# Statistical inference for the discovery of hidden interactions in complex networks
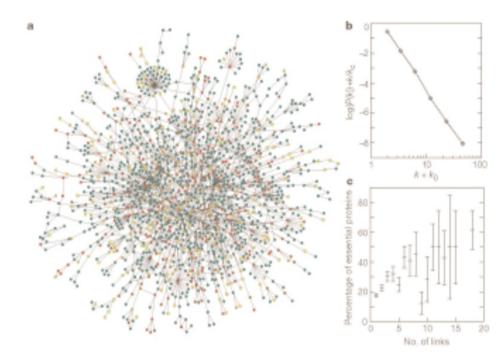
**Roger Guimerà**

ICREA and
Chemical Engineering, Universitat Rovira i Virgili

NetSci'13
Copenhagen, June 4, 2013

# One billion dollars to map the human proteome



Jeong, et al., *Nature* (2001)

# Accuracy and coverage are a concern for protein interaction (and most other) datasets



von Mering et al., *Nature* (2002)

# All network data is subject to noise

# Network properties are often sensitive to even low error rates

# Network properties are often sensitive to even low error rates

For the most part, we ignore(d) the issue of network data reliability and pretend(ed) that there is no problem

# Outline



Can we help to clean up noisy network data?

Can we uncover unknown drug interactions?

Can we predict human decisions?

Can we predict conflict in small teams?

**NETWORK INFERENCE**

Statistical physics

Statistical & machine learning

Network theory

Team Work

# What is to be done?

➜ Given a single noisy observation of a network, determine:

  ➜ *Missing interactions*   Interactions that exist but are not captured in our observation of the system

  ➜ *Spurious interactions*   Interactions that do *not* exist but, for some reason, are included in our observation

➜ *Reconstruct the network*, so that our reconstruction has properties that are closer to the properties of the true network

# What is to be done?

➔ Given a single noisy observation of a network, determine:

  ➔ *Missing interactions*   Interactions that exist but are not captured in our observation of the system

  ➔ *Spurious interactions*   Interactions that do *not* exist but, for some reason, are included in our observation

➔ *Reconstruct the network*, so that our reconstruction has properties that are closer to the properties of the true network

➔ But:

  ➔ We want to be able to do this for arbitrary real networks about which we don't know anything

  ➔ There seems to be a paradox in trying to identify what is wrong in a network observation–from *the network observation itself* !

# There are two possible scenarios when in comes to solving the paradox

➔ Scenario 1: We *don't* have a clue about what the network should look like, or where does it come from (mechanistically or statistically):

  ➔ We cannot do anything

➔ Scenario 2: We *do* have some ideas about the structure of the network:

  ➔ We can formalize these ideas into a set of models

  ➔ We can use the models to assess what is likely to be missing/wrong

# The "reliability formalism"

➜ We assume our network is the outcome of an undetermined model $M$ from a (potentially infinite) collection of models $\mathcal{M}$

➜ We observe a network $A^O$

➜ Given my observation $A^O$, what is the probability that a property $X$ takes the value $X=x$ if we generate a new network (with the same model)?

$$
\begin{aligned}
p(X = x | A^O) &= \int_{\mathcal{M}} dM\, p(X = x | M)\, p(M | A^O) \\
&= \frac{\int_{\mathcal{M}} dM\, p(X = x | M)\, p(A^O | M)\, p(M)}{\int_{\mathcal{M}} dM\, p(A^O | M)\, p(M)}
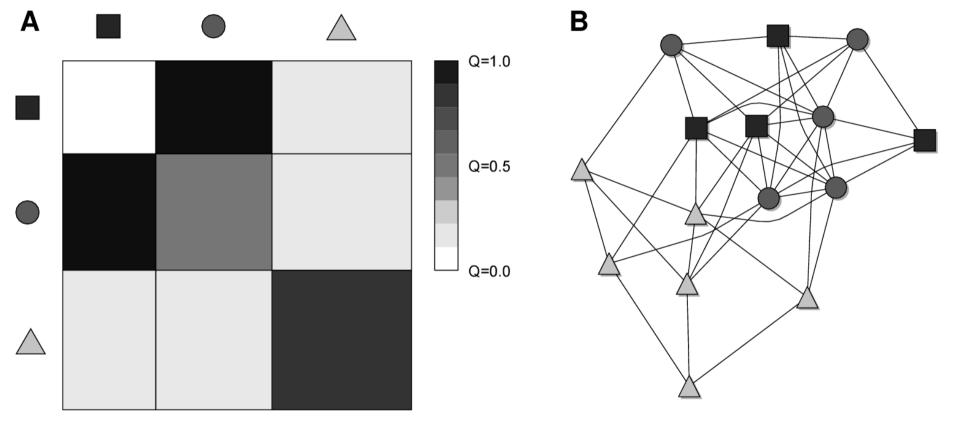\end{aligned}
$$

➜ We call $p(X=x|A^O)$ the reliability of the $X=x$ measurement

# In particular, one can use the formalism to infer missing and spurious interactions

$$p(A_{ij} = 1|A^O) = \frac{\int_{\mathcal{M}} dM \, p(A_{ij} = 1|M) \, p(A^O|M) \, p(M)}{\int_{\mathcal{M}} dM \, p(A^O|M) \, p(M)}$$

➔ What property of networks is general enough that applies to *all* complex networks?

  ➔ Broad (scale-free) connectivity distribution? No

  ➔ Small world property? Yes—but no realistic/tractable model

  ➔ Modularity? Group structure? YES

Clauset, Moore, Newman, *Nature* (2008)

Guimera, Sales-Pardo, *PNAS* (2009)

# Stochastic block models (SBM) are *general,* *empirically grounded* and analytically *tractable*



➔ A stochastic block model is fully determined by a partition of the nodes into groups and the probabilities $Q$ that a node in a group is connected to a node in any other group

White, Boorman, Breiger, *AJS* (1976)
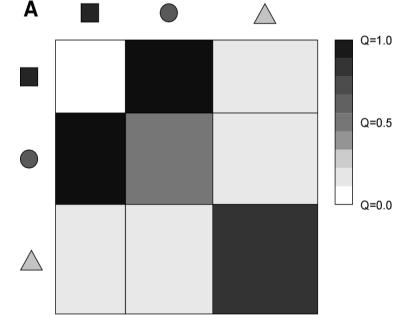
Holland, Laskey, Leinhardt, *Soc. Networks* (1983)

Nowicki, Snijders, *JASA* (2001)

# Stochastic block models (SBM) are *general*, *empirically grounded* and analytically *tractable*

$$p(A_{ij} = 1 | A^O) = \frac{\int_{\mathcal{M}} dM \, p(A_{ij} = 1 | M) \, p(A^O | M) \, p(M)}{\int_{\mathcal{M}} dM \, p(A^O | M) \, p(M)}$$

$$p(A_{ij} = 1 | M) = Q_{\sigma_i \sigma_j}$$

$$p(A^O | M) = \prod_{\alpha \leq \beta} Q_{\alpha\beta}^{n_{\alpha\beta}^1} (1 - Q_{\alpha\beta})^{n_{\alpha\beta}^0}$$

$$p(M) = \text{constant}$$

$$\int_{\mathcal{M}} dM \rightarrow \sum_{P \in \mathcal{P}} \prod_{\alpha \leq \beta} \left( \int_0^1 dQ_{\alpha\beta} \right)$$

# The link reliability is an ensemble average over all possible partitions of the nodes into groups
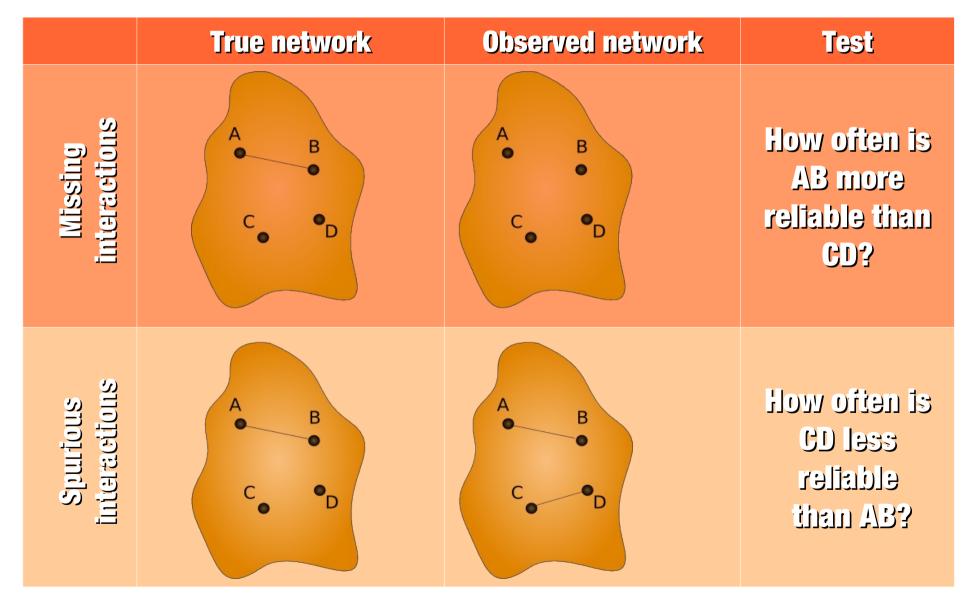
➜ In the end, the reliability of a link is

$$p(A_{ij} = 1 | A^O) = \frac{1}{Z} \sum_{P \in \mathcal{P}} \left( \frac{n^1_{\sigma_i \sigma_j} + 1}{n^0_{\sigma_i \sigma_j} + n^1_{\sigma_i \sigma_j} + 2} \right) \exp[-\mathcal{H}(\mathcal{P})]$$
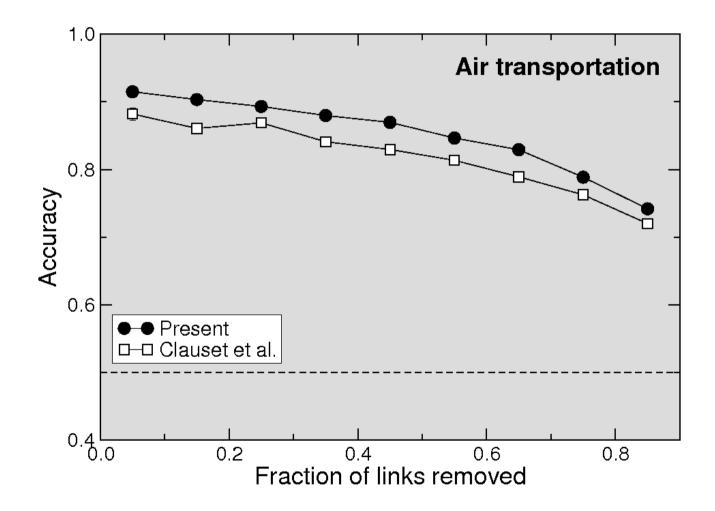
➜ Where:

$$\mathcal{H}(\mathcal{P}) = \sum_{\alpha \leq \beta} \left[ \ln(n_{\alpha\beta} + 1)! - \ln(n^0_{\alpha\beta})! - \ln(n^1_{\alpha\beta})! \right]$$

# We test our algorithm to see if it can identify missing and spurious interactions in real networks



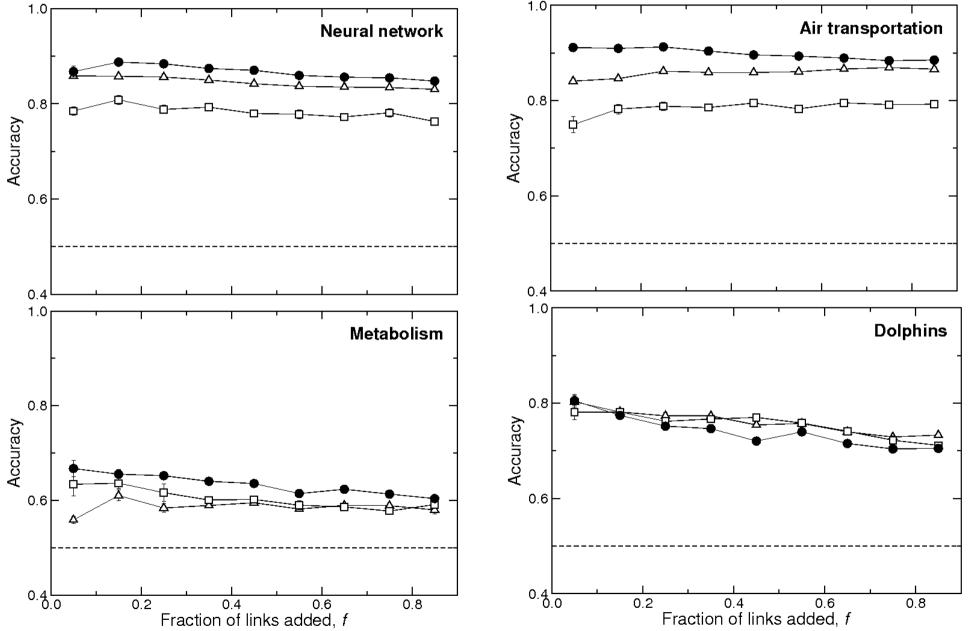| | True network | Observed network | Test |
|---|---|---|---|
| **Missing interactions** | | | **How often is AB more reliable than CD?** |
| **Spurious interactions** | | | **How often is CD less reliable than AB?** |

# Our approach accurately recovers missing interactions

# Our approach accurately recovers missing interactions

# Our approach accurately recovers spurious interactions

# Wonkish interlude I: H, module identification, maximum likelihood block models and all that

$$p(A_{ij} = 1|A^O) = \frac{1}{Z} \sum_{P \in \mathcal{P}} \left( \frac{n^1_{\sigma_i \sigma_j} + 1}{n^0_{\sigma_i \sigma_j} + n^1_{\sigma_i \sigma_j} + 2} \right) \exp[-\mathcal{H}(\mathcal{P})]$$

➜ What is this "energy"?

$$\mathcal{H}(P) = -\ln p(P|A^O)$$

➜ Therefore, the partition that minimizes **this** energy is the most likely given the data (except for priors, degree correction of the block model...):

  ➜ More appropriate "modularity" function

  ➜ No need to play with the number of groups

  ➜ No over-fitting

# Wonkish interlude II

Unipartatite unweighted: $\mathcal{H}(\mathcal{P}) = \sum_{\alpha \le \beta} \left[ \ln(n_{\alpha\beta} + 1)! - \ln(n_{\alpha\beta}^0)! - \ln(n_{\alpha\beta}^1)! \right]$

Unipartite weighted: $\mathcal{H}(\mathcal{P}) = \sum_{\alpha \le \beta} \left[ \ln(n_{\alpha\beta} + K - 1)! - \sum_{k=1}^{K} \ln(n_{\alpha\beta}^k)! \right]$

Bipartite weighted: $\mathcal{H}(\mathcal{P}_\mathcal{U}, \mathcal{P}_\mathcal{I}) = \sum_{\alpha, \beta} \left[ \ln(n_{\alpha\beta} + K - 1)! - \sum_{k=1}^{K} \ln(n_{\alpha\beta}^k)! \right]$

Guimera, Sales-Pardo, *PNAS* (2009)

Guimera, Sales-Pardo, *PLOS ONE* (2011)

Guimera, Llorente, Moro, Sales-Pardo, *PLOS ONE* (2012)

Rovira-Asenjo, Gumi, Sales-Pardo, Guimera, *in press* (2013)

# Reconstructing a network is more complicated than just adding missing interactions and removing spurious interactions

➔ Challenges:

  ➔ We don't know how many links need to be added and removed

  ➔ Links cannot be added and removed independently of each other

# We define a network reliability

➜ The reliability of a network is

$$p(A|A^O) = \frac{1}{Z} \sum_{P \in \mathcal{P}} f(A; A^O, P) \exp[-\mathcal{H}(\mathcal{P})]$$

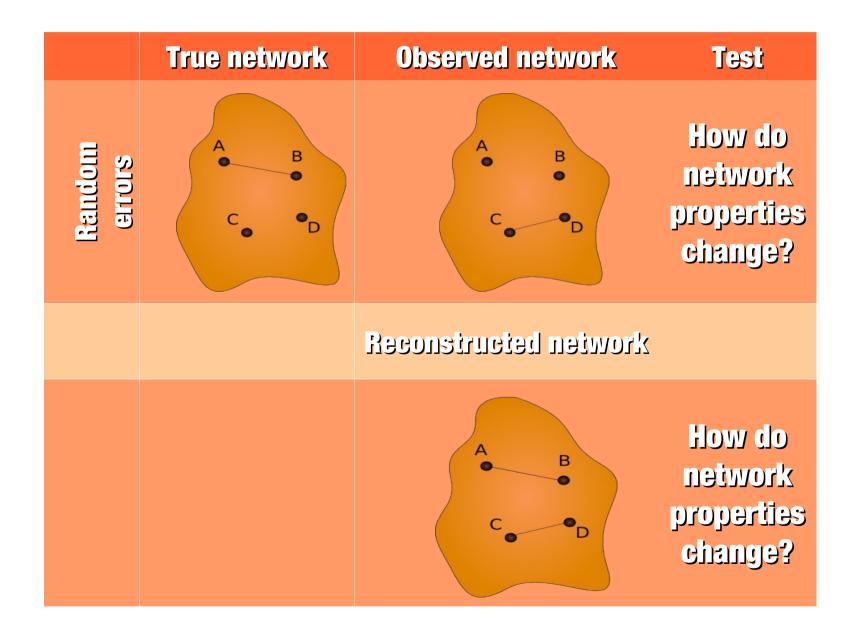Guimera, Sales-Pardo, *PNAS* (2009)

# The *network reconstruction* is the most reliable network
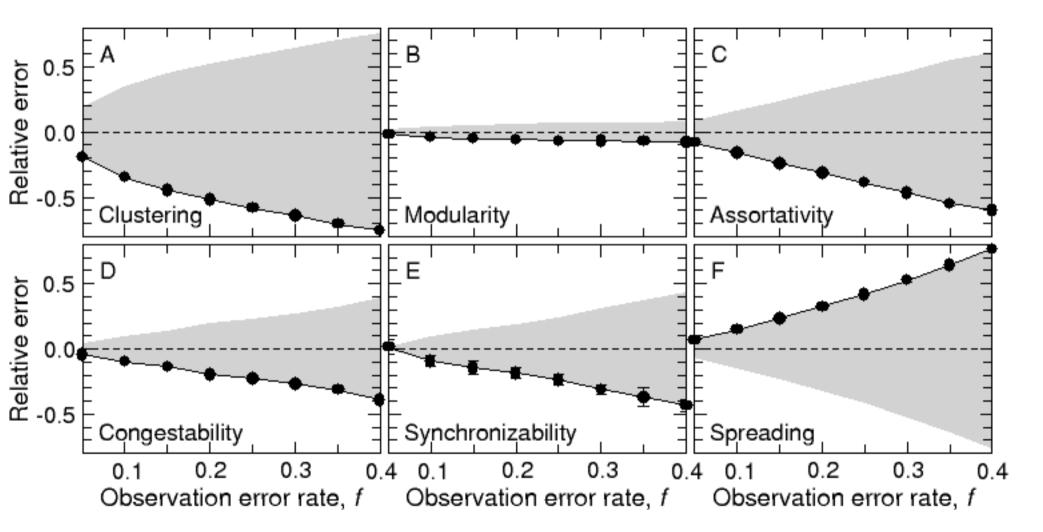
➜ The reliability of a network is

$$p(A|A^O) = \frac{1}{Z} \sum_{P \in \mathcal{P}} f(A; A^O, P) \exp[-\mathcal{H}(\mathcal{P})]$$

➜ The reconstruction $A^R$ is the network that maximizes this probability
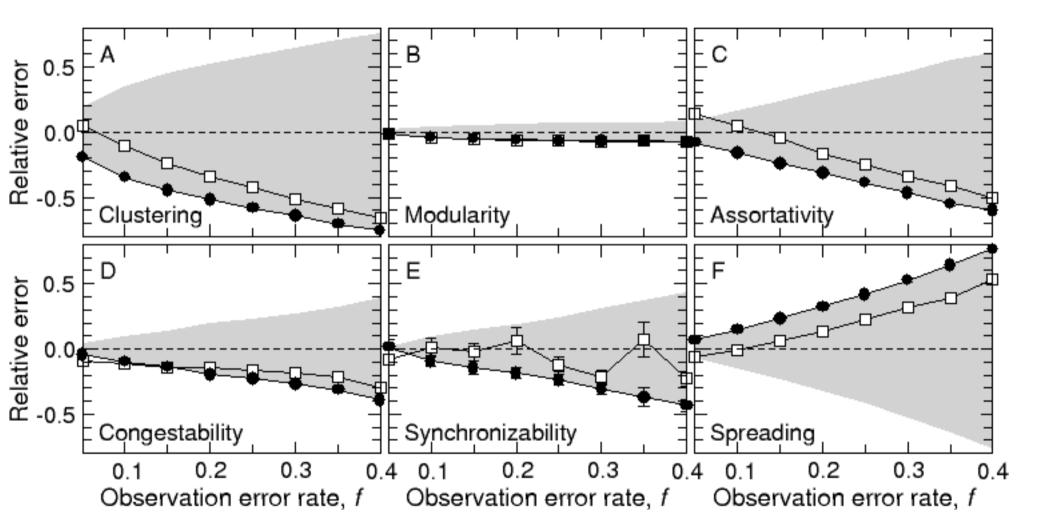
➜ We obtain $A^R$ using uphill search

# We can test what is the effect of random errors in our network observations

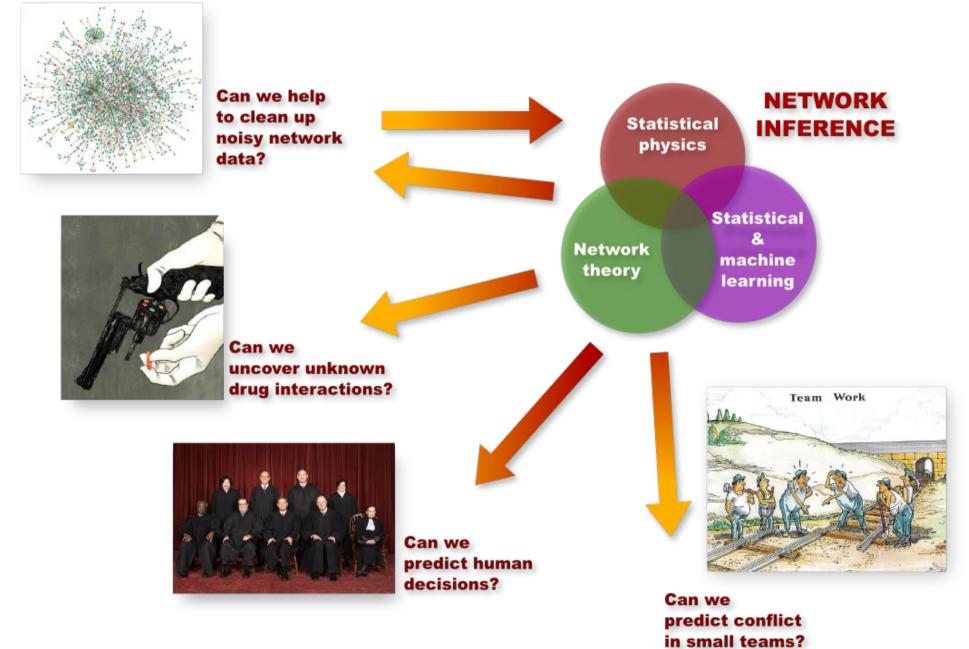# Network reconstructions provide better estimates of global network properties than the observations themselves

# Network reconstructions provide better estimates of global network properties than the observations themselves



Guimera, Sales-Pardo, *PNAS* (2009)

# Outline



Can we help to clean up noisy network data?

Can we uncover unknown drug interactions?

Can we predict human decisions?

**NETWORK INFERENCE**

Statistical physics

Network theory

Statistical & machine learning

Can we predict conflict in small teams?

# The challenge of discovering novel drug-drug interactions



→ Twenty-nine percent [of U.S. population aged 57-85] used at least 5 prescription medications concurrently.

→ Overall, 4% of individuals were potentially at risk of having a major drug-drug interaction.

Qato et al. *JAMA* (2008)

# Can we predict which severe drug interactions will be dded to / removed from a database?
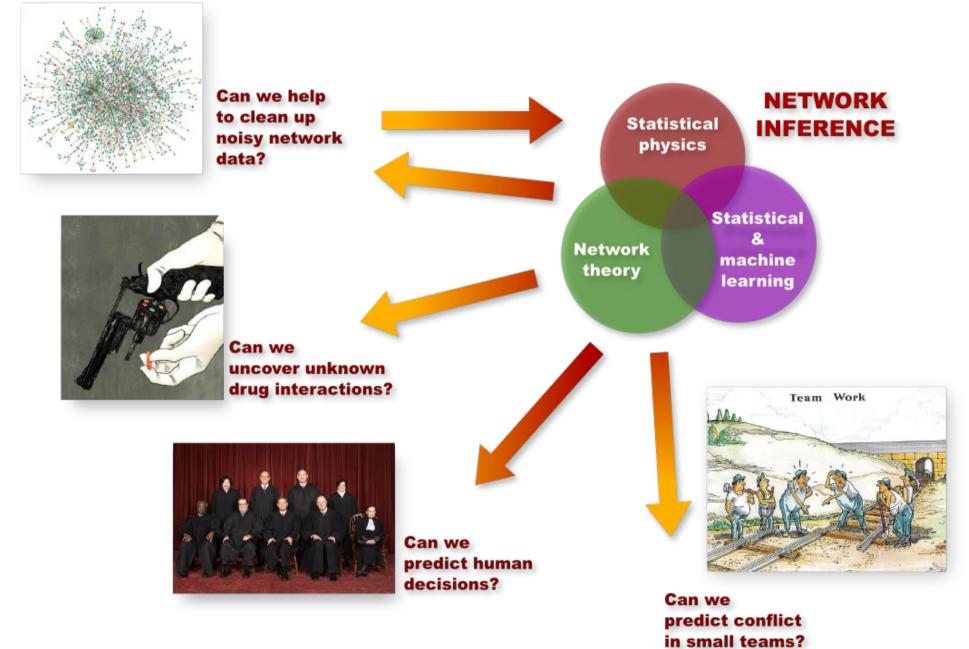


→ Two snapshots of the drug-interaction database available at *drugs.com*:

– May 10th, 2010

– February 22nd, 2012

→ Between the snapshots:

– 1349 interactions added

– 165 interactions removed

Guimera, Sales-Pardo, *submitted* (2013)

# We can predict which severe drug interactions will be removed from and added to a database



Guimera, Sales-Pardo, *submitted* (2013)

# Outline



Can we help to clean up noisy network data?

Can we uncover unknown drug interactions?

Can we predict human decisions?

Can we predict conflict in small teams?

**NETWORK INFERENCE**

Statistical physics

Statistical & machine learning

Network theory
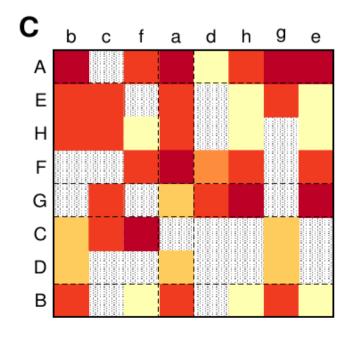
Team Work

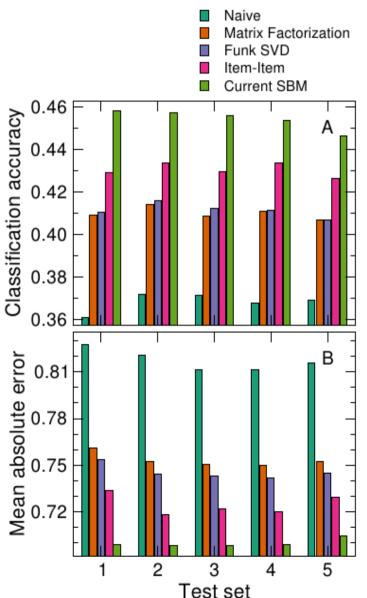# Predicting human preferences can be reformulated as a problem of network inference

# Our approach predicts human preferences better than state-of-the-art collaborative filtering algorithms

➜ MovieLens set: 100,000 real 1-5 movie ratings by ~1,000 users

➜ 5 independent splits of the data into 80,000 observed ratings and 20,000 validation ratings

Guimera, Llorente, Moro, Sales-Pardo (*PLOS ONE* 2012)

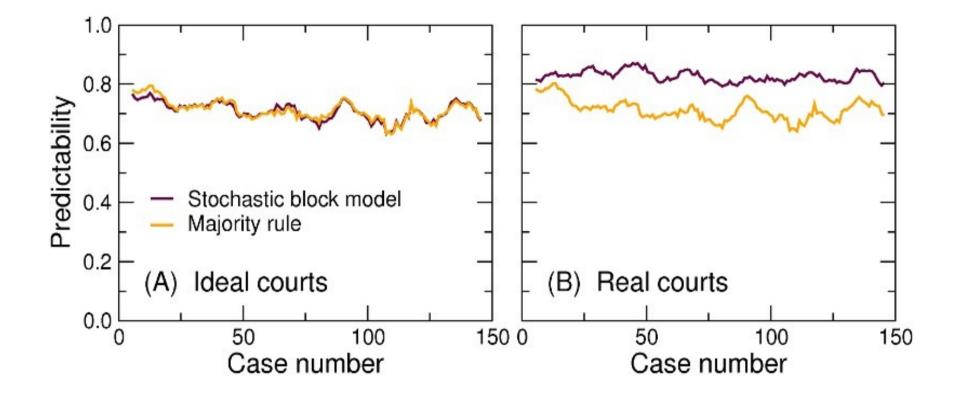# Our approach predicts human preferences better than state-of-the-art collaborative filtering algorithms

➔ MovieLens set: 100,000 real 1-5 movie ratings by ~1,000 users

➔ 5 independent splits of the data into 80,000 observed ratings and 20,000 validation ratings
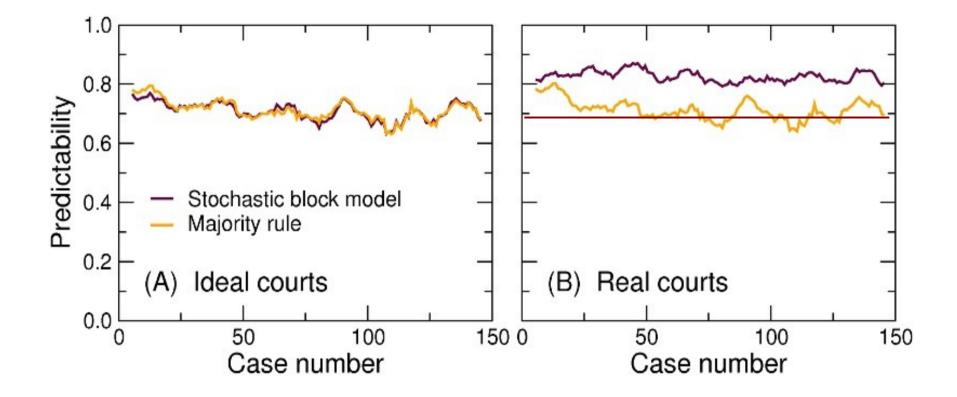
Guimera, Llorente, Moro, Sales-Pardo (*PLOS ONE* 2012)

Can we predict what a US Supreme Court justice votes based on what the others did?

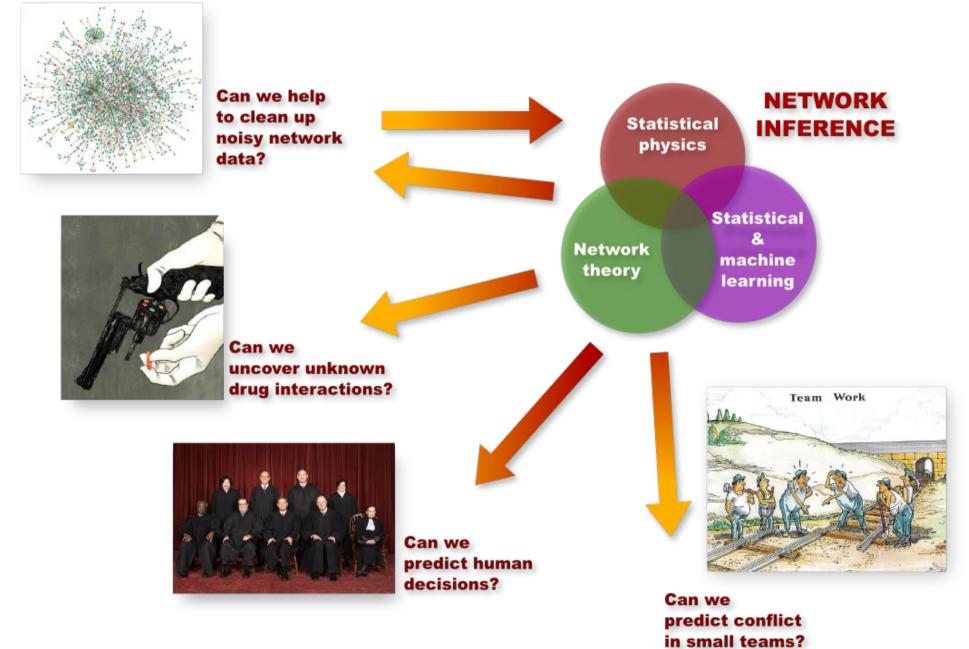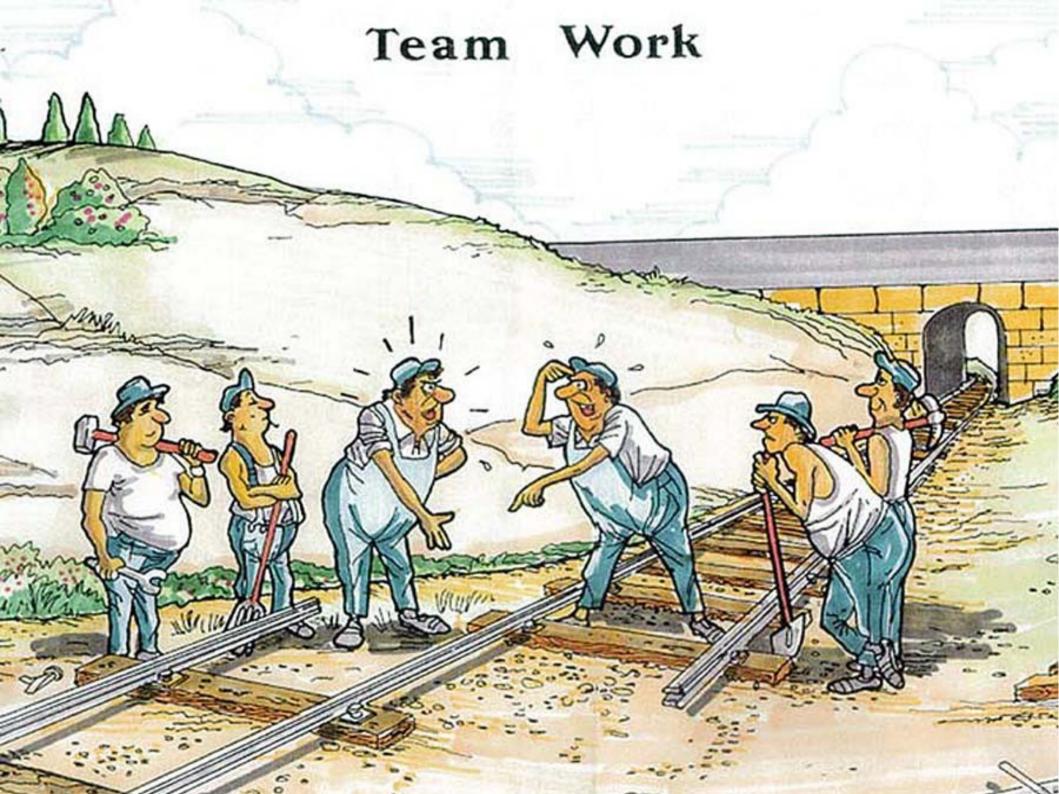# Supreme Court votes are more predictable than expected from ideal courts

# Supreme Court votes are more predictable than expected from ideal courts

# Outline



Can we help to clean up noisy network data?

Can we uncover unknown drug interactions?

Can we predict human decisions?

Can we predict conflict in small teams?

**NETWORK INFERENCE**

Statistical physics

Network theory
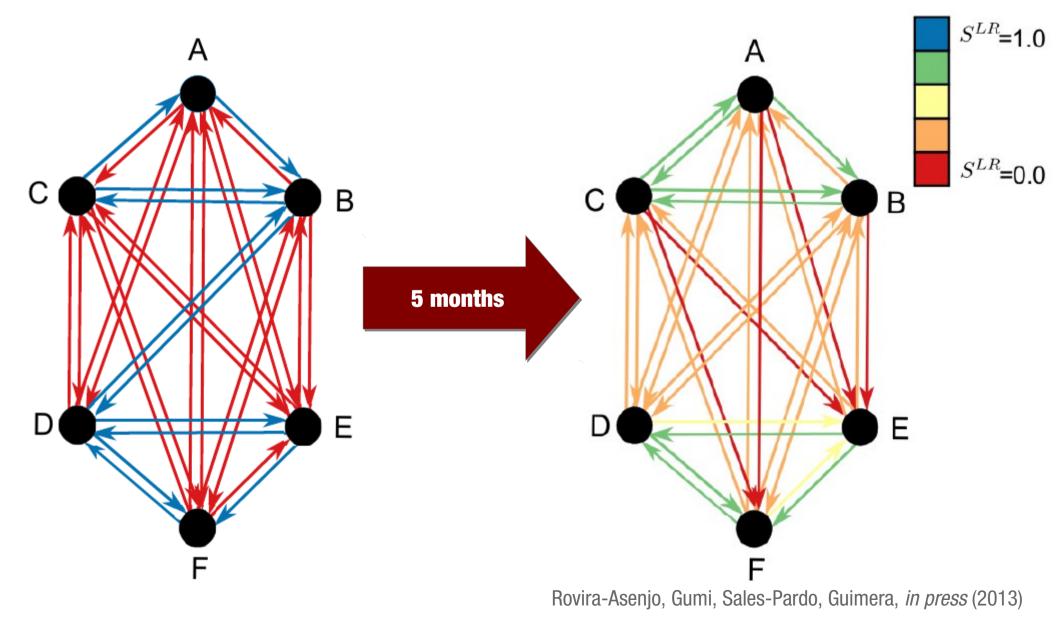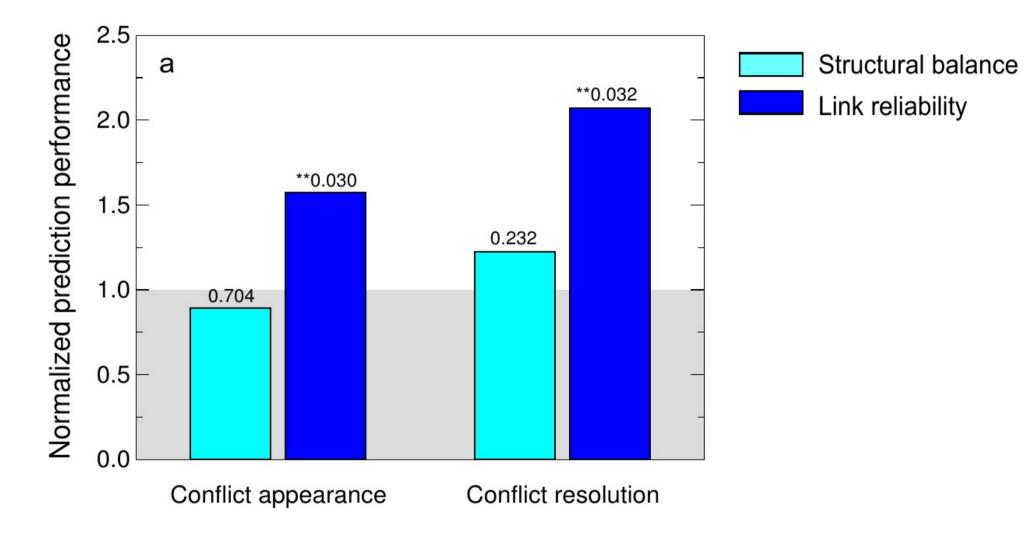
Statistical & machine learning

Team Work

Team Work

# Tracking team conflict in the real world

➔ 16 teams with ~6 people, working on a real project during 9 months

➔ We administer 2 surveys:

   ➔ First: After 4 months working together

   ➔ Second: At the end of the project

➔ "Would you like to work with this person again in the future"

# Can we predict where conflict is going to arise and where it is going to resolve?



Rovira-Asenjo, Gumi, Sales-Pardo, Guimera, *in press* (2013)

# Our approach predicts conflict appearance and conflict resolution whereas structural balance does not



Rovira-Asenjo, Gumi, Sales-Pardo, Guimera, *in press* (2013)

# Outline



Can we help to clean up noisy network data?

Can we uncover unknown drug interactions?

Can we predict human decisions?

Can we predict conflict in small teams?

**NETWORK INFERENCE**

Statistical physics

Network theory

Statistical & machine learning

Team Work

# Thank you

➔ T. Gumí, A. Llorente, E. Moro, N. Rovira-Asenjo, M. Sales-Pardo

➔ Funding

➔ More information:

– http://seeslab.info

– @sees_lab