

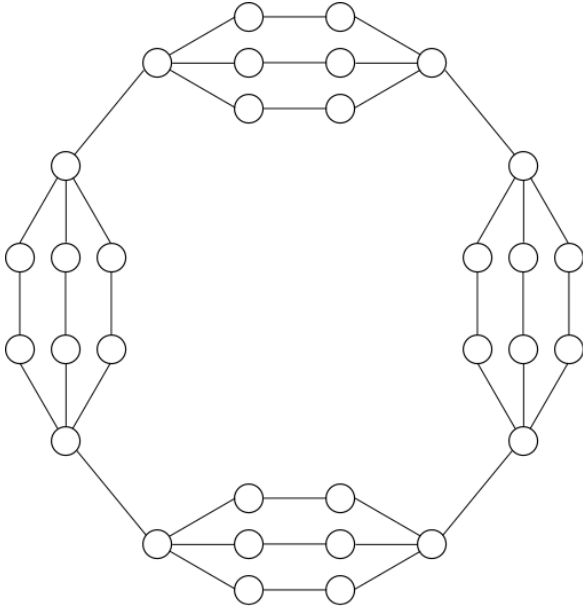
# **Beyond networks: Incorporating node metadata into network analysis**

**Leto Peel**

**Université catholique de Louvain**

**@PiratePeel**

**Here is a network  $G=(V,E)$**



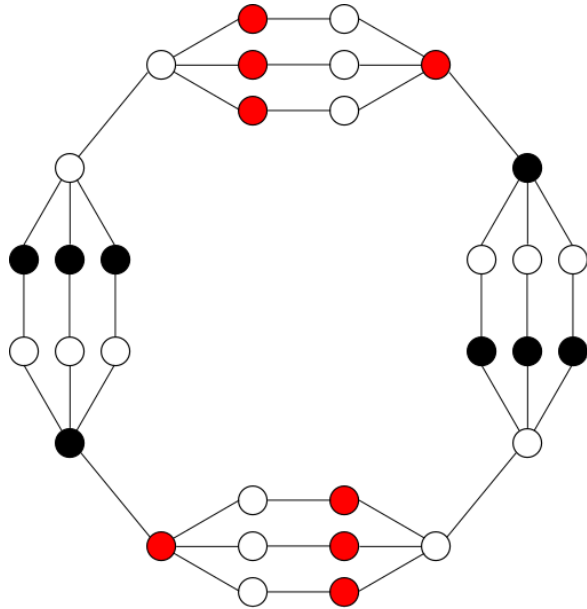
**social networks**

**food webs**

**internet**

**protein interactions**

# Network nodes can have properties or attributes (metadata)



- Metadata (M) values
- Metadata (M) unknown

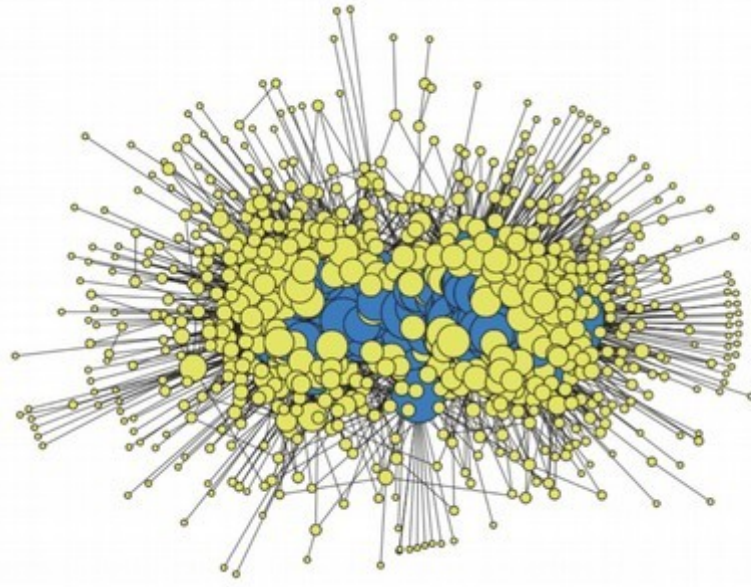
**social networks** *age, sex, ethnicity, race, etc.*

**food webs** *feeding mode, species body mass, etc.*

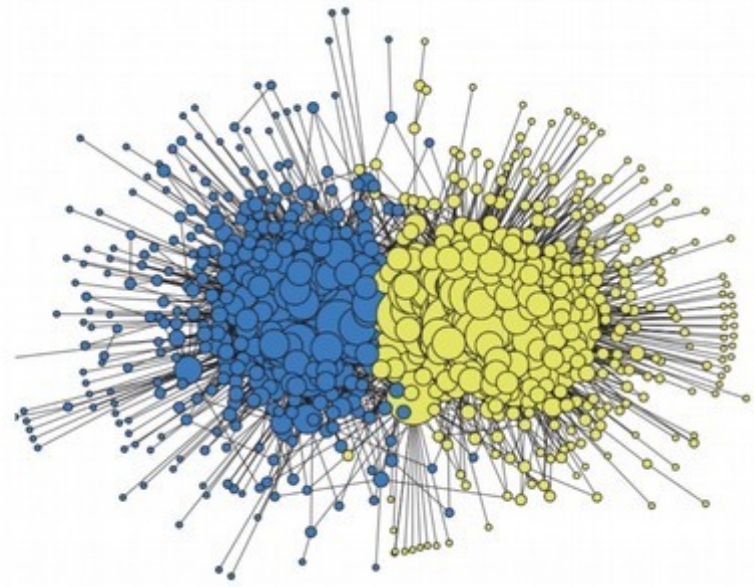
**internet** *data capacity, physical location, etc.*

**protein interactions** *molecular weight, association with cancer, etc.*

# Recovering metadata implies sensible methods



stochastic block model

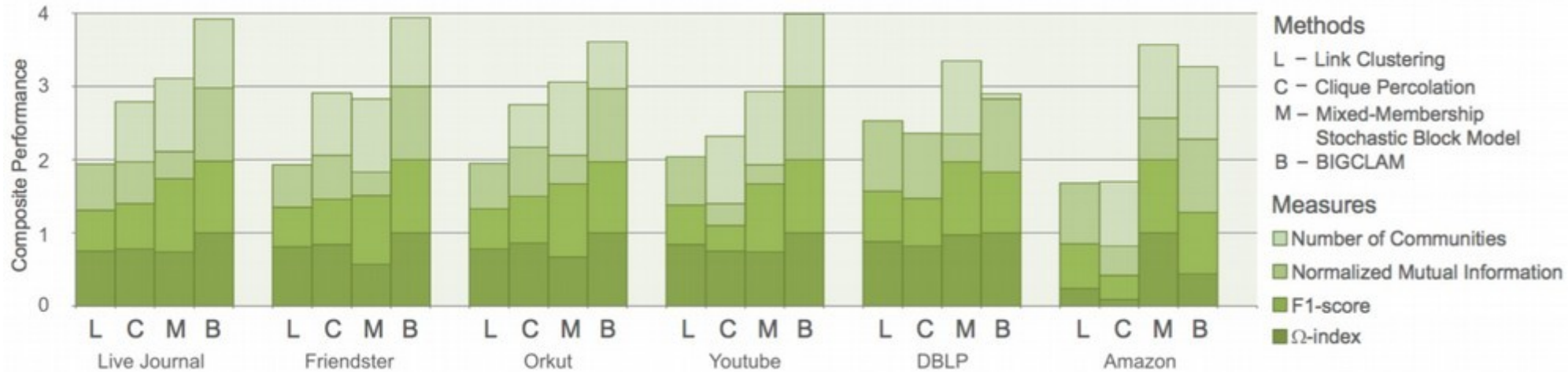


stochastic block model  
with degree correction

Karrer, Newman. Stochastic blockmodels and community structure in networks. Phys. Rev. E 83, 016107 (2011).

Adamic, Glance. The political blogosphere and the 2004 US election: divided they blog. 36–43 (2005).

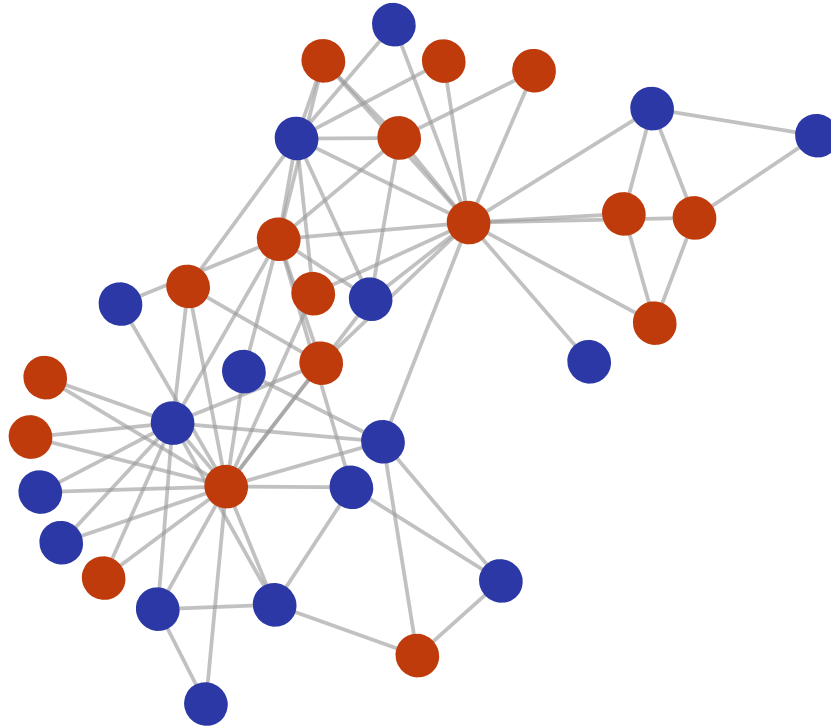
# Metadata is often treated as *ground truth*



*You know, I know these  
communities aren't real...*

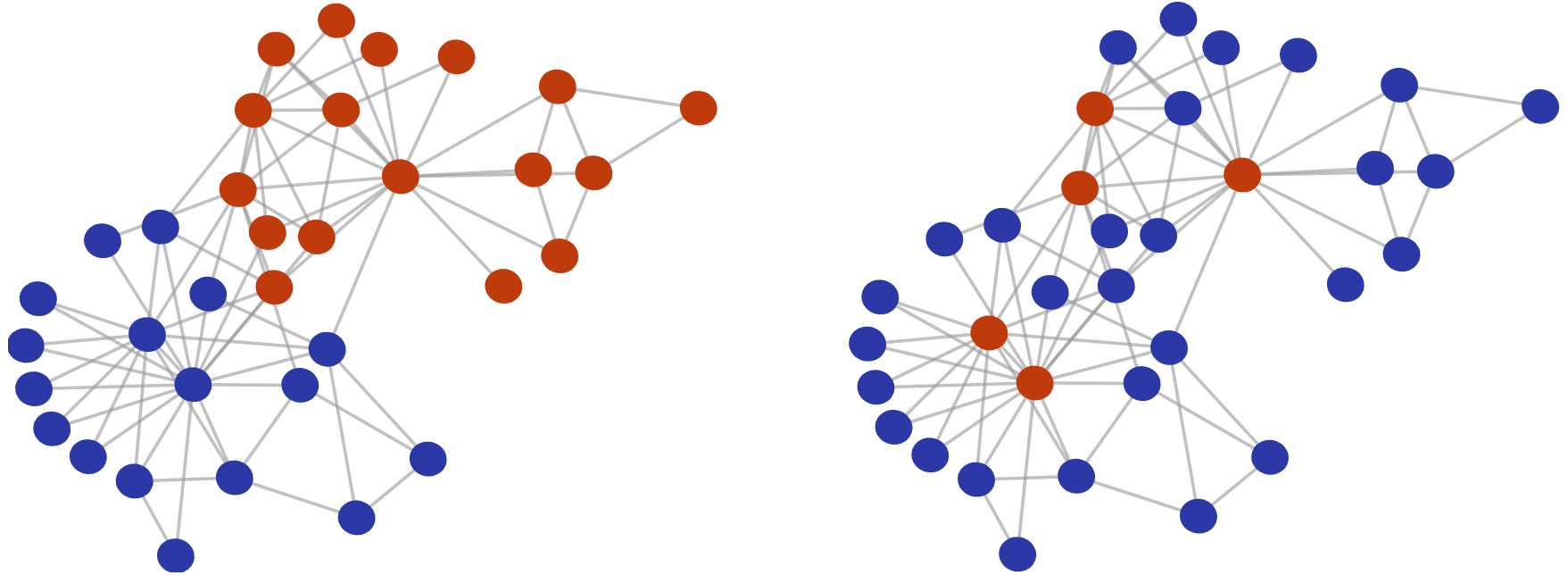


# When communities $\neq$ metadata...



(i) the metadata do not relate to the network structure,

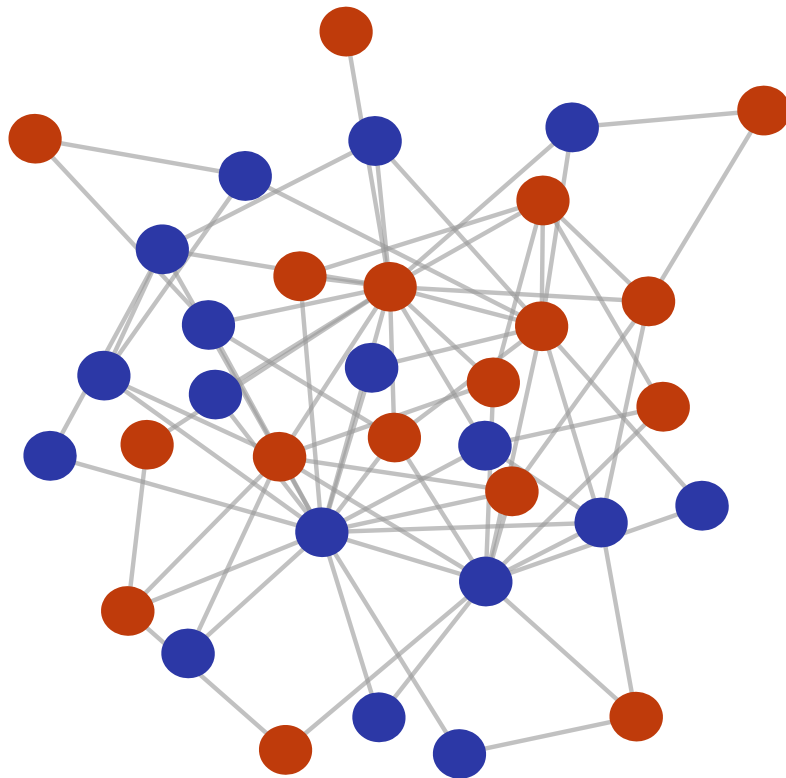
# When communities $\neq$ metadata...



(ii) the detected communities and the metadata capture different aspects of the network's structure,

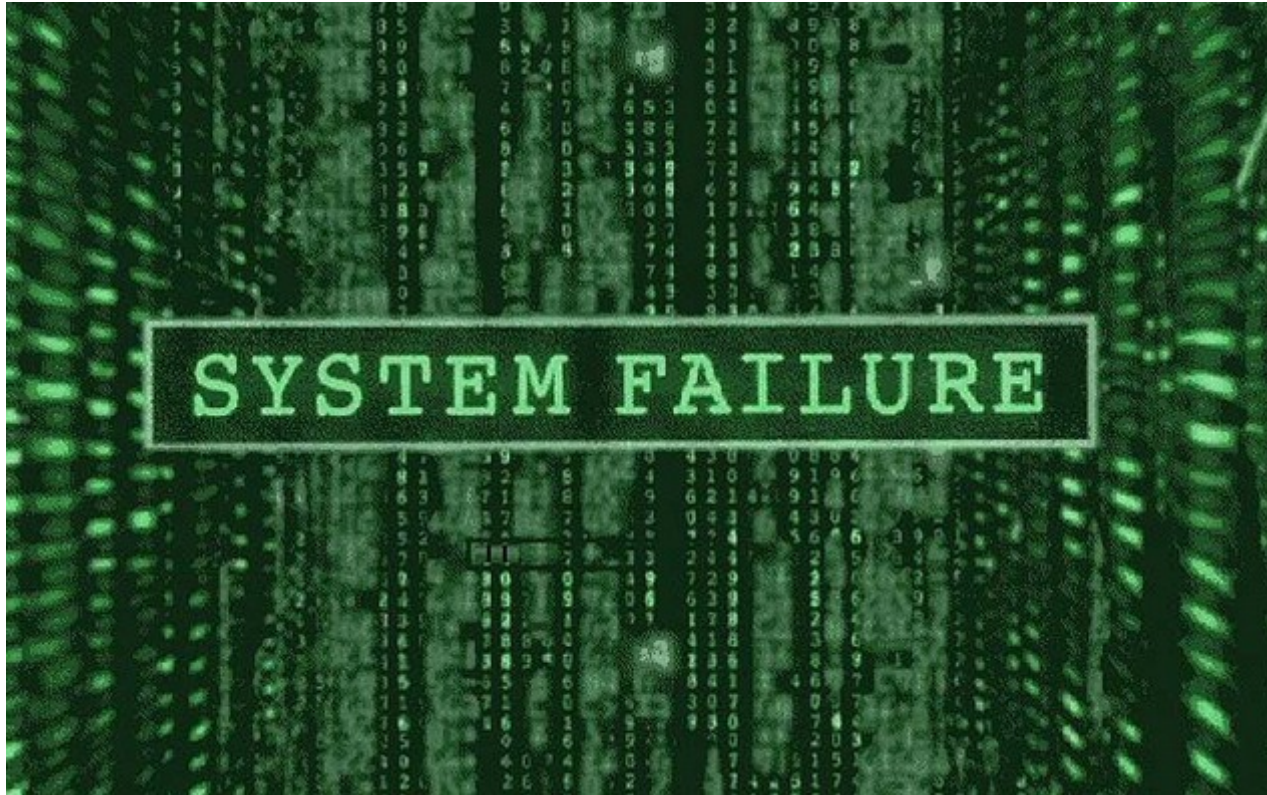


# When communities $\neq$ metadata...



(iii) the network contains no structure (e.g., an E-R random graph)

# When communities $\neq$ metadata...



(iv) the community detection algorithm does not perform well.

Typically we assume this is the only possible cause

## **Metadata are not ground truth for community detection**

### **No interpretability of negative results.**

- (i) M unrelated to network structure
- (ii) C and M capture different aspects of network structure
- (iii) the network has no structure
- (iv) the algorithm does not perform well

## Metadata are not ground truth for community detection

### **No interpretability of negative results.**

- (i) M unrelated to network structure
- (ii) C and M capture different aspects of network structure
- (iii) the network has no structure
- (iv) the algorithm does not perform well

### **Multiple sets of metadata exist.**

Which set is ground truth?

## **Metadata are not ground truth for community detection**

### **No interpretability of negative results.**

- (i) M unrelated to network structure
- (ii) C and M capture different aspects of network structure
- (iii) the network has no structure
- (iv) the algorithm does not perform well

### **Multiple sets of metadata exist.**

Which set is ground truth?

### **We see what we look for.**

Confirmation bias. Publication bias.

## **Metadata are not ground truth for community detection**

### **No interpretability of negative results.**

- (i) M unrelated to network structure
- (ii) C and M capture different aspects of network structure
- (iii) the network has no structure
- (iv) the algorithm does not perform well

### **Multiple sets of metadata exist.**

Which set is ground truth?

### **We see what we look for.**

Confirmation bias. Publication bias.

### **“Community” is model dependent.**

Do we expect all networks across all domains to have the same relationship with communities?

**DON'T TRY TO FIND THE GROUND TRUTH**

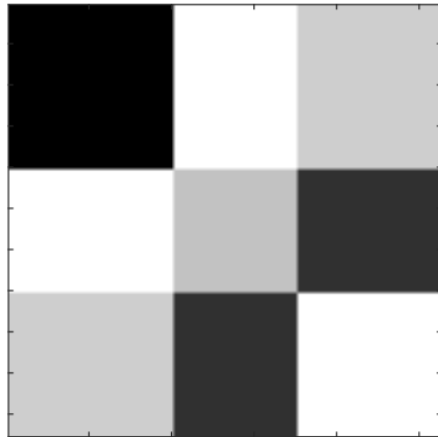
**INSTEAD... TRY TO REALIZE THERE IS NO GROUND TRUTH**



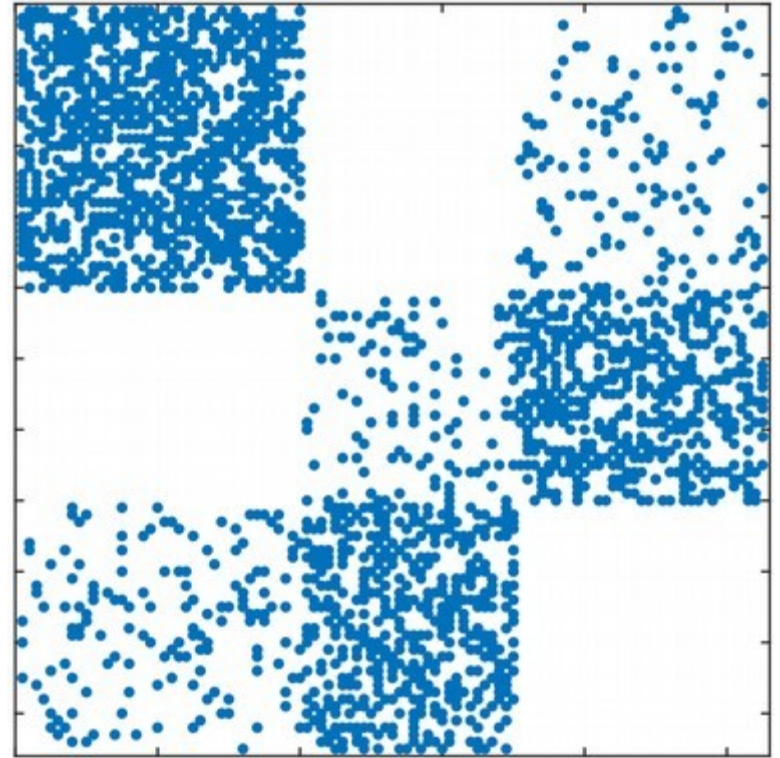
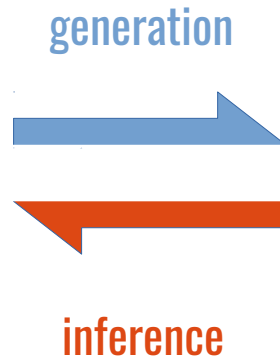
**How should we use metadata?**



# The Stochastic Blockmodel (SBM)



Mixing Matrix



Adjacency Matrix

# The Stochastic Blockmodel (SBM)



# Blockmodel Entropy Significance Test

*How well do the metadata explain the network?*

# Blockmodel Entropy Significance Test

*How well do the metadata explain the network?*

1. Divide the network  $G$  into groups according to metadata labels  $M$ .
2. Fit the parameters of an SBM and compute the entropy  $H(G,M)$
3. Compare this entropy to a distribution of entropies of networks partitioned using permutations of the metadata labels.

# Blockmodel Entropy Significance Test

*How well do the metadata explain the network?*

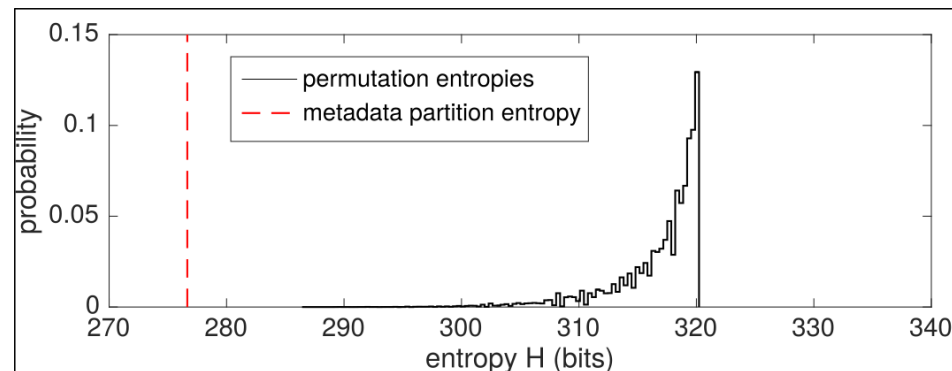
1. Divide the network  $G$  into groups according to metadata labels  $M$ .

2. Fit the parameters of an SBM and compute the entropy  $H(G,M)$

3. Compare this entropy to a distribution of entropies of networks partitioned using permutations of the metadata labels.

metadata is randomly assigned  
→ model gives no explanation, high  $H$

metadata correlates with structure  
→ model gives good explanation, low  $H$

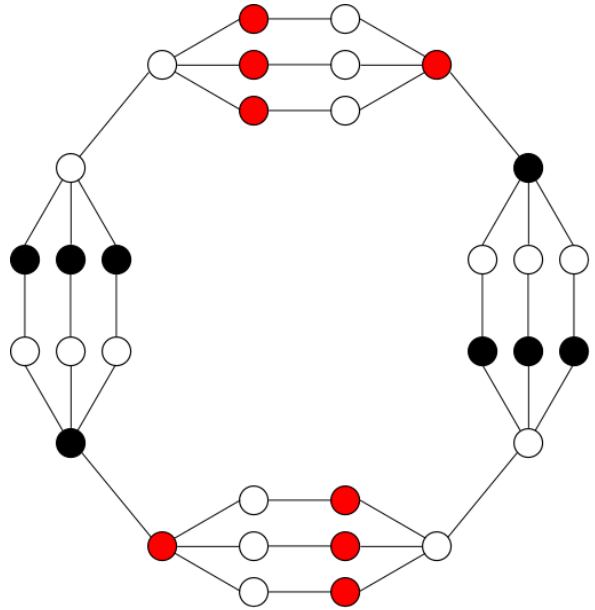


# Multiple networks; multiple metadata attributes

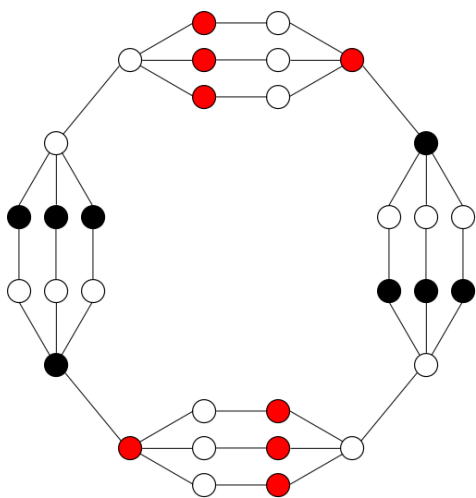
Network	Status	Gender	Office	Practice	Law School
Friendship	$< 10^{-6}$	0.034	$< 10^{-6}$	0.033	0.134
Cowork	$< 10^{-3}$	0.094	$< 10^{-6}$	$< 10^{-6}$	0.922
Advice	$< 10^{-6}$	0.010	$< 10^{-6}$	$< 10^{-6}$	0.205

Multiple sets of metadata provide a significant explanation for multiple networks.

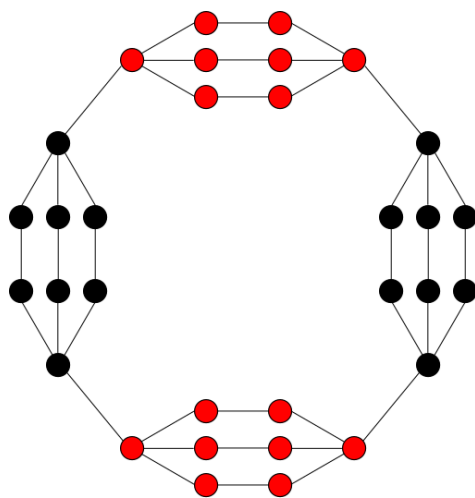
# Can we predict unknown metadata values?



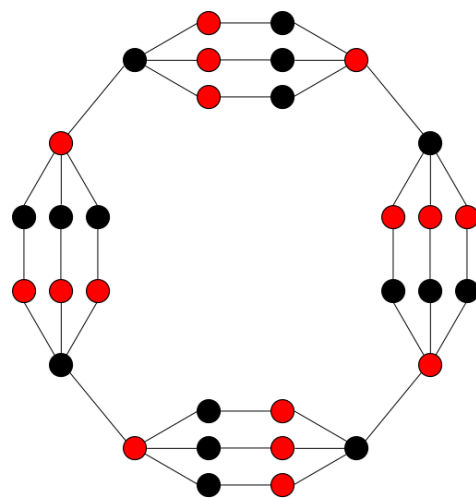
- Metadata values
- Metadata unknown



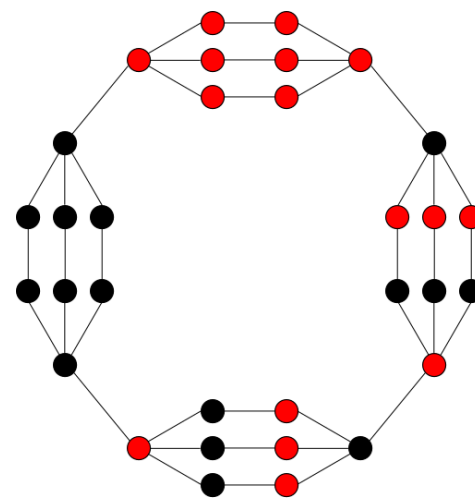
*assortative*



*disassortative*

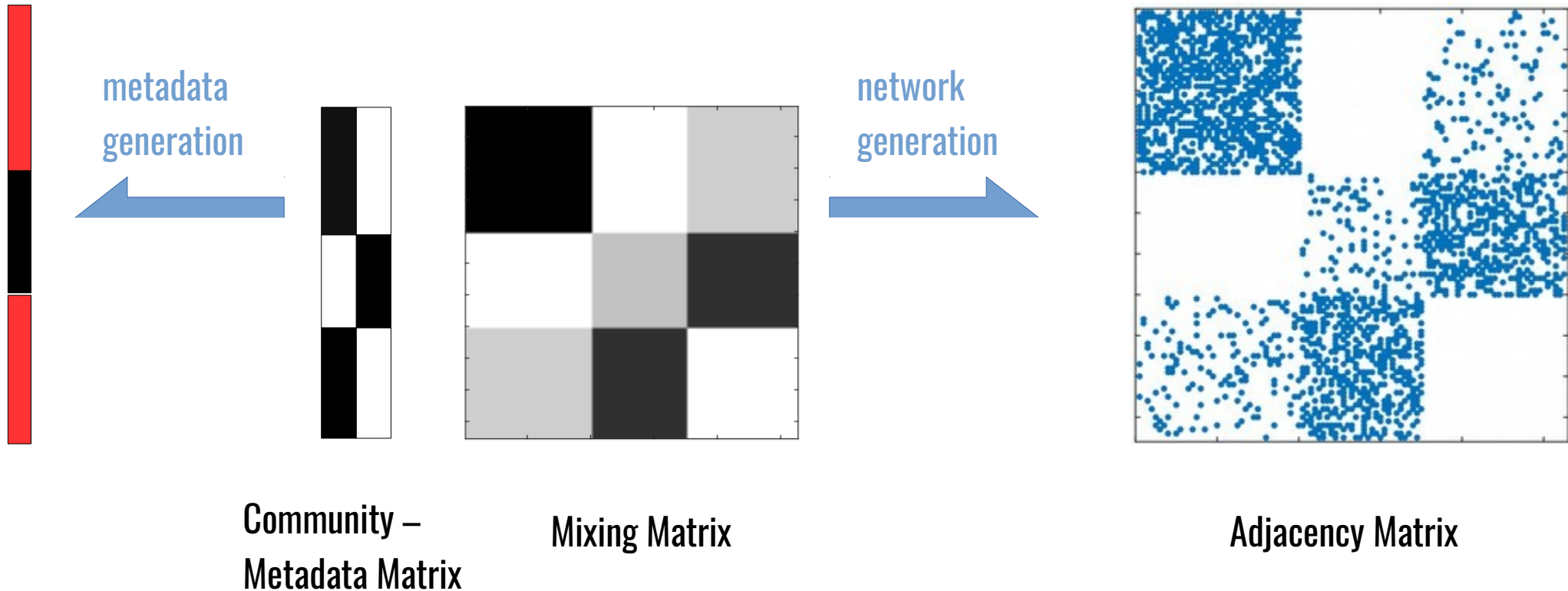


*mixed*

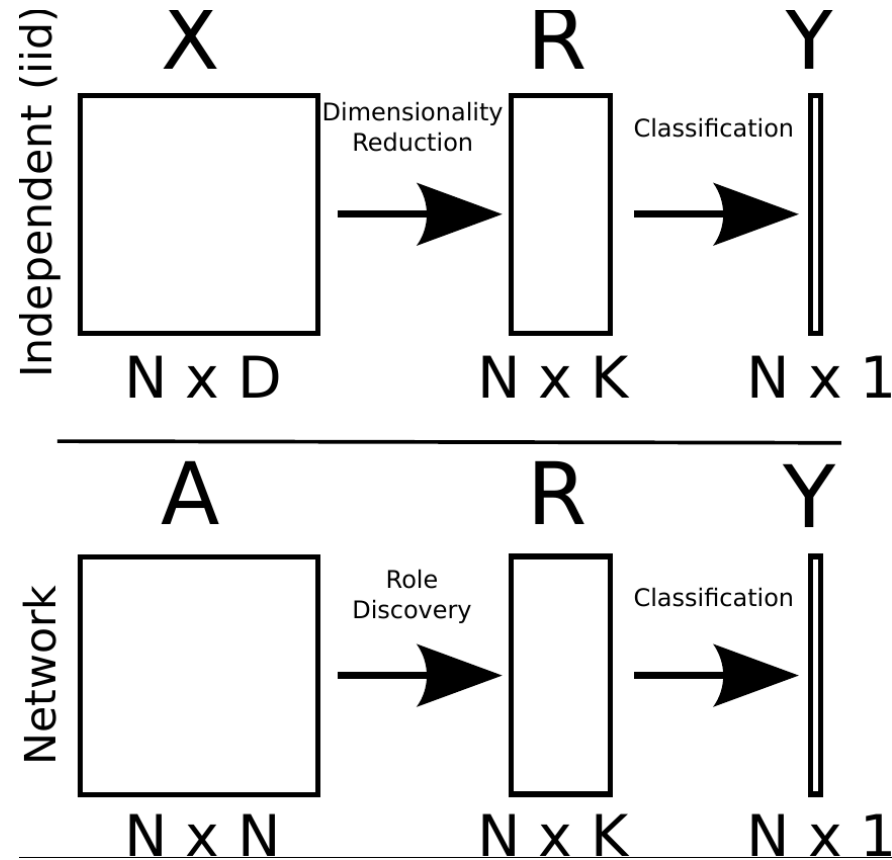


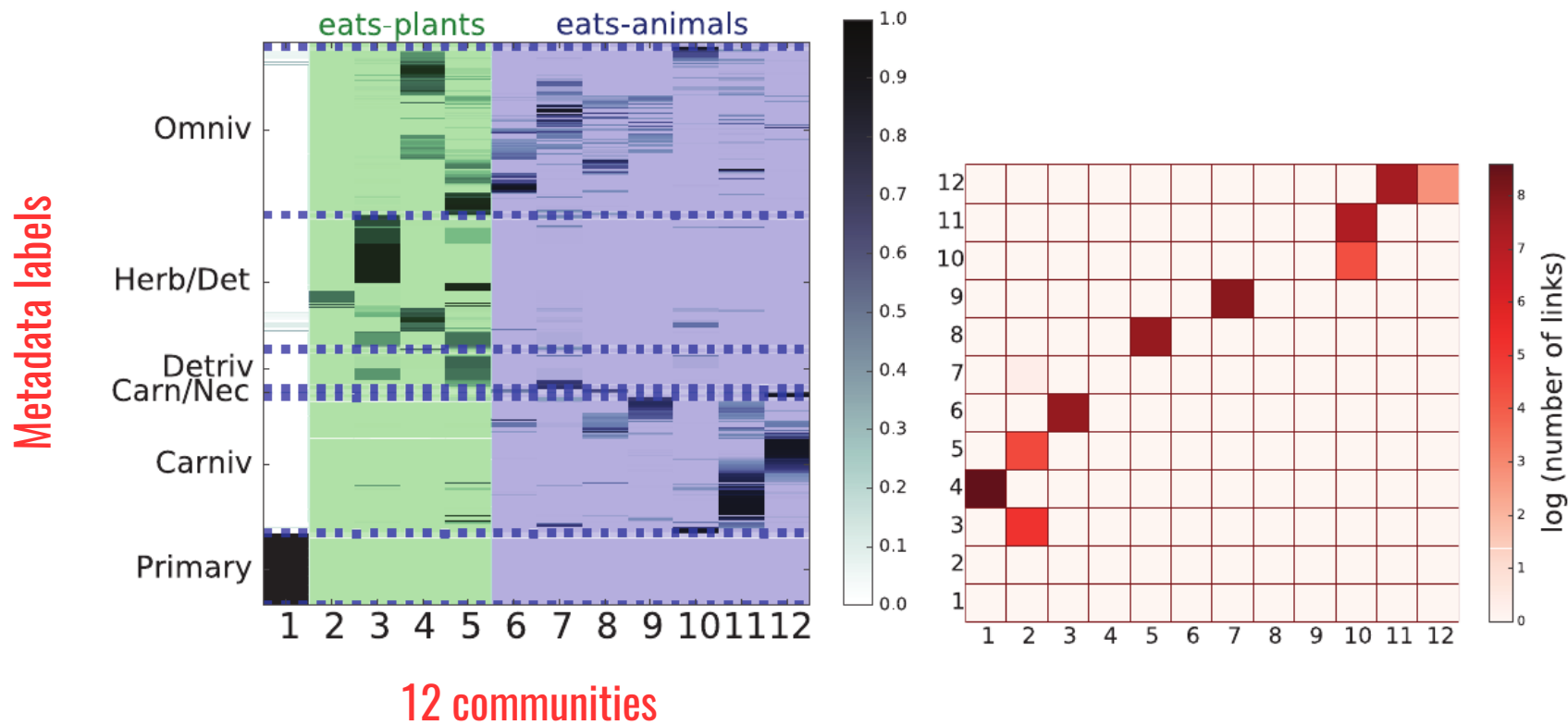


# A stochastic block model with metadata

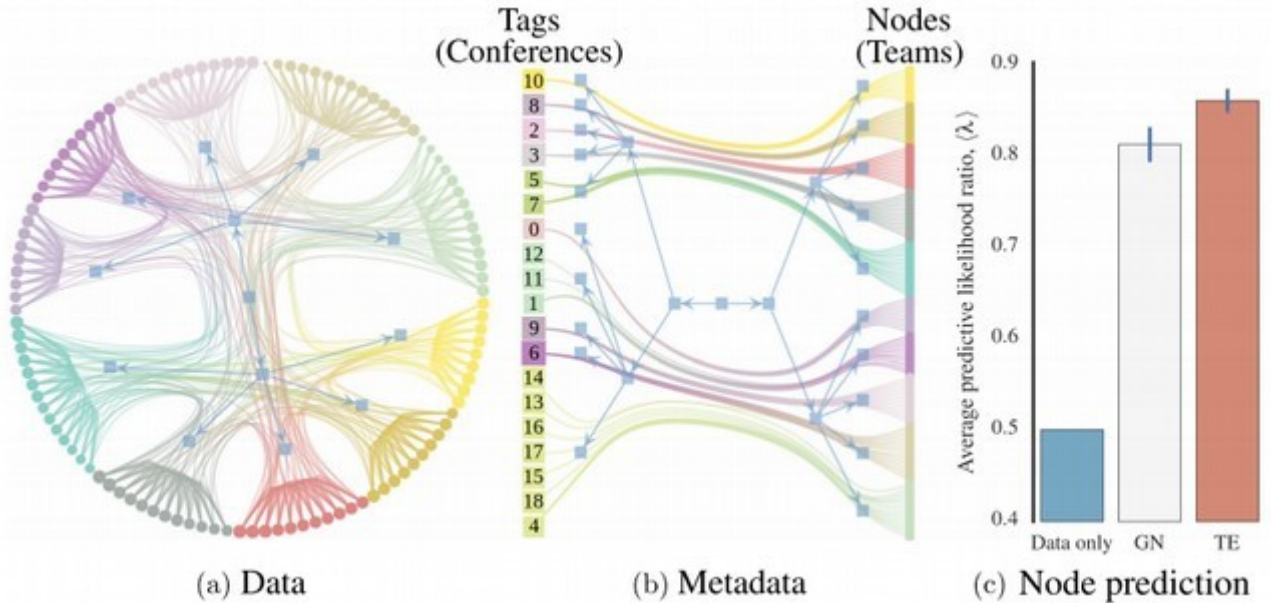
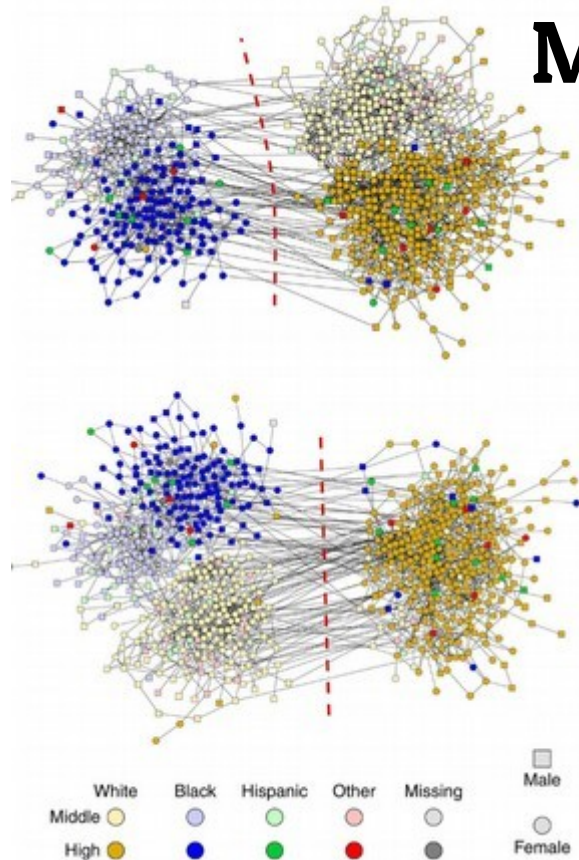


# Dimensionality reduction + classification





# More SBMs + metadata

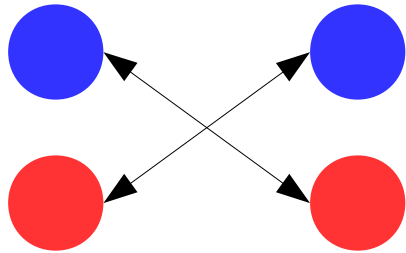


Newman, Clauset. "Structure and inference in annotated networks." Nat. Comms. 7 (2016).

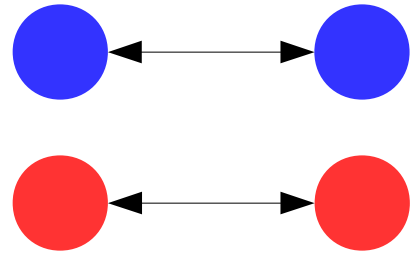
Hric, Peixoto, Fortunato. "Network structure, metadata, and the prediction of missing nodes and annotations." Phys. Rev. X 6.3: 031038 (2016)

# Mixing patterns in networks

$$r_{\text{global}} = \frac{\sum_g e_{gg} - \sum_g a_g b_g}{1 - \sum_g a_g b_g}$$



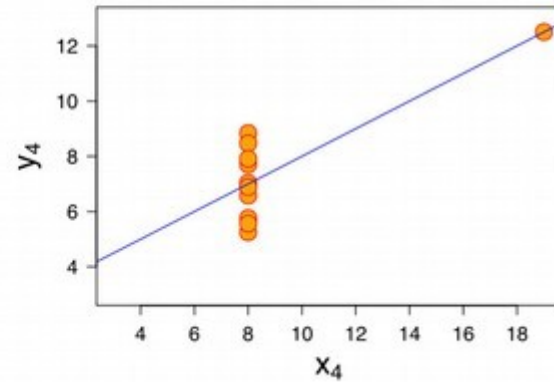
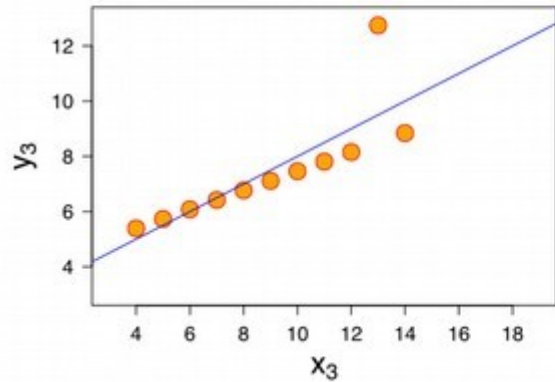
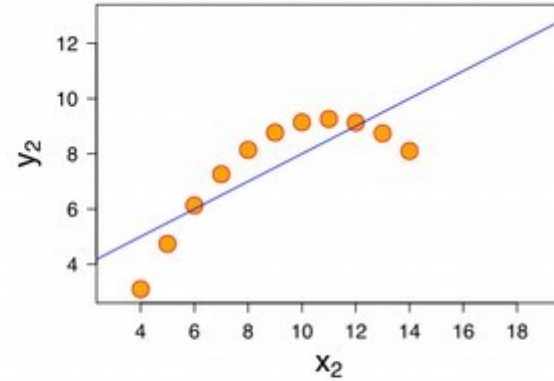
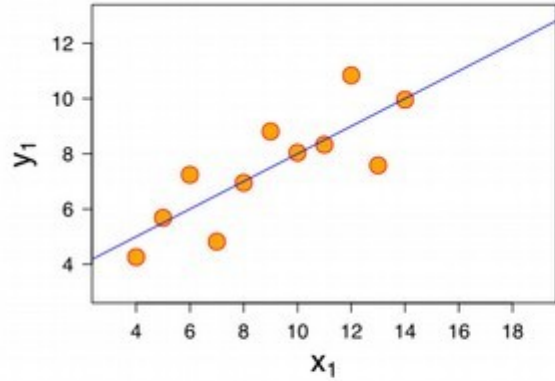
disassortative



assortative

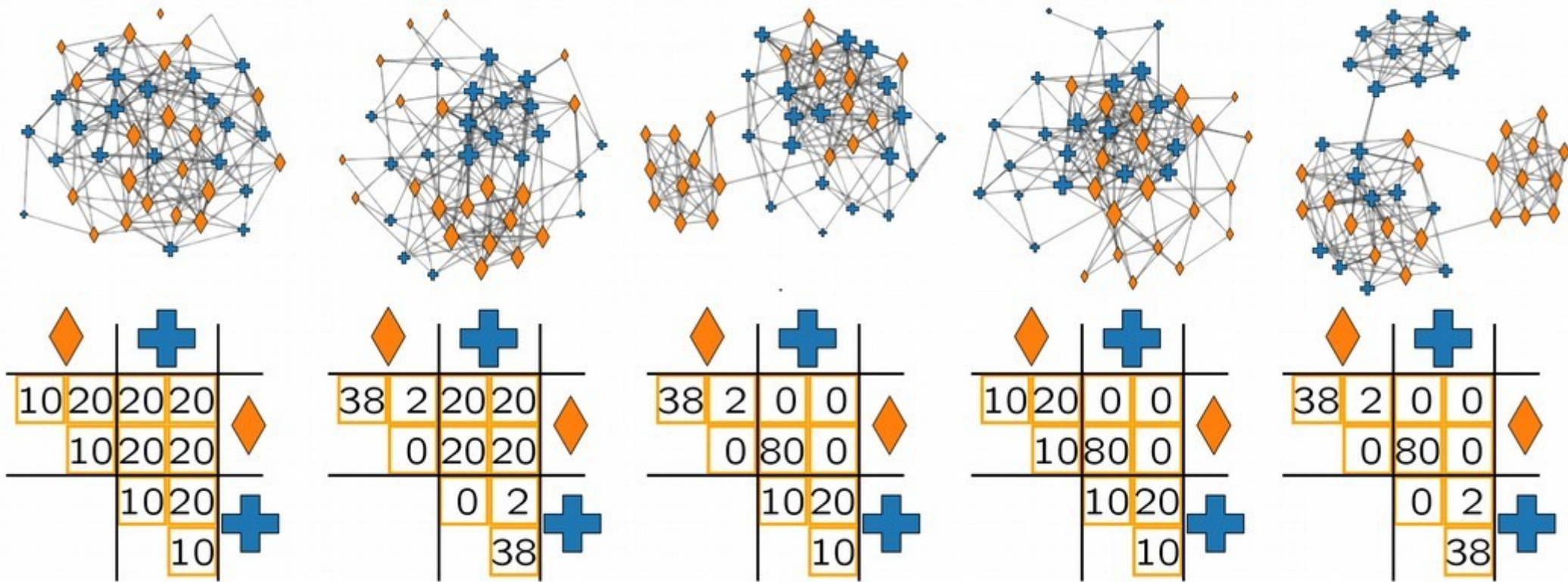


# Assortativity is correlation across edges

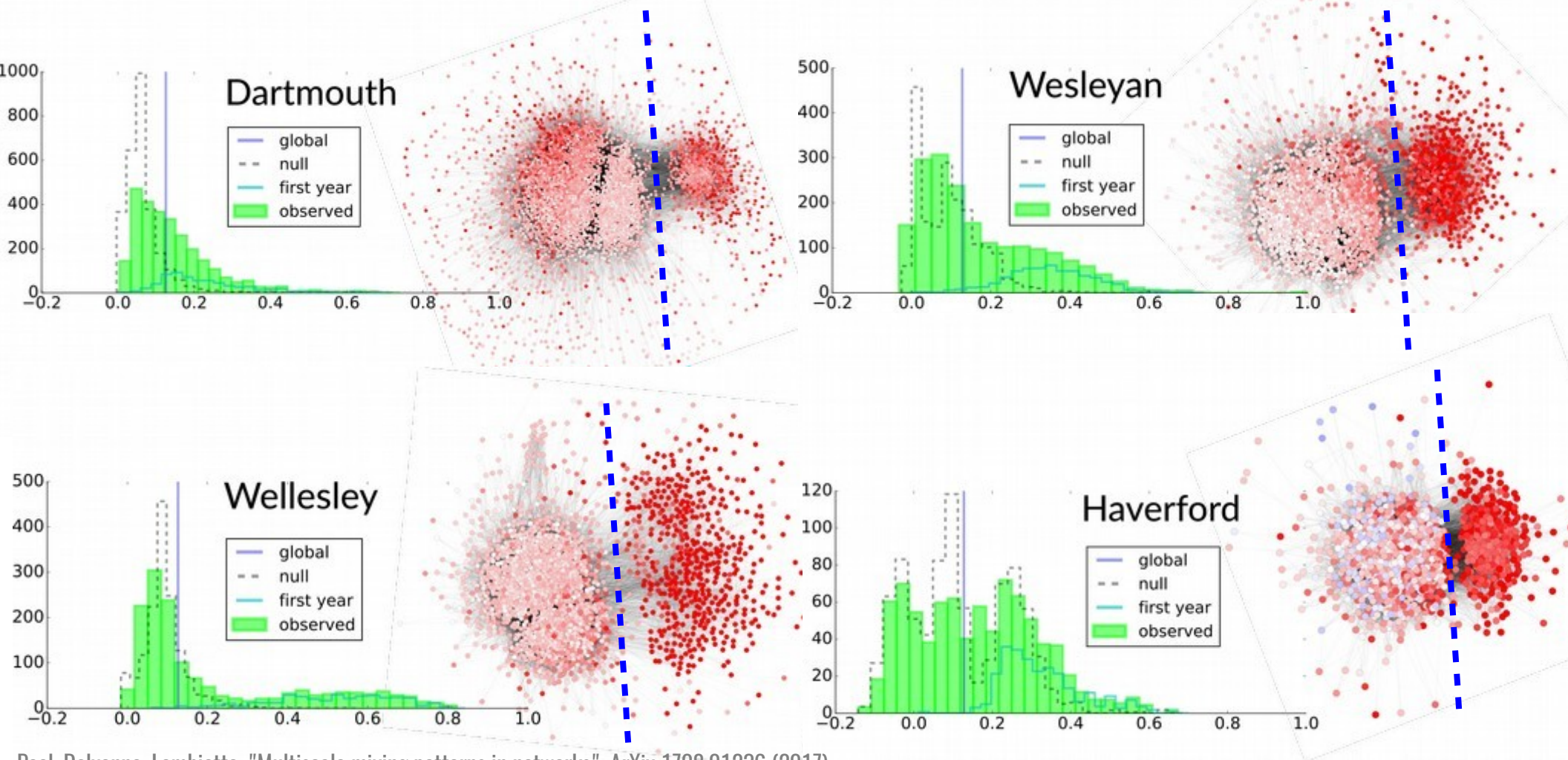




# All these networks have assortativity $r=0$



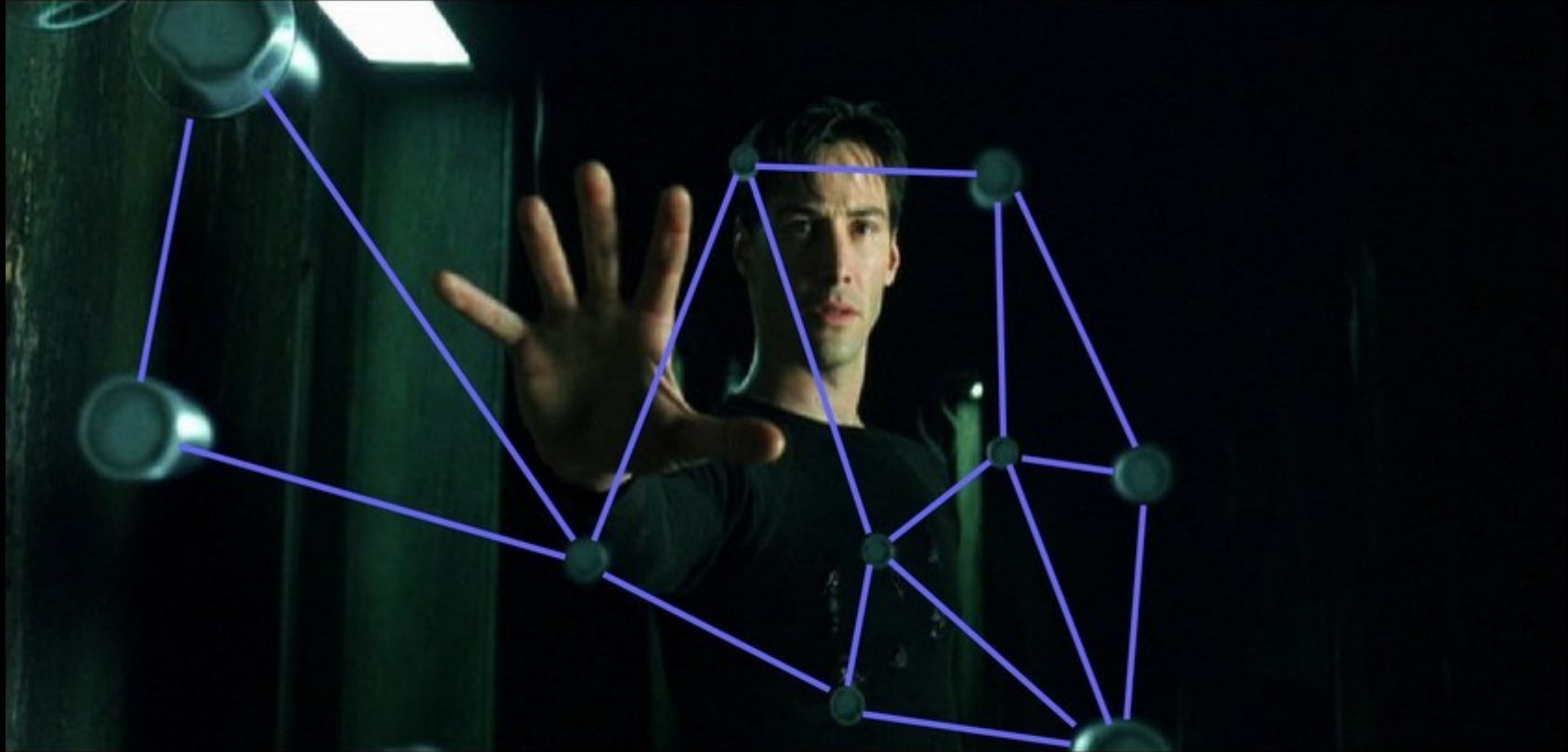
# Facebook 100 – residence





**Final thoughts...**

Am I telling you that you can detect ground truth communities?



No, I'm telling you that when you're ready, you won't have to.

# **Final thoughts...**

**The relationship between network structure and metadata is important:**

- Determine if the relationship is significant**
- Predict missing values**
- Understand how assortativity varies over a network**

# References & collaborators...

Peel, L., **Topological Feature Based Classification** 14th International Conference on Information Fusion (FUSION) 2011

Peel, L., **Supervised Blockmodelling** ECML/PKDD Workshop on Collective Learning and Inference on Structured Data (CoLISD) 2012

Peel, L., **Active Discovery of Network Roles for Predicting the Classes of Network Nodes** Journal of Complex Networks 3 (3): 431-449, 2015

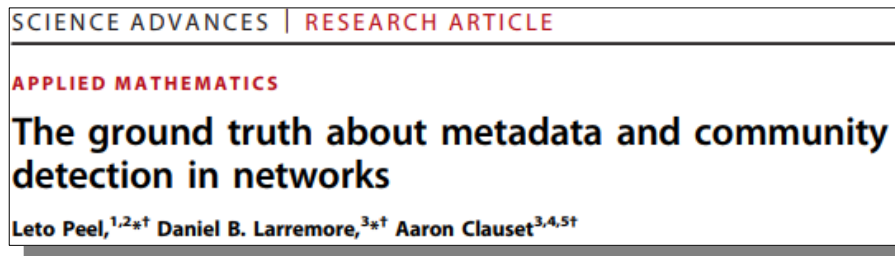
Peel, L., **Graph-based semi-supervised learning for relational networks** SIAM International Conference on Data Mining (SDM) 2017



Daniel B.  
Larremore



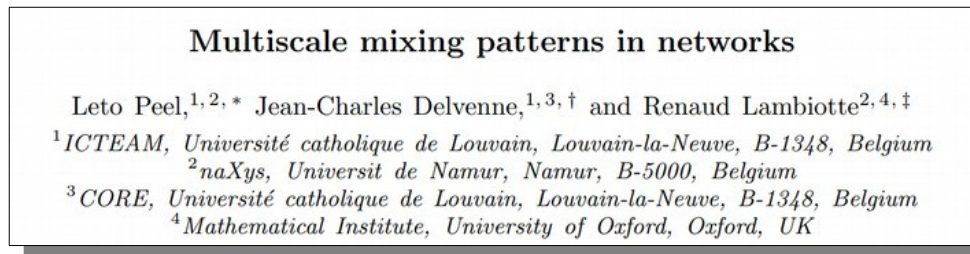
Aaron  
Clauset



Jean-Charles  
Delvenne



Renaud  
Lambiotte



pre-print arXiv:1708.01236

@PiratePeel

Contact: leto.peel@uclouvain.be