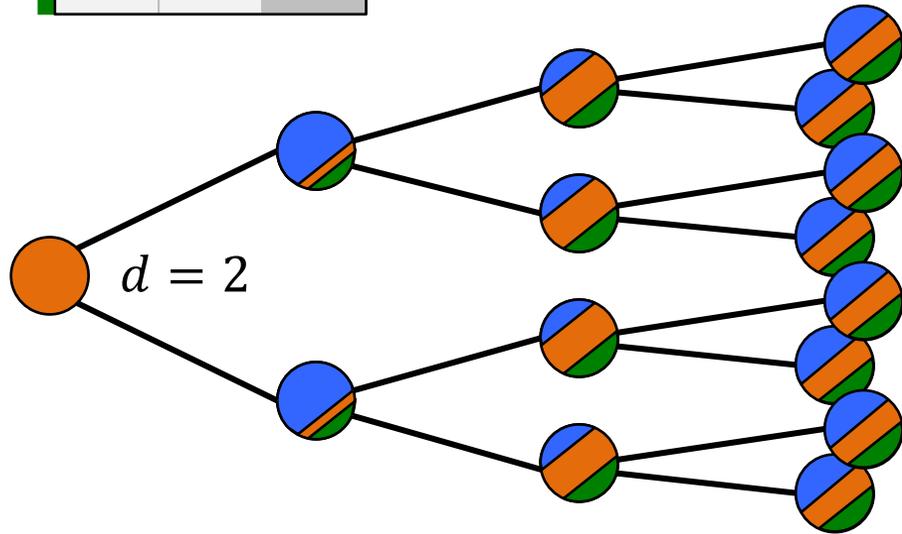
<del>0.258</del>	0.608	0.134	$\Sigma=1$
0.608	<del>0.034</del>	0.358	
0.134	0.358	<del>0.508</del>	

# Approximate Agenda

- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation [VLDB'15]
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations [SIGMOD'20]
  - How well does it work?
  - What is the magic sauce?
  - What we would like to do (but it does not work)
  - What we actually do (Distant Compatibility Estimation)
- Discussion

# Distant compatibility estimation (DCE)

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



$$\mathbf{H}^2 = \begin{array}{|c|c|c|} \hline 0.44 & 0.28 & 0.28 \\ \hline 0.28 & 0.44 & 0.28 \\ \hline 0.28 & 0.28 & 0.44 \\ \hline \end{array}$$

$$\mathbf{H}^3 = \begin{array}{|c|c|c|} \hline 0.31 & 0.38 & 0.31 \\ \hline 0.38 & 0.31 & 0.31 \\ \hline 0.31 & 0.31 & 0.38 \\ \hline \end{array}$$

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

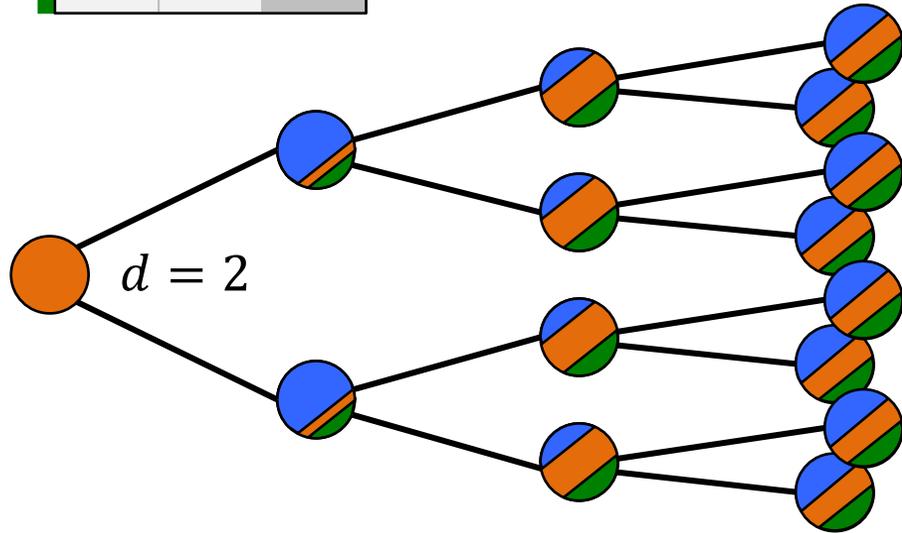
Expected signals for neighbors

0.6, 0.44, 0.38, 0.35, ...

(maximal entries)

# Distant compatibility estimation (DCE)

$$\mathbf{H} = \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.6 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}$$



	$\ell = 1$	$\ell = 2$	$\ell = 3$
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

graph with:

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

Expected # of labeled neighbors of distance  $\ell$  ?

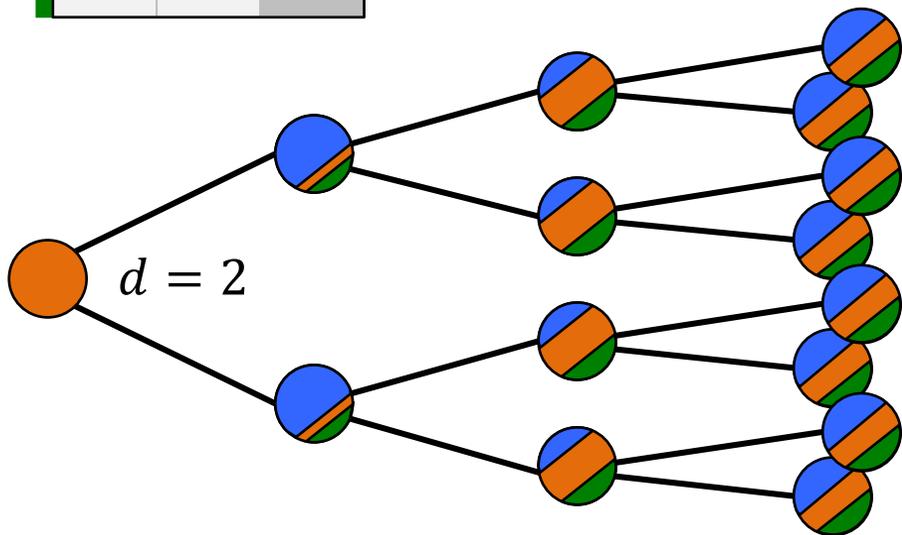
$d^{\ell-1} m f^2$  expected neighbors of distance  $\ell$

Idea: amplify the signal from observed length- $\ell$  paths 😊

# Distant compatibility estimation (DCE)

observed path- $\ell$ -row  
stochastic compatibilities

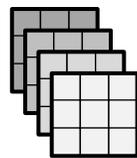
DETAILS

$$\mathbf{H} = \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.6 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}$$


	$\ell = 1$	$\ell = 2$	$\ell = 3$
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

distance-smoothed energy function



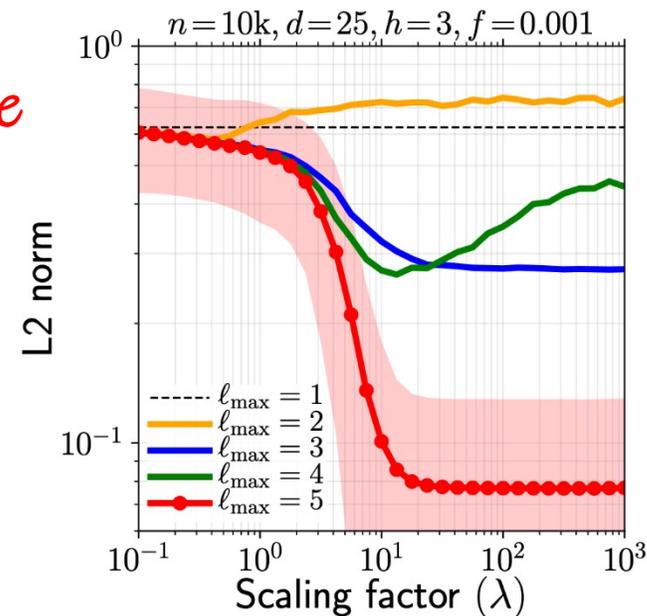
Statistics for path lengths 1, 2, ...

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

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

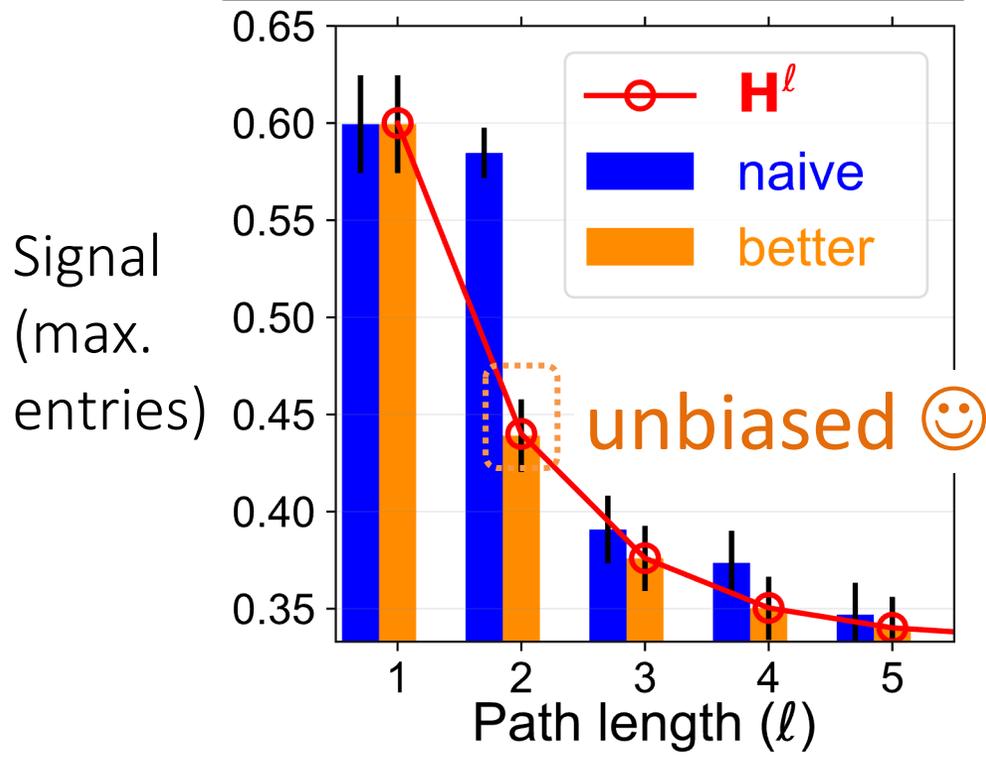
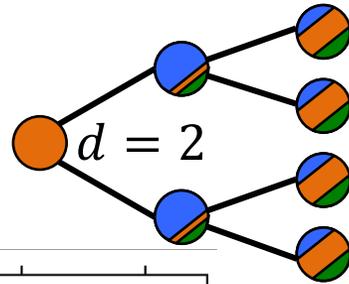
one free parameter (like in PageRank) 😊

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



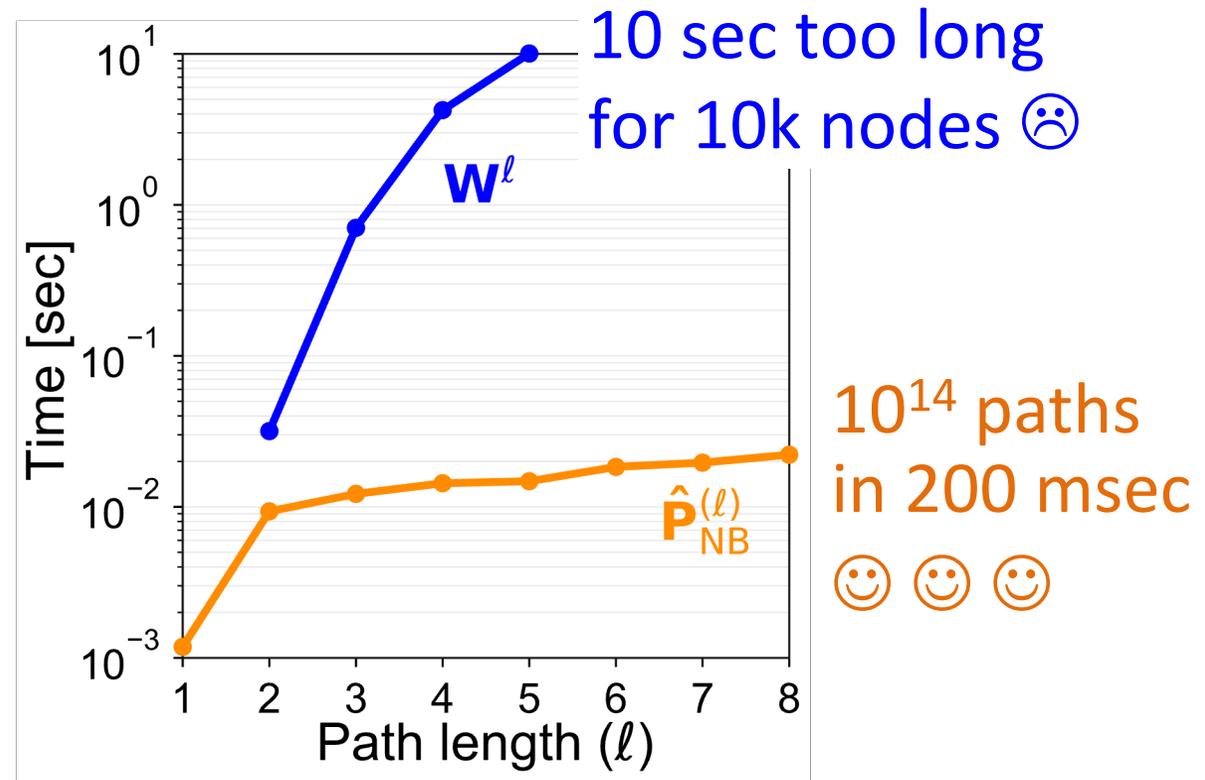
# Two technical difficulties

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

## Intuition

### Relational algebra

$$\pi_x(R(x) \bowtie S(x, y))$$

$$\Rightarrow R(x) \bowtie \pi_x S(x, y)$$

*W: sparse  $n \times n$  matrix*

*X: thin  $n \times k$  ( $k \ll n$ )  
label matrix*

### Linear algebra

$$(W \cdot W) \cdot X$$

$$\Rightarrow W \cdot (W \cdot X)$$

*small  $n \times k$  intermediate results*

## 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  $\ell$  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  $\ell \in [\ell_{\text{max}}]$  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  $\ell \in [\ell_{\text{max}}]$  in  $\mathcal{O}(mk\ell_{\text{max}})$ .

# Scalable factorized path summation

## Intuition

### Relational algebra

$$\pi_x(R(x) \bowtie S(x, y))$$

$$\Rightarrow R(x) \bowtie \pi_x S(x, y)$$

*W: sparse  $n \times n$  matrix*

### Linear algebra

$$(W \cdot W) \cdot X$$

*X: thin  $n \times k$  ( $k \ll n$ )  
label matrix*

$$\Rightarrow W \cdot (W \cdot X)$$

*small  $n \times k$  intermediate results*

Similar ideas of factorized calculation:

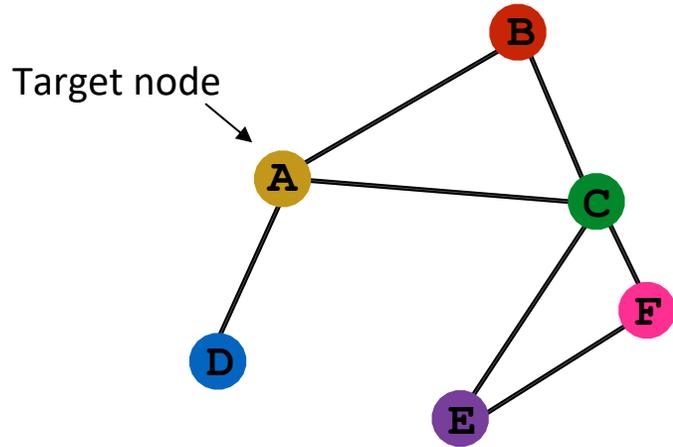
- Generalized distributive law  
[Aji-McEliece IEEE TIT '00]
- Algebraic path problems  
[Mohri JALC'02]
- Provenance semirings  
[Green+ PODS'07]
- Valuation algebras  
[Kohlas-Wilson AI'08]
- Factorized databases  
[Olteanu-Schleich Sigmod-Rec'16]
- FAQ (Functional Aggregate Queries)  
[AboKhamis-Ngo-Rudra PODS'16]
- Associative arrays  
[Kepner, Janathan MIT-press'18]
- Optimal ranked enumeration  
[Tziavelis+ VLDB'20]

# Approximate Agenda

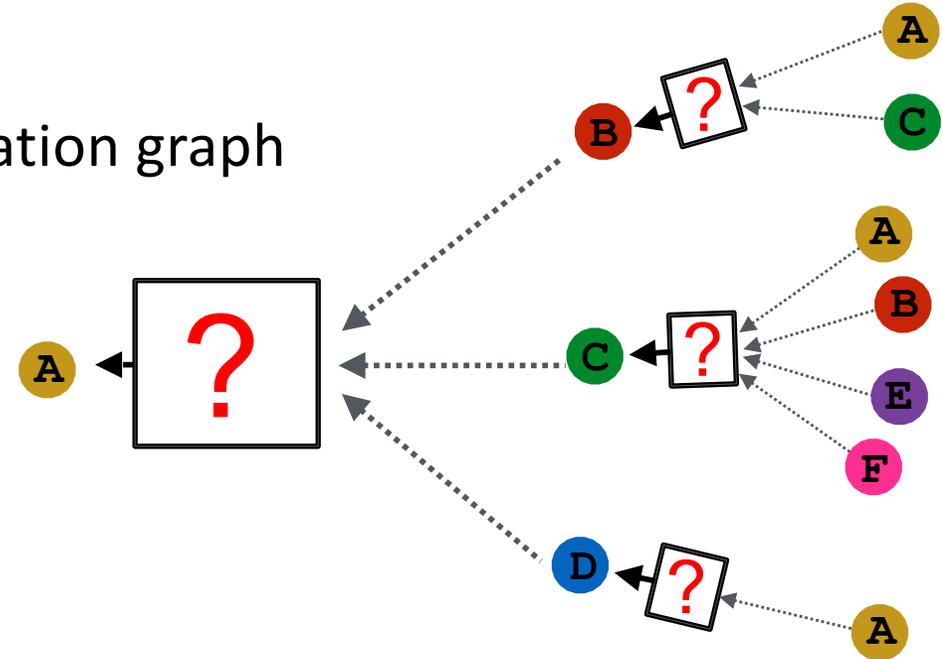
- Problem 1: How to propagate compatibilities?  
Linearized Belief Propagation
- Problem 2: How to learn/estimate compatibilities?  
Factorized graph representations
- Discussion

# Neighborhood aggregation in GNNs

Input graph



Computation graph



*simplification*

*simplification*

*simplification*

## GCNs

[Kipf, Welling 2017]  
Graph convolution,  
Supervised learning  
via cross entropy

## SGCs

[Wu+ 2019]  
Simplified graph convolution,  
Non-linearity only at the end

## LinBP / DCE

[VLDB'15], [SIGMOD'20]  
No non-linearities  $\Sigma$ ,  
Infinite layers, No oversmoothing  
because spectral radius  $< 1$ ,  
Structured regression via  $\ell_2$ -norm

## Belief Propagation

[Pearl 1986]  
Multiplication  $\Pi$

# Open topics

1. Network information only in semi-supervised setting: how much can one learn without node features?
  - a unified information theoretic framework (#parameters vs. #labeled data)
2. Phenomenology of network effects: label bias, degree distributions, long-distance interactions ("triangles"), combinatorial properties,...
  - how to create "unbiased" synthetic graph generators
3. What formalism can learn those phenomena "well enough"?
  - and how well "factorizable"

Thank you 😊