Iterative Methods for Network Alignment

David F. Gleich Computer Science Purdue University

with

Work supported by DOE CSCAPES Institute grant (DE-FC02-08ER25864), NSF CAREER grant 1149756-CCF, and the Center for Adaptive Super Computing Software Multithreaded Architectures (CASS-MT) at PNNL. Stanford's CADS grant from the Library of Congress. PNNL is operated by Battelle Memorial Institute under contract DE-AC06-76RL01830 Arif Khan, Alex Pothen Purdue University, Computer Science Mahantesh Halappanavar Pacific Northwest National Labs Mohsen Bayati, Amin Saberi Stanford University Ying Wang Google

What is the best way of matching graph *A* to *B*?







From Sharan and Ideker, Modeling cellular machinery through biological network comparison. *Nat. Biotechnol. 24,* 4 (Apr. 2006), 427–433.

Desserts

URI: <http://id.loc.gov/authorities/sh85037243#concept>

Type: Topical Term

Broader Terms:

• Confectionery

Narrower Terms:

- Ambient desserts
- Banana splits
- Charlottes (Desserts)
- Chocolate desserts
- Frozen desserts
- Ice cream cones
- Mousses
- Puddings
- <u>Refrigerated desserts</u>
- Sundaes
- Whipped toppings

LC Classification: TX773

Created: 1986-02-11

Last Modified: 1988-01-15 17:36:44



297,266 vertices, 248,230 edges 205,948 vertices, 382,353 edges 4,971,629 edges

Category: Desserts

From Wikipedia, the free encyclopedia

The main article for this category is **dessert**.

Desserts are sweet foods eaten purely for pleasure, typically at the end of a meal.

Subcategories

This category has the following 16 subcategories, out of 16 total.

	*	C cont.	Р
	[+] Desserts by country (20)	[+] Cookies (1)	[+] Pastry (4)
	-	[+] Custard desserts (0)	[+] Pies (9)
	В	_	[+] Puddings (2)
	[+] Brand name desserts (3)	D	-
	-	[+] Dessert sauces (0)	S
	C	[+] Doughnuts (1)	[+] Sweet breads (2)
	[+] Cakes (0)		
	[+] Chocolate desserts (0)	F	μ
	[+] Confectionery (9)	[+] Frozen desserts (2)	[+] Dessert stubs (1)
		I	
		[+] Ice cream (5)	
Γ.			



Wikimedia Commons has media related to: **Desserts**

Sometimes small data becomes big ...

Dataset	Size	Nonzeros
LCSH-2 WC-3	$59,849 \\ 70,509$	227,464 403,960
Product graph	4,219,893,141	91,886,357,440

... Ananth has some better techniques to work with these large problems ...

What is the best way of matching graph *A* to *B* using only edges in *L*?



Matching? 1-1 relationship *Best?* highest weight and overlap



- ... is NP-hard
- ... has no approximation algorithm



- Computer Vision
- Ontology matching
- Database matching
- Bioinformatics

objective = α matching + β overlap

via mathematical programming



Find a 1-1 matching between vertices with as many overlaps as possible.

via mathematical programming



Find a 1-1 matching between vertices with as **many overlaps** as possible.

Our contributions

A new belief propagation method (Bayati et al. 2009, 2013) Outperformed state-of-the-art PageRank and optimizationbased heuristic methods

High performance C++ implementations (Khan et al. 2012) 40 times faster (C++ ~ 3, complexity ~ 2, threading ~ 8) 5 million edge alignments ~ 10 sec

www.cs.purdue.edu/~dgleich/codes/netalignmc



Iterative methods

for network alignment

Each iteration involves

Matrix-vector-ish computations with a sparse matrix, e.g. sparse matrix vector products in a semiring, dot-products, axpy, etc.

Bipartite max-weight matching using a different weight vector at each iteration

No "convergence" 100-1000 iterations

Let x[i] be the score for each pair-wise match in L

```
for i=1 to ...
update x[i] to y[i]
compute a
   max-weight match
   with y
update y[i] to x[i]
   (using match in MR)
```

Open question 1

Any sort of property of these methods beyond ...

(i) Principled derivation and(ii) "David and Ananth say they work"?

Belief propagation methods

Summary

- Construct a probability model where the most likely state is the solution!
- Locally update information
- Like a generalized dynamic program

It works

 Most likely, it won't converge

History

- BP used for computing marginal probabilities and maximum aposterori probability
- Wildly successful at solving satisfiability problems
- Convergent algorithm for max-weight matching

Belief propagation for network alignment





variable i tells function j what it thinks about being in state s. This is just the product of what all the other functions tell i about being in state s.



function j tells variable i what it thinks about being in state s. This means that we have to locally maxamize f_j among all possible choices. Note $y_i = s$ always (too cumbersome to include in notation.)

Belief propagation for network alignment

For $t \ge 1$, the messages in iteration *t* are obtained from the messages in iteration t - 1 recursively. In particular for all $ii' \in E_L$

$$m_{ii' \to f_i}^{(t)} = \alpha W_{ii'} - \left(\max_{k \neq i} \left[m_{ki' \to g_{i'}}^{(t-1)} \right] \right)^+ + \sum_{ii' \neq i' \in V_S} \min\left(\frac{\beta}{2}, \max(0, \frac{\beta}{2} + m_{jj' \to h_{ii'jj'}}^{(t-1)}) \right).$$
(1)

The update rule for $m_{ii' \rightarrow g_{i'}}^{(t)}$ is similar, and

$$m_{ii' \to h_{ii'jj'}}^{(t)} = \alpha W_{ii'} - \left(\max_{k \neq i} \left[m_{ki' \to g_{i'}}^{(t-1)} \right] \right)^{+} - \left(\max_{k' \neq i'} \left[m_{ik' \to f_{i}}^{(t-1)} \right] \right)^{+} + \sum_{\substack{k \neq i' \neq j' \\ ii' kk' \in V_{S}}} \min \left(\frac{\beta}{2}, \max(0, m_{kk' \to h_{ii'kk'}}^{(t-1)} + \frac{\beta}{2}) \right).$$
(2)

Synthetic evaluation of network alignment



Open question 2

When could we hope to solve such synthetic problems in asymptotic regimes?

Does it work?

LCSH – Library of Congress subject headings Rameau – French National Library subject headings Manually matched

Obj.	Alg.	Weight	Overlap	Time (s)	Correct	Rec.	Prec.	Triangles
	Sol. MWM	$36332.42 \\93279.0$	$39847 \\ 16990$	 29.6	$57645 \\ 29098$	$100\% \\ 50.5\%$	$100\%\ 23.3\%$	$2073 \\ 350$
$\alpha = 1, \beta = 1$	BP BP++ MR	84622.0 85810.1 87588.6	$\begin{array}{c} 46400 \\ 46942 \\ 48367 \end{array}$	$23522.0 \\ 27115.6 \\ 33366.9$	$32585 \\ 32857 \\ 33225$	$56.5\%\ 57.0\%\ 57.6\%$	$27.6\%\ 27.4\%\ 27.0\%$	$1515 \\ 1548 \\ 1617$
$\alpha=1,\beta=2$	BP BP++ MR	81752.6 84615.7 85438.4	$\begin{array}{c} 46569 \\ 46656 \\ 48934 \end{array}$	$23427.1 \\ 26673.1 \\ 56961.6$	$31724 \\ 31952 \\ 32303$	$55.0\%\ 55.4\%\ 56.0\%$	$27.6\%\ 26.7\%\ 26.3\%$	$1483 \\ 1531 \\ 1604$
$\alpha=0,\beta=1$	BP BP++ MR	60617.9 60502.8 65994.2	$ 45247 \\ 41592 \\ 46163 $	$14284.8 \\ 13979.5 \\ 10384.4$	$24794 \\ 24498 \\ 25455$	$\begin{array}{c} 43.0\% \\ 42.5\% \\ 44.2\% \end{array}$	$23.2\% \\ 23.0\% \\ 21.5\%$	$ 1467 \\ 1484 \\ 1602 $

Open question 3

How can we evaluate alignments? What are possible null-models?

LCSH	WC
Science fiction television series	Science fiction television programs
Turing test	Turing test
Machine learning	Machine learning
Hot tubs	Hot dog



Higher-order network alignment





Network alignment via mathematical programming



Find a 1-1 matching between vertices with as many overlaps as possible.

Triangle alignment via mathematical programming



Find a 1-1 matching between vertices with as many overlaps **and triangles** as possible.

Tensor eigenvalues and a power method



Human protein interaction networks 48,228 triangles Yeast protein interaction networks 257,978 triangles The tensor T has ~100,000,000 nonzeros We work with it implicitly

Synthetic evaluation of network alignment



Open question 4

When do we need triangles?