

Covariance Matrix Adaptation Evolution Strategy for Link Prediction in Dynamic Social Networks

Catherine A. Bliss, Morgan R. Frank, Christopher M. Danforth, & Peter Sheridan Dodds

Department of Mathematics & Statistics
Vermont Complex Systems Center
Computational Story Lab
Vermont Advanced Computing Core
University of Vermont



- Background
 - Data
 - Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions



- Background
 - Data
 - Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions



Background

Data

Reciprocal reply networks

Link prediction

Similarity indices

Evolutionary computation

Results

Conclusions



Background

Data

Reciprocal reply networks

Link prediction

Similarity indices

Evolutionary computation

Results

Conclusions



- Background
- Data
 - Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions



- Background
- Data
 - Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions

Background

Data

Reciprocal reply networks

Link prediction

Similarity indices

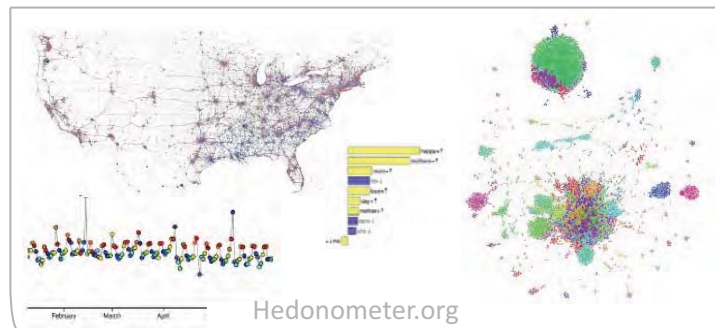
Evolutionary computation

Results

Conclusions



40,000 tweets (100MB) / min.
50 million tweets (150GB) / day
50 billion tweets (100TB) / 5+ years



Hedonometer.org



40,000 tweets (100MB) / min.
50 million tweets (150GB) / day
50 billion tweets (100TB) / 5+ years



Background

Data

Reciprocal reply networks

Link prediction

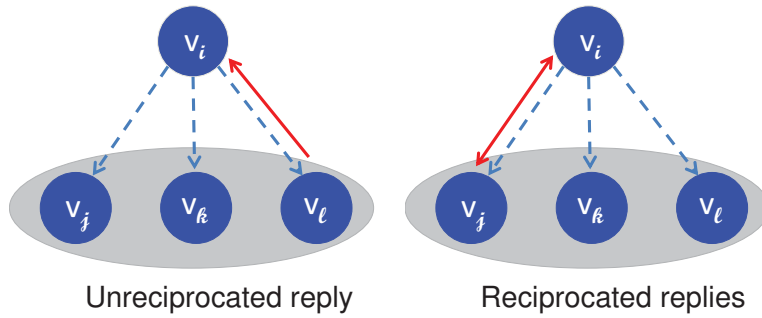
Similarity indices

Evolutionary computation




Results

Conclusions

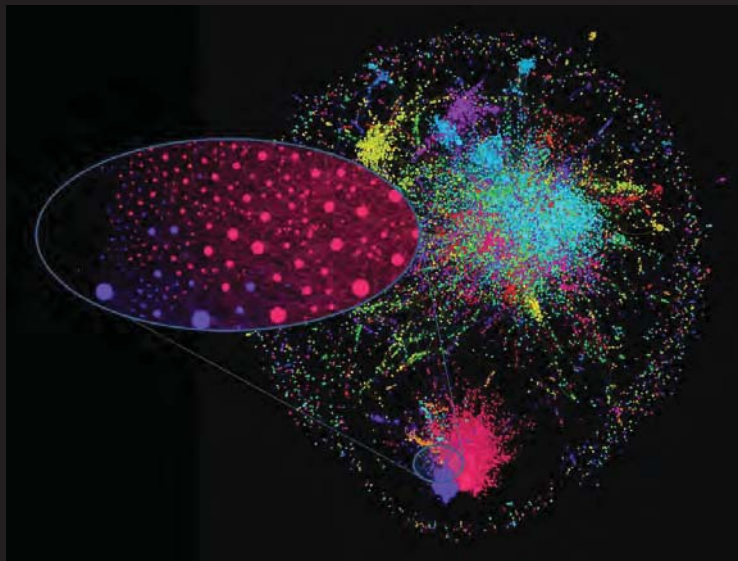
Reciprocal-reply networks



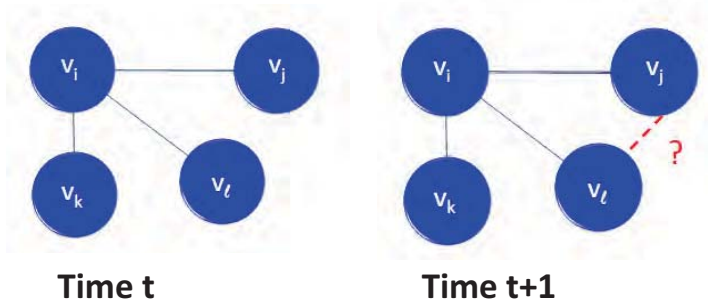
We define a *reciprocal-reply network* as an unweighted, undirected network in which a link is established between nodes v_i and v_j if we observe reciprocal replies between these nodes **during the unit of time under analysis**.¹

¹ C. A. Bliss, I. M. Kloumann, K. D. Harris, C. M. Danforth & P. S. Dodds. 2012. Twitter reciprocal reply networks exhibit assortativity with respect to happiness, *Journal of Computational Science*   

Twitter RRN Link Prediction



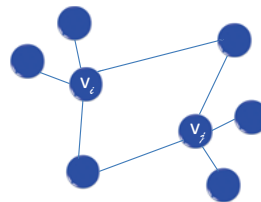
The link prediction problem: *Given a snapshot of the network at time= t , can we predict links which will appear in time= $t+1$?*



- Background
- Data
- Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions

Similarity indices

- ▶ Liben-Nowell & Kleinberg (2007) author collaboration networks ($N \propto 10^3$)
- ▶ Use similarity indices to rank the most likely occurring top N links



Common neighbors
(Newman, 2001)



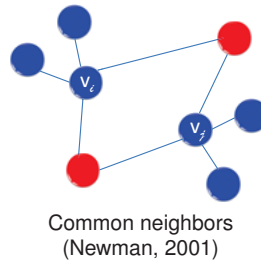
- Background
- Data
- Reciprocal reply networks
- Link prediction
 - Similarity indices
 - Evolutionary computation
- Results
- Conclusions

Similarity indices

- ▶ Liben-Nowell & Kleinberg (2007) author collaboration networks ($N \propto 10^3$)

- ▶ Use similarity indices to rank the most likely occurring top N links

$$C(u, v) = |\Gamma(u) \cap \Gamma(v)|$$



Background

Data
Reciprocal reply networks

Link prediction

Similarity indices
Evolutionary computation

Results

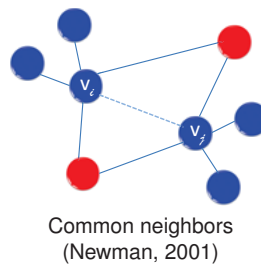
Conclusions

Similarity indices

- ▶ Liben-Nowell & Kleinberg (2007) author collaboration networks ($N \propto 10^3$)

- ▶ Use similarity indices to rank the most likely occurring top N links

$$C(u, v) = |\Gamma(u) \cap \Gamma(v)|$$



Background

Data
Reciprocal reply networks

Link prediction

Similarity indices
Evolutionary computation

Results

Conclusions

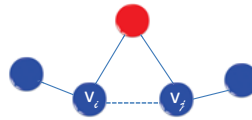
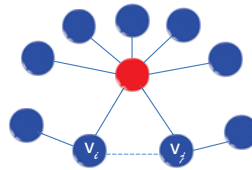
Similarity indices

- ▶ Liben-Nowell & Kleinberg (2007) author collaboration networks ($N \propto 10^3$)

- ▶ Use similarity indices to rank the most likely occurring top N links

$$C(u, v) = |\Gamma(u) \cap \Gamma(v)|$$

$$R(u, v) = \sum_{z \in \Gamma(u) \cap \Gamma(v)} \frac{1}{|\Gamma(z)|}$$



Resource Allocation
(Zhou, Lu, & Zhang, 2009)



Background

Data
Reciprocal reply networks

Link prediction

Similarity indices
Evolutionary computation

Results

Conclusions

Similarity indices

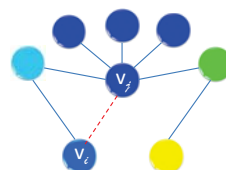
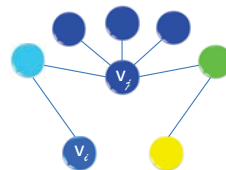
- ▶ Liben-Nowell & Kleinberg (2007) author collaboration networks ($N \propto 10^3$)

- ▶ Use similarity indices to rank the most likely occurring top N links

$$C(u, v) = |\Gamma(u) \cap \Gamma(v)|$$

$$R(u, v) = \sum_{z \in \Gamma(u) \cap \Gamma(v)} \frac{1}{|\Gamma(z)|}$$

$$W(u, v) = 1 - \frac{1}{2} \sum |f_{u,n} - f_{v,n}|$$



Word similarity
(Bliss et al., 2012)



Background

Data
Reciprocal reply networks

Link prediction

Similarity indices
Evolutionary computation

Results

Conclusions

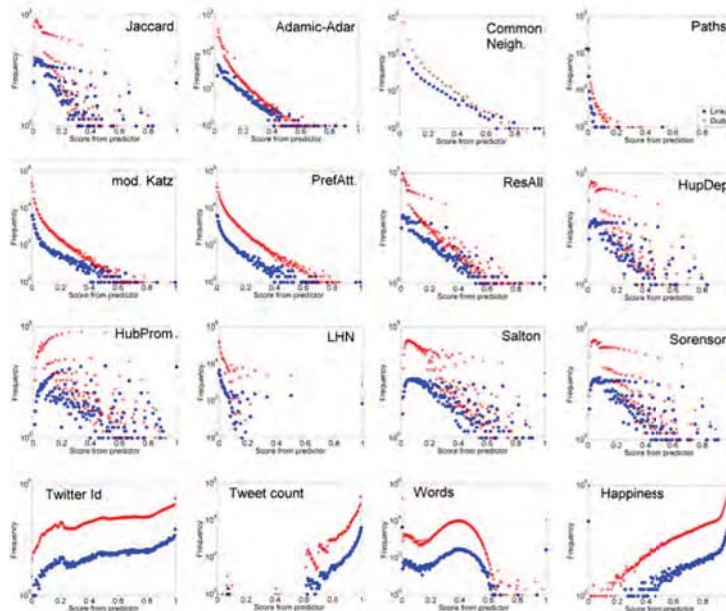
Similarity indices

| | |
|------------------------|---|
| Common neighbors | $C(u, v) = \Gamma(u) \cap \Gamma(v) $ |
| Jacard | $J(u, v) = \frac{ \Gamma(u) \cap \Gamma(v) }{ \Gamma(u) \cup \Gamma(v) }$ |
| Adamic-Adar | $A(u, v) = \sum_{z \in \Gamma(u) \cap \Gamma(v)} \frac{1}{\log(\Gamma(z))}$ |
| Pref Attachment | $Pr(u, v) = k_u \times k_v$ |
| Hub promoted | $Hp(x, y) = \frac{ \Gamma(u) \cap \Gamma(v) }{\min\{k_u, k_v\}}$ |
| Hub depressed | $Hd(u, v) = \frac{ \Gamma(u) \cap \Gamma(v) }{\max\{k_u, k_v\}}$ |
| LHN | $L(u, v) = \frac{ \Gamma(u) \cap \Gamma(v) }{k_u k_v}$ |
| Salton | $Sa(u, v) = \frac{ \Gamma(u) \cap \Gamma(v) }{\sqrt{k_u k_v}}$ |
| Sorenson | $So(u, v) = \frac{2 \Gamma(u) \cap \Gamma(v) }{k_u + k_v}$ |
| Resource Allocation | $R(u, v) = \sum_{z \in \Gamma(u) \cap \Gamma(v)} \frac{1}{ \Gamma(z) }$ |
| Average Path Weight | $P(u, v) = \frac{\sum_{p \in \mathcal{P}_2(u, v) \cup \mathcal{P}_3(u, v)} w_p}{ \mathcal{P}_2(u, v) + \mathcal{P}_3(u, v) }$ |
| Katz | $K = \sum_{n=1}^{\infty} \beta^n A^n$ |
| Tweet count similarity | $T(u, v) = 1 - \frac{ \Gamma(u) - \Gamma(v) }{\max\{ \Gamma(a) - \Gamma(b) \}}_{a, b \in V}$ |
| Word similarity | $W(u, v) = 1 - \frac{1}{2} \sum_{n=1}^{50000} f_{u,n} - f_{v,n} $ |
| Happiness similarity | $H(u, v) = 1 - \frac{ h(u) - h(v) }{\max\{ h(a) - h(b) \}}_{a, b \in V}$ |
| Id similarity | $I(u, v) = 1 - \frac{ d(u) - d(v) }{\max\{ d(a) - d(b) \}}_{a, b \in V}$ |

Navigation icons: back, forward, search, etc.

Background
 Data
 Reciprocal reply networks
 Link prediction
 Similarity indices
 Evolutionary computation
 Results
 Conclusions

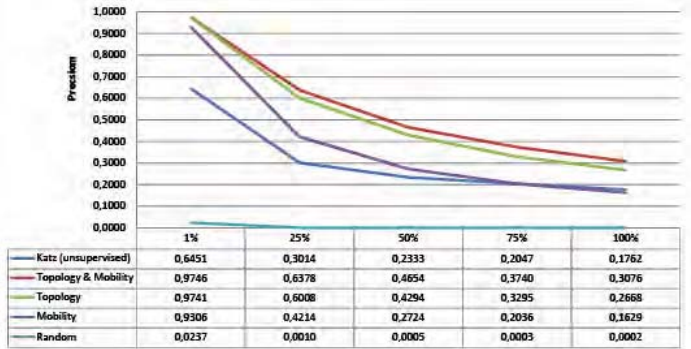
Signal detection



Navigation icons: back, forward, search, etc.

Background
 Data
 Reciprocal reply networks
 Link prediction
 Similarity indices
 Evolutionary computation
 Results
 Conclusions

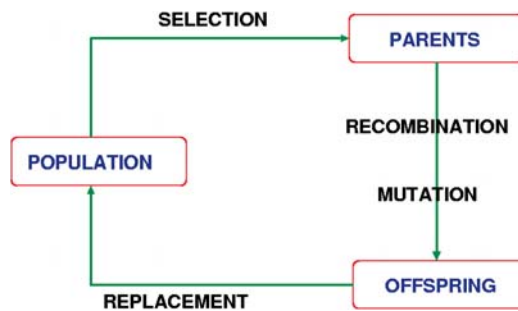
Combining indices



Wang et al., (2011) examine mobile call graphs $N \propto 10^4$

- ▶ Combination of topological and node based similarity indices outperform single indices
- ▶ Use supervised learning
- ▶ Examine a small subset of node-node pairs, $\propto 10^3$
- ▶ Challenge - methods for large, sparse networks

Evolutionary computation



CMA-ES implementation²

Individual

An individual or candidate solution is a vector, $\vec{w} \in \mathbf{R}^{16}$.

| J | A | C | P | K | H | W | Pr | R | Hd | Hp | L | Sa | So | I | T |
|----|----|----|----|----|----|----|----|----|-----|-----|----|----|-----|----|----|
| .4 | .9 | .5 | .8 | .3 | .2 | .1 | .6 | .7 | -.1 | .01 | .6 | .2 | .04 | .8 | .1 |

²Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larrañaga, I. Inza and E. Bengoetxea (eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102, Springer.



CMA-ES implementation²

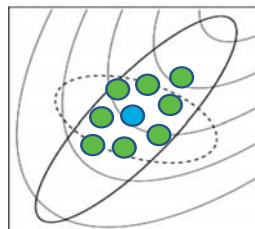
Individual

An individual or candidate solution is a vector, $\vec{w} \in \mathbf{R}^{16}$.

| J | A | C | P | K | H | W | Pr | R | Hd | Hp | L | Sa | So | I | T |
|----|----|----|----|----|----|----|----|----|-----|-----|----|----|-----|----|----|
| .4 | .9 | .5 | .8 | .3 | .2 | .1 | .6 | .7 | -.1 | .01 | .6 | .2 | .04 | .8 | .1 |

Reproduction & Mutation

CMA-ES
From 1 individual, generate a Gaussian cloud of candidate solutions in \mathbf{R}^{16} using the covariance matrix.



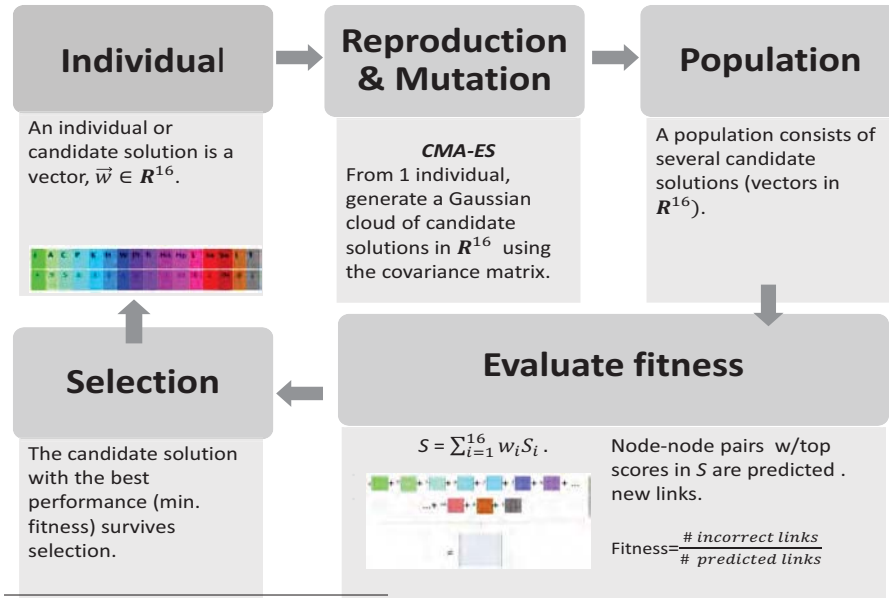
Population

A population consists of several candidate solutions (vectors in \mathbf{R}^{16}).

²Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larrañaga, I. Inza and E. Bengoetxea (eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102, Springer.

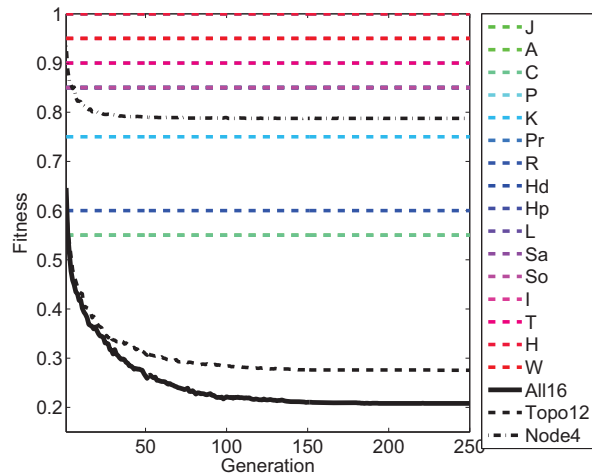


CMA-ES implementation²



²Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larrañaga, I. Inza and E. Bengoetxea (eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102, Springer.

Fitness



Background

Data
Reciprocal reply networks

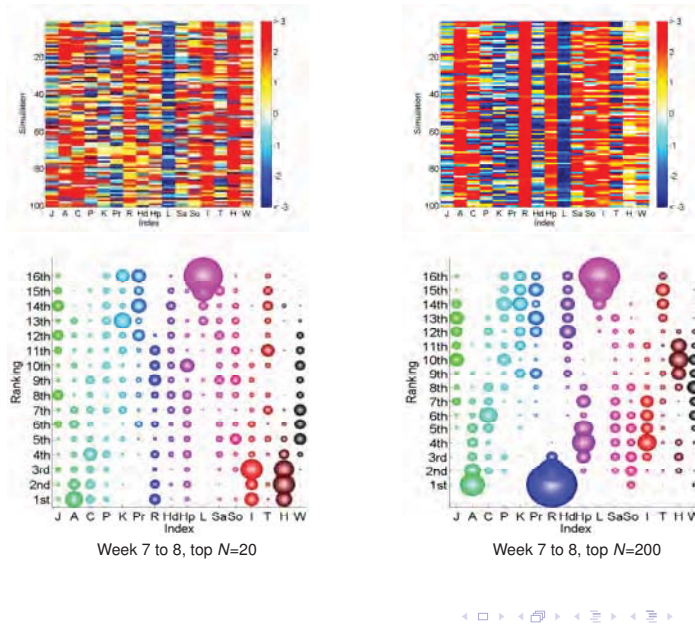
Link prediction

Similarity indices
Evolutionary computation

Results

Conclusions

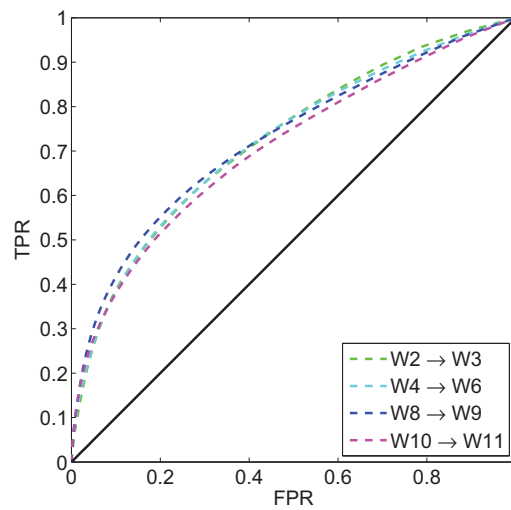
Best predictors



Navigation icons: back, forward, search, etc.

Background
 Data
 Reciprocal reply networks
 Link prediction
 Similarity indices
 Evolutionary computation
 Results
 Conclusions

Receiver Operating Curve

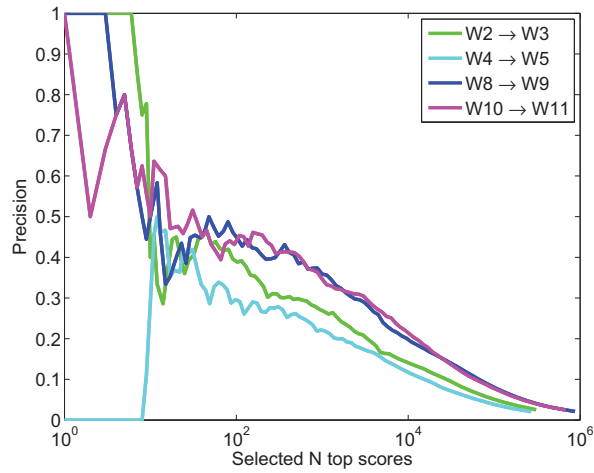


ROC depicts $TPR = \frac{TP}{TP+FN}$ as a function of $FPR = \frac{FP}{FP+TN}$. The area under the curve (AUC) represents the chance that our predictor assigns a higher score to user-user pairs who exhibit new links than user-user pairs who do not exhibit new links. AUC>0.70

Navigation icons: back, forward, search, etc.

Background
 Data
 Reciprocal reply networks
 Link prediction
 Similarity indices
 Evolutionary computation
 Results
 Conclusions

Precision



Precision depicts $\frac{TP}{TP+FP}$. High precision is achieved for $topN < 20$, which is often the region of interest. The precision for predicted links in $W4 \rightarrow W5$'s is lower than the other weeks and this may be due to missing data for those weeks



Background

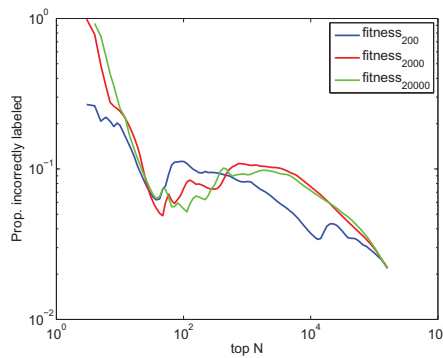
Data
Reciprocal reply networks
Link prediction
Similarity indices
Evolutionary computation

Results

Conclusions

Missing data

- ▶ Sample 50% of tweets that we have
- ▶ Build nets & re-run CMA-ES on this smaller sample
- ▶ Compare # of links that are labeled false-positives to our fuller set of tweets
- ▶ Count # of FPs that are TPs (had we had seen the fuller set of data)



Background

Data
Reciprocal reply networks
Link prediction
Similarity indices
Evolutionary computation

Results

Conclusions

Conclusions

- ▶ Evolutionary algorithms show promise
- ▶ Many additional questions in link prediction (e.g., prediction of weights, prediction of link decay)
- ▶ Leveraging link prediction to understand network dynamics
- ▶ Further investigation of the role of incomplete data on network inference



Background
Data
Reciprocal reply networks
Link prediction
Similarity indices
Evolutionary computation
Results
Conclusions

Thank you

- ▶ **Manuscript:** In press at the *Journal of Computational Science*. Pre-print available at <http://arxiv.org/abs/1304.6257>
- ▶ **Contact:** www.cems.uvm.edu/~cabliss
- ▶ **Lab:** www.onehappybird.com



Background
Data
Reciprocal reply networks
Link prediction
Similarity indices
Evolutionary computation
Results
Conclusions