# Network structure, metadata and the prediction of missing nodes

Tiago P. Peixoto

*Universität Bremen, Germany & ISI Foundation, Italy*

Darko Hric

*Aalto University, Finland*

Santo Fortunato

*Indiana University, USA*

Seoul, May 2016

# NETWORKS WITH METADATA

Many network datasets contain *metadata*: Annotations that go beyond the mere adjacency between nodes.

Often assumed as indicators of topological structure, and used to *validate* community detection methods. A.k.a. "ground-truth".

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Metadata (Conferences)

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



SBM fit

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

Why the discrepancy?

Some hypotheses:

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

## Why the discrepancy?

Some hypotheses:

- ▶ The model is not sufficiently descriptive.

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

## Why the discrepancy?

Some hypotheses:

- ► The model is not sufficiently descriptive.
- ► The metadata is not sufficiently descriptive or is inaccurate.

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

Why the discrepancy?

Some hypotheses:

- ▶ The model is not sufficiently descriptive.
- ▶ The metadata is not sufficiently descriptive or is inaccurate.
- ▶ Both.

# EXAMPLE: AMERICAN COLLEGE FOOTBALL



Discrepancy

Why the discrepancy?

Some hypotheses:

- ▶ The model is not sufficiently descriptive.
- ▶ The metadata is not sufficiently descriptive or is inaccurate.
- ▶ Both.
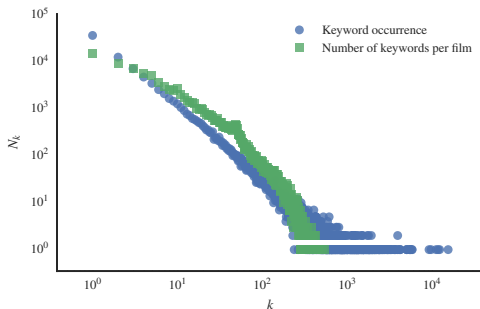- ▶ Neither.

# METADATA IS OFTEN VERY HETEROGENEOUS

## EXAMPLE: IMDB FILM-ACTOR NETWORK

Data: $96,982$ Films, $275,805$ Actors, $1,812,657$ Film-Actor Edges

Film metadata: Title, year, genre, production company, country, user-contributed keywords, etc.

Actor metadata: Name, Age, Gender, Nationality, etc.

User-contributed keywords ($93,448$)

# METADATA IS OFTEN VERY HETEROGENEOUS
## EXAMPLE: IMDB FILM-ACTOR NETWORK

| Keyword | Occurrences |
| --- | --- |
| 'independent-film' | 15513 |
| 'based-on-novel' | 12303 |
| 'character-name-in-title' | 11801 |
| 'murder' | 11184 |
| 'sex' | 9759 |
| 'female-nudity' | 9239 |
| 'nudity' | 5846 |
| 'death' | 5791 |
| 'husband-wife-relationship' | 5568 |
| 'love' | 5560 |
| 'violence' | 5480 |
| 'police' | 5463 |
| 'father-son-relationship' | 5063 |

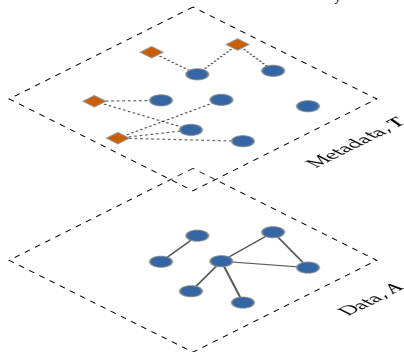# METADATA IS OFTEN VERY HETEROGENEOUS

## EXAMPLE: IMDB FILM-ACTOR NETWORK

| Keyword | Occurrences | Keyword | Occurrences |
| --- | --- | --- | --- |
| 'independent-film' | 15513 | 'discriminaton-against-anteaters' | 1 |
| 'based-on-novel' | 12303 | 'partisan-violence' | 1 |
| 'character-name-in-title' | 11801 | 'deliberately-leaving-something-behind' | 1 |
| 'murder' | 11184 | 'princess-from-outer-space' | 1 |
| 'sex' | 9759 | 'reference-to-aleksei-vorobyov' | 1 |
| 'female-nudity' | 9239 | 'dead-body-on-the-beach' | 1 |
| 'nudity' | 5846 | 'liver-failure' | 1 |
| 'death' | 5791 | 'hit-with-a-skateboard' | 1 |
| 'husband-wife-relationship' | 5568 | 'helping-blind-man-cross-street' | 1 |
| 'love' | 5560 | 'abandoned-pet' | 1 |
| 'violence' | 5480 | 'retired-clown' | 1 |
| 'police' | 5463 | 'resentment-toward-stepson' | 1 |
| 'father-son-relationship' | 5063 | 'mutilating-a-plant' | 1 |

# BETTER APPROACH: METADATA AS DATA

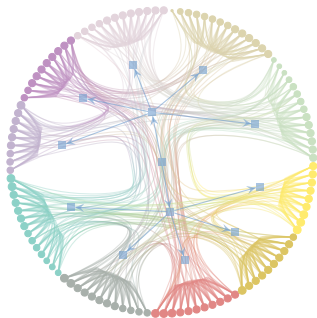Main idea: Treat metadata as data, not "ground truth".

<u>Generalized annotations</u>

$$A_{ij} \to \text{Data layer}$$
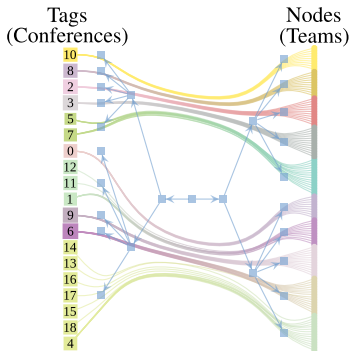$$T_{ij} \to \text{Annotation layer}$$



- ▶ Joint model for data and metadata (the layered SBM [1]).
- ▶ Arbitrary types of annotation.
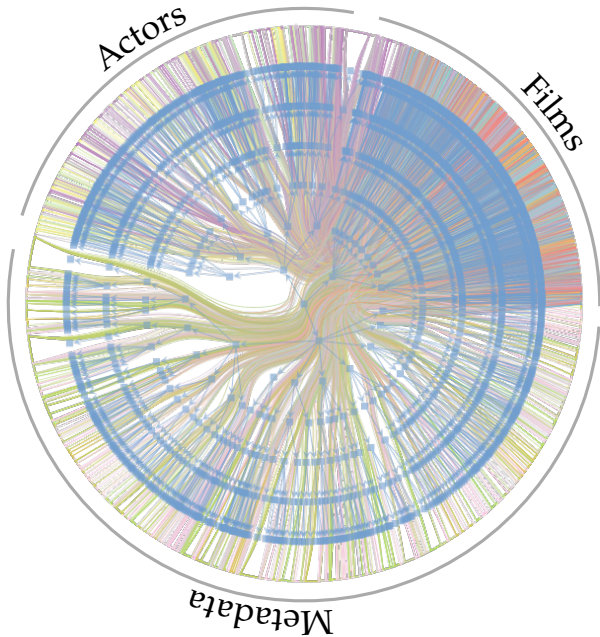- ▶ Both data and metadata are clustered into groups.
- ▶ Fully nonparametric.

[1] T.P.P, Phys. Rev. E 92, 042807 (2015)

# Example: American college football



(a) Data

(b) Metadata

# EXAMPLE: IMDB FILM-ACTOR NETWORK

# PREDICTION OF MISSING EDGES

Drug-drug interactions



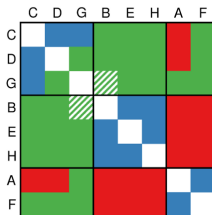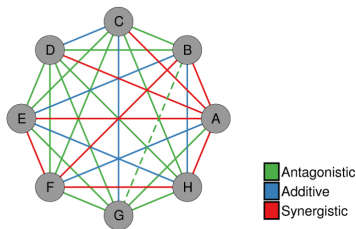$$G' = \underbrace{G}_{\text{Observed}} \cup \underbrace{\delta G}_{\text{Missing}}$$

Posterior probability of missing edges

$$P(\delta G | G, \{b_i\}) = \frac{\sum_\theta P(G \cup \delta G | \{b_i\}, \theta) P(\theta)}{\sum_\theta P(G | \{b_i\}, \theta) P(\theta)}$$

A. Clauset, C. Moore, MEJ Newman, Nature,
2008
R. Guimerà, M Sales-Pardo, PNAS 2009

Antagonistic
Additive
Synergistic

R. Guimerà, M. Sales-Pardo, PLoS Comput
Biol, 2013

# METADATA AND PREDICTION OF *missing nodes*

Node probability, with known group membership:

$$P(\boldsymbol{a}_i|\boldsymbol{A}, b_i, \boldsymbol{b}) = \frac{\sum_\theta P(\boldsymbol{A}, \boldsymbol{a}_i|b_i, \boldsymbol{b}, \theta)P(\theta)}{\sum_\theta P(\boldsymbol{A}|\boldsymbol{b}, \theta)P(\theta)}$$

Node probability, with unknown group membership:

$$P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{b}) = \sum_{b_i} P(\boldsymbol{a}_i|\boldsymbol{A}, b_i, \boldsymbol{b})P(b_i|\boldsymbol{b}),$$

Node probability, with unknown group membership, but known metadata:

$$P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c}) = \sum_{b_i} P(\boldsymbol{a}_i|\boldsymbol{A}, b_i, \boldsymbol{b})P(b_i|\boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c}),$$

Group membership probability, given metadata:

$$P(b_i|\boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c}) = \frac{P(b_i, \boldsymbol{b}|\boldsymbol{T}, \boldsymbol{c})}{P(\boldsymbol{b}|\boldsymbol{T}, \boldsymbol{c})} = \frac{\sum_\gamma P(\boldsymbol{T}|b_i, \boldsymbol{b}, \boldsymbol{c}, \gamma)P(b_i, \boldsymbol{b})P(\gamma)}{\sum_{b'_i} \sum_\gamma P(\boldsymbol{T}|b'_i, \boldsymbol{b}, \boldsymbol{c}, \gamma)P(b'_i, \boldsymbol{b})P(\gamma)}$$

Predictive likelihood ratio:

$$\lambda_i = \frac{P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c})}{P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c}) + P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{b})} \qquad \begin{array}{l} \lambda_i > 1/2 \rightarrow \text{the metadata improves} \\ \text{the prediction task} \end{array}$$

# METADATA AND PREDICTION OF MISSING NODES



(a) Data  (b) Metadata  (c) Node prediction

$$\lambda_i = \frac{P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c})}{P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{T}, \boldsymbol{b}, \boldsymbol{c}) + P(\boldsymbol{a}_i|\boldsymbol{A}, \boldsymbol{b})}$$

# METADATA AND PREDICTION OF MISSING NODES



Aligned          Misaligned          Random

# METADATA PREDICTIVENESS

Neighbor probability:

$$P_e(i|j) = k_i \frac{e_{b_i,b_j}}{e_{b_i}e_{b_j}}$$

Neighbour probability, given metadata tag:

$$P_t(i) = \sum_j P(i|j)P_m(j|t)$$

Null neighbor probability (no metadata tag):

$$Q(i) = \sum_j P(i|j)\Pi(j)$$

Kullback-Leibler divergence:
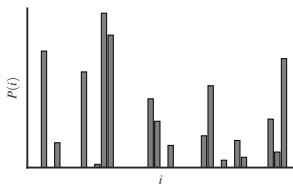
$$D_{KL}(P_t||Q) = \sum_i P_t(i) \ln \frac{P_t(i)}{Q(i)}$$

Relative divergence:

$$\mu_r \equiv \frac{D_{KL}(P_t||Q)}{H(Q)} \to \text{Metadata group predictiveness}$$

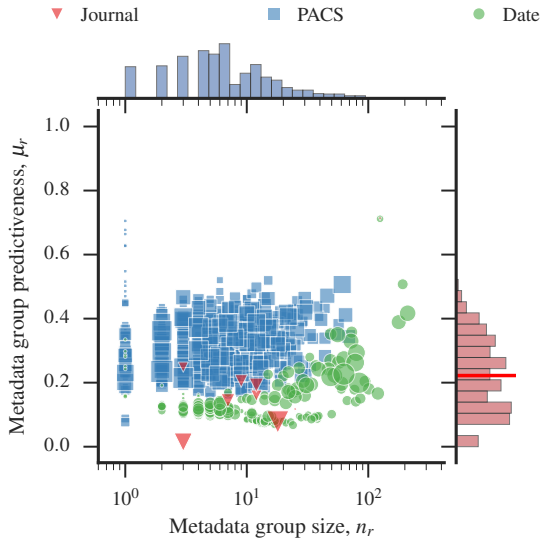Neighbour prob. without metadata



Neighbour prob. with metadata

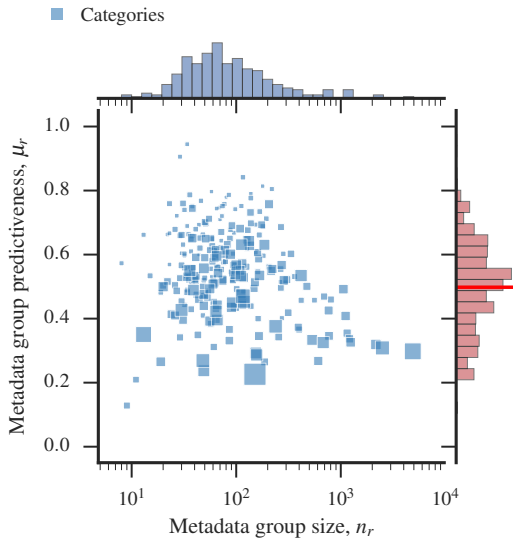# METADATA PREDICTIVENESS

## IMDB FILM-ACTOR NETWORK

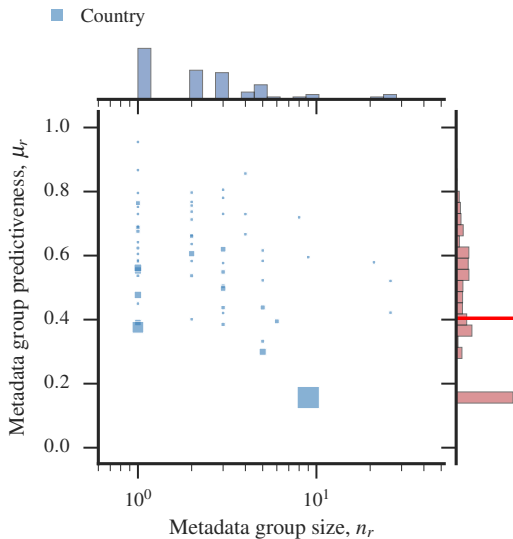# METADATA PREDICTIVENESS

## APS CITATION NETWORK
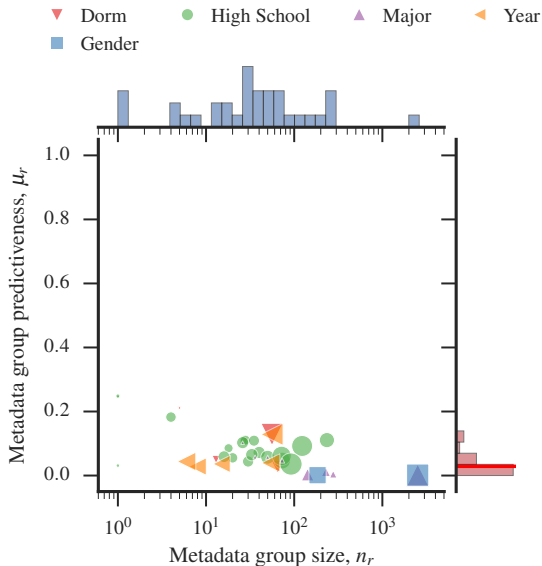
# METADATA PREDICTIVENESS

AMAZON CO-PURCHASES

# METADATA PREDICTIVENESS
INTERNET AS

# METADATA PREDICTIVENESS

## FACEBOOK PENN STATE

# THE END

Main Message:

- ► Metadata is often structured, heterogeneous and noisy.
- ► It is in general not trivially descriptive of network structure
  ($\neq$ "ground truth").
- ► It should be treated as part of the data, and modeled.

> Darko Hric, T. P. P., Santo Fortunato, arXiv:1604.00255

Other talks:
**"The Trouble with Community Detection"**
M. E. J. Newman and Aaron Clauset
Wed. 14:00, Dongkang B, 3F

**"The Ground Truth about Metadata and Community Detection in Networks"**
Leto Peel, Daniel B. Larremore and Aaron Clauset
Wed. 15:00, Dongkang B, 3F

 graph-tool

> Very fast, freely available C++ code as part of the
> graph-tool Python library.
> `http://graph-tool.skewed.de`

# EFFICIENT INFERENCE ALGORITHMS

### Smart MCMC

▶ Choose a random vertex $v$ (happens to belong to block $r$).

▶ Move it to a random block $s \in [1, B]$, chosen with a probability $p(r \rightarrow s|t)$ proportional to $e_{ts} + \epsilon$, where $t$ is the block membership of a randomly chosen neighbour of $v$.

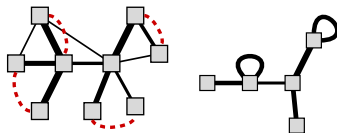▶ Accept the move with probability

$$a = \min\left\{ e^{-\beta\Delta\mathcal{S}} \frac{\sum_t p_t^i p(s \rightarrow r|t)}{\sum_t p_t^i p(r \rightarrow s|t)}, 1 \right\}.$$

▶ Repeat.



Fast mixing times.

### Agglomerative initialization



Avoids metastable states.

Algorithmic complexity:

$O(N \ln^2 N)$
(independent of $B$)

Scales up to $10^7 - 10^8$ edges.

✿ graph-tool

Freely available efficient implementation
http://graph-tool.skewed.de