Scalable Network Analysis

Inderjit S. Dhillon University of Texas at Austin COMPUTERSCIENCE

COMAD, Ahmedabad, India Dec 20, 2013

Outline

- Unstructured Data Scale & Diversity
 - Evolving Networks
- Machine Learning Problems arising in Networks
 - Recommender Systems
 - Link Prediction
 - Sign Prediction
- Formulation as Missing Value Estimation
- Scalable Algorithms
 - NOMAD: Distributed matrix completion algorithm
- Results on Applications
- Conclusions

Structured Data



- Data organized into fields: Relational databases, spreadsheets, XML
- Highly optimized for storage & retrieval (e.g. using SQL)

Structured Data



- Focus is on data format, efficient storage & search
- Less or no uncertainty in semantics: e.g. businesses know the fields of the data

Modern data is unstructured and diverse



Much greater growth rate

Total Archived Capacity, by Content Type, Worldwide, 2008-2015 (Petabytes)



Source: Enterprise Strategy Group, 2010.

- Dynamic aspects of unstructured data:
 - Constantly evolving
 - Uncertainties abound: What should I ask of the data?
 - Seek insights
- Heterogeneity renders traditional database models
 inadequate

- Buzzwords "Big Data" & "Data Science"
 - *Machine Learning*: Predictive models for data
 - Engineering perspective: Scale matters



Graph Evolution

- Social networks are highly dynamic
- Constantly grow, change quickly over time
 - Users arrive/leave, relationships form/dissolve
- Understanding graph evolution is important

Graphs meet Machine Learning

- Network analysis: Understanding structure & evolution of networks
- Formulate predictive problems on the adjacency matrix of the graph
- Confluence of graph theory & machine learning

Scalable Network Analysis



Recommender Systems



Link Prediction in Social Networks



• Problem: Infer *missing* relationships from a given snapshot of the network

Predicting gene-disease links





• Sign Prediction Problem: Given a snapshot of the signed social network, predict the signs of missing edges

Formulation as Missing Value Estimation









 $\min_{W \in \mathbb{R}^{m \times k}, H \in \mathbb{R}^{n \times k}} \sum_{(i,j) \in \Omega} (A_{ij} - W_i H_j^T)^2 + \lambda (\|W\|_F^2 + \|H\|_F^2)$



 $\min_{W \in \mathbb{R}^{m \times k}, H \in \mathbb{R}^{n \times k}} \sum_{(i,j) \in \Omega} (A_{ij} - W_i H_j^T)^2 + \lambda (\|W\|_F^2 + \|H\|_F^2)$

-0.46 -0.07 -0.35 -0.19 -0.07 -0.11 -0.53 -0.06 -0.05 -0.53 -0.14HT -0.17 0.13 -0.42 0.45 0.17 -0.25 -0.18 0.27 -0.59 0.05 0.14 -0.23 0.16 0.08 0.57 -0.39 -0.37 -0.08 -0.21 -0.43 0.17 -0.15 W 5 5 3 2 1 -8.72 0.03 -1.03 5 2 3 2 5 -7.56 -0.79 0.62 3.44 5 4.07 -3.95 2.55 3 3 5 3 2 4 2 -3.52 -3.32 3.73 5 5 1 -7.78 2.34 2.33 5 -2.44 -5.29 -3.92 1 5 1 2 -1.78 -1.68 1 4 1.90

 $\min_{W \in \mathbb{R}^{m \times k}, H \in \mathbb{R}^{n \times k}} \sum_{(i,j) \in \Omega} (A_{ij} - W_i H_j^T)^2 + \lambda (\|W\|_F^2 + \|H\|_F^2)$

Link Prediction

• Can be posed as matrix completion problem

$$A^{(t)} = \begin{bmatrix} . & 1 & 1 & . & . \\ 1 & . & 1 & 1 & 1 \\ . & 1 & . & . & . \\ . & 1 & 1 & . & . \end{bmatrix} \approx \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$
$$\frac{\text{Test Link}}{1 \text{ Score}}$$

• Issue: Only positive relationships are observed

Iest Link	score
(4,5)	1
(1,4)	1
(3,4)	1
(1,5)	1

Link Prediction

• Formulate Biased Matrix Completion Problem:

 $\min_{W \in \mathbb{R}^{n \times k}, H \in \mathbb{R}^{n \times k}} \sum_{(i,j) \in \Omega} (A_{ij} - W_i H_j^T)^2 + \alpha \sum_{(i,j) \notin \Omega} (A_{ij} - W_i H_j^T)^2 + \lambda (\|W\|_F^2 + \|H\|_F^2)$



- Social Balance [Harary, 1953]:
 - In real-world signed networks, triangles tend to be balanced



Theorem: All triangles in a network are balanced if and only if there exist two antagonistic groups.



- Relaxation: Weak balance
- Allow triangles with all negative edges



Theorem: All triangles in a network are weakly balanced if and only if there exist k antagonistic groups.

Sign Prediction

Theorem: A *k*-weakly balanced signed network has rank at most *k*.

 Sign inference can be posed as low-rank matrix completion

Theorem: If there are no "small" groups, the underlying network can be *exactly* recovered, under certain conditions.

K. Chiang et al. Prediction and Clustering in Signed Networks: A Local to Global Perspective. To appear in JMLR.

Scalable Algorithms

Stochastic Gradient

Sample random index (i,j) and update corresponding factors:

$$\mathbf{w}_{i} \leftarrow \mathbf{w}_{i} - \eta((A_{ij} - \mathbf{w}_{i}^{T}\mathbf{h}_{j})\mathbf{h}_{j} + \lambda\mathbf{w}_{i})$$
$$\mathbf{h}_{j} \leftarrow \mathbf{h}_{j} - \eta((A_{ij} - \mathbf{w}_{i}^{T}\mathbf{h}_{j})\mathbf{w}_{i} + \lambda\mathbf{h}_{j})$$

- Time per update O(k)
- Effective for very large-scale problems

Distributed Stochastic Gradient Descent (DSGD) [Gemulla *et al.* KDD 2011]

- Decoupled updates
- Easy to parallelize
- But communication & computation are interleaved

DSGD





Inderjit S. Dhillon University of Texas at Austin 34

Scalable Network Analysis

DSGD

Curse of the last reducer





DSGD



Inderjit S. Dhillon University of Texas at Austin 36

• **Goal**: Keep CPU & network simultaneously busy.

Non-locking stOchastic Multi-machine algorithm for Asynchronous & Decentralized matrix factorization

• Asynchronous distributed solution.



Inderjit S. Dhillon University of Texas at Austin 38 Scalable Network Analysis





Inderjit S. Dhillon University of Texas at Austin 40







Inderjit S. Dhillon University of Texas at Austin 42









Inderjit S. Dhillon University of Texas at Austin 44







Inderjit S. Dhillon University of Texas at Austin 46















NOMAD Algorithm

- **Initialize:** Randomly assign \mathbf{h}_i columns to worker queues 1.
- 2. Parallel Foreach q in $\{1, 2, \dots, p\}$
- 3. If queue[q] not empty then
- $(j, \mathbf{h}_j) \leftarrow \text{queue}[q].\text{pop}()$ for (i, j) in $\Omega_j^{(q)}$ do 4.
- 5.
- Do SGD úpdates 6.
- 7. end for
- 8. Sample q' uniformly from {1,2,...,p}
- queue[q'].push((j, \mathbf{h}_j)) 9.
- 10. end if

```
11. Parallel End
```

NOMAD Algorithm

- **Initialize**: Randomly assign h_i columns to worker queues 1.
- **Parallel Foreach** q in $\{1, 2, ..., p\}$ 2.
- 3. **If** queue[*q*] not empty **then**
- $(j, \mathbf{h}_j) \leftarrow \text{queue}[q].pop()$ 4.
- for (i, j) in $\Omega_i^{(q)}$ do 5.

10.

- Do SGD 6. Distributed setting: Write over the network end for
- 7. Sample q' unif rmly from {1,2,...,p} 8. 9.



Algorithm Complexity

• Average space required per worker:

 $O(mk/p + nk/p + |\Omega|/p)$

• Average time for one sweep:

 $O(|\Omega|k/p)$

Results on Applications

Recommender Systems

Netflix dataset: 2,649,429 users, 17,770 movies, ~100M ratings



Recommender Systems

Synthetic dataset: ~85M users, 17,770 items, ~8.5B observations



machines=32, cores=4, k = 100, $\lambda = 0.01$

Link Prediction

- Flickr dataset: 1.9M users & 42M links.
- Test set: sampled 5K users.



*using network-based features.

Predicting Gene-Disease Links



Sign Prediction

- Epinions dataset (+ve & -ve reviews):
 - 131K nodes, 840K edges, 15% edges negative.
 - MF-ALS is faster and achieves higher accuracy.



Conclusions

- Rapid growth of unstructured data demands scalable machine learning solutions for analysis
- Machine learning problems arising in network analysis can be cast in the matrix completion framework
- Our proposed asynchronous distributed algorithm NOMAD outperforms state-of-the-art matrix completion solvers
- Beyond Matrix Completion: Asynchronous distributed framework for solving machine learning problems