

A COMPARAISON OF COMMUNITY DETECTION ALGORITHMS ON ARTIFICIAL NETWORKS

Günce Keziban Orman

*Galatasaray University C.S. Department
TUBITAK UEKAE*

keziban.orman@uekae.tubitak.gov.tr

Vincent Labatut

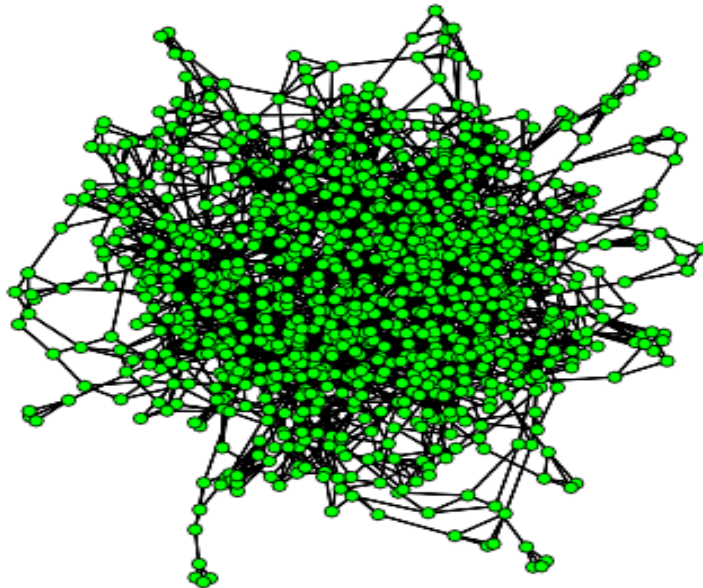
Galatasaray University C.S Department

vlabatut@gsu.edu.tr

1. Introduction
 1. Complex networks
 2. Community detection
 3. Testing algorithms
2. Method
 1. Selected algorithms
 2. Lancichinetti *et. al* generative model
 3. Normalized mutual information
3. Results & Discussions
4. Conclusion

1.1 COMPLEX NETWORKS

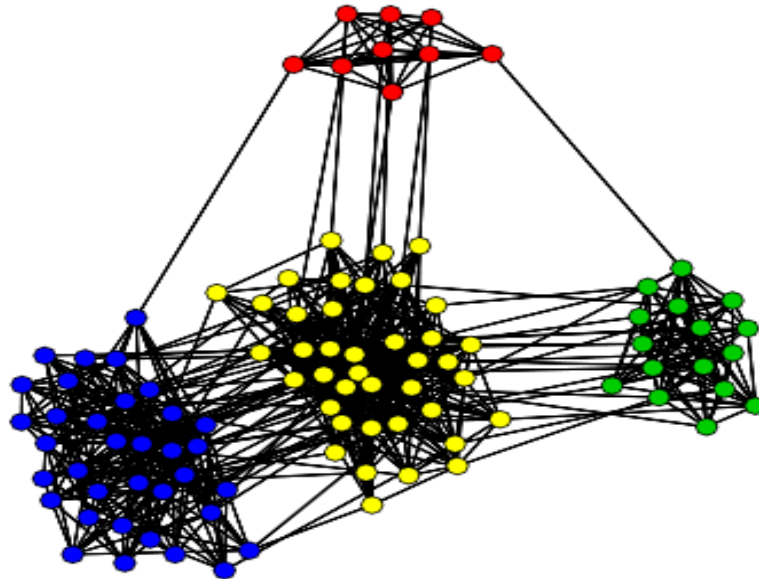
- Large graphs with non-trivial topological features
- Model systems of interacting objects.
- ex: Internet, www, protein web, friendship networks ...



1.2.1 COMMUNITY DETECTION

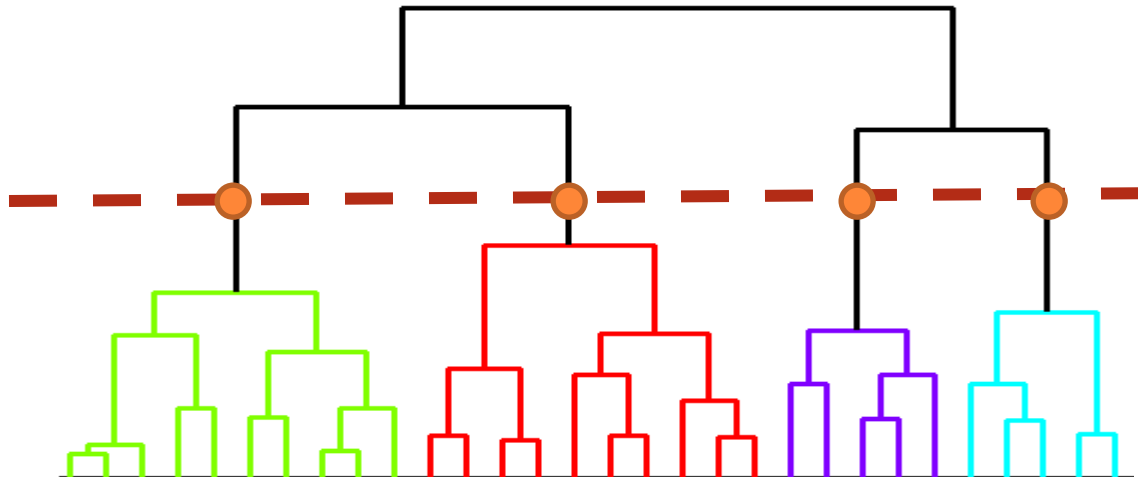
DEFINITION

- Community: group of nodes with dense inner links and sparse outer links
- Community detection: find the best graph partition according to this definition



1.2.2 COMMUNITY DETECTION SOLUTION APPROACHES

- Hierarchical Approaches
 - Divisive vs. Agglomerative
- Optimization Approaches
 1. Partition: stochastic vs heuristic
 2. Quality evaluation
- Others



1.2.3 COMMUNITY DETECTION - MODULARITY

- Newman's modularity measure:

$$Q = \sum_i (e_{ii} - a_i^2)$$

- e_{ii} : observed fraction of links inside the i^{th} community
 - a_i : estimation of e_{ii} under the hypothesis of uniformly randomly distributed links.
- Values:
 - $Q=0$: networks without community structure and/or random partition
 - $Q \approx 1$: network with strong community structure and good partition
 - $Q \in [0.3, 0.7]$ is generally considered a good result

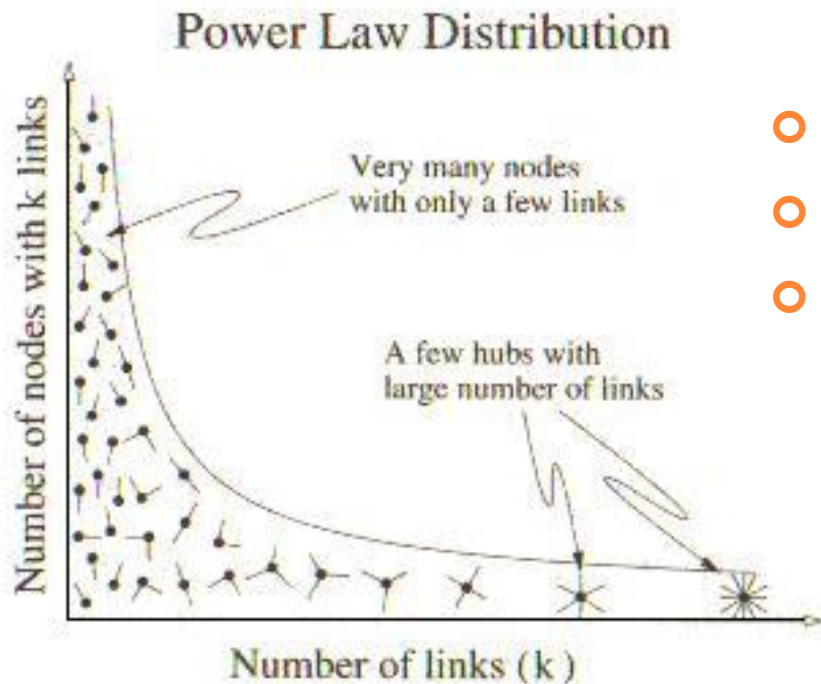
1.3 TESTING ALGORITHMS

	Real networks	Artificial Networks
Size	Usually small	No theoretical limit
Network properties	Uncontrollable, depends on the modeled system	depend on the generation model parameters
Construction	Expensive and/or difficult to build	Computer generated
Community structure	Possibly subjective or unknown community structure	Communities created and controlled

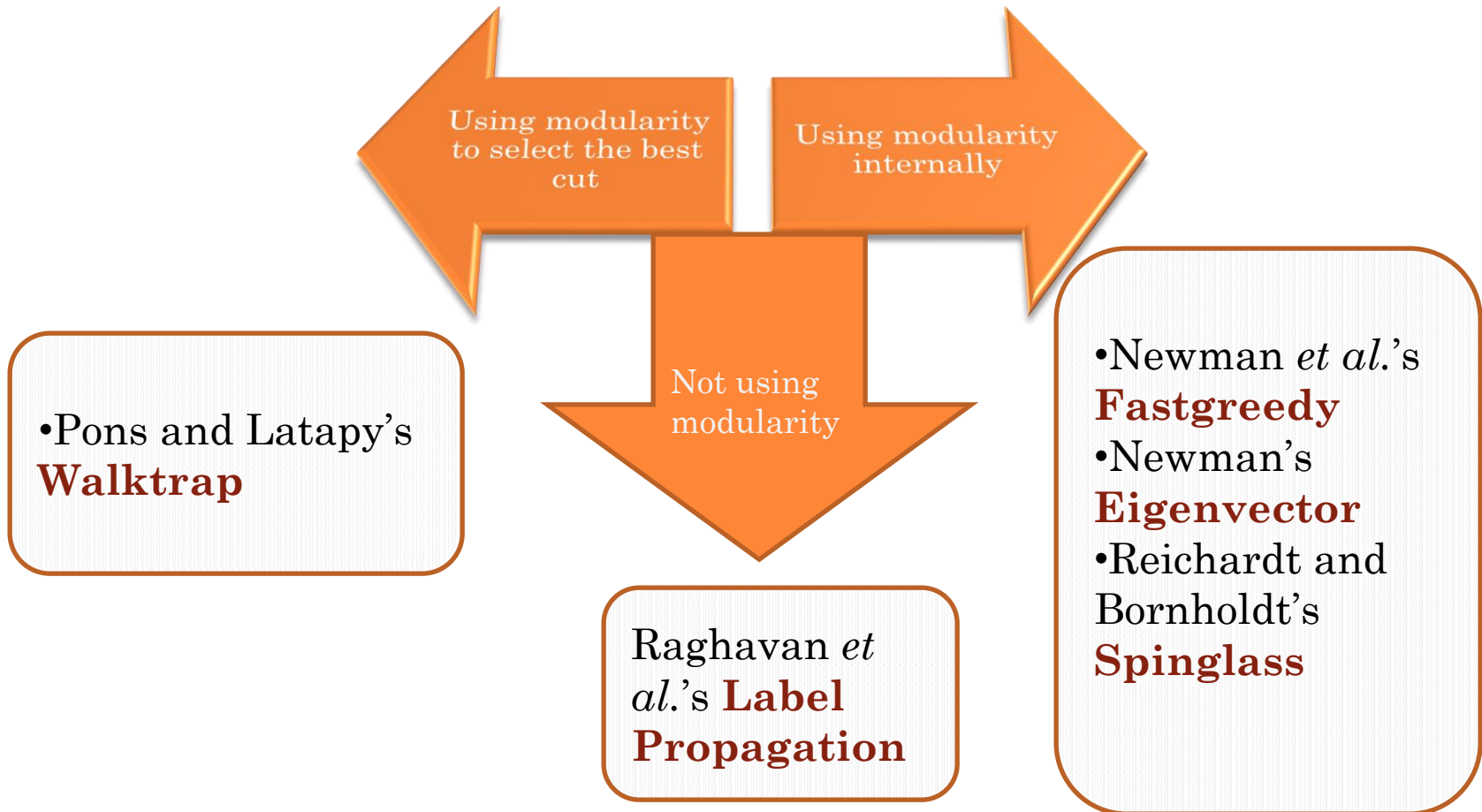
1.3 TESTING ALGORITHMS

Real Network Properties:

- Scale-free (power-law degree distribution)
- Power-law community size distribution
- Small average distance
- High transitivity
- High degree correlation



2.1 SELECTED ALGORITHMS



2.2 LANCICHINETTI *ET. AL* GENERATIVE MODEL

$n, \beta, \gamma, \langle k \rangle, k_{max}, \mu$

● Apply configuration model with average degree $\langle k \rangle$, max. degree k_{max} and power-law exponent γ

● Draw community sizes with power-law exponent β and affect each node to a community

● Rewire some links in order to respect μ , without changing the nodes degrees

2.2 LANCICHINETTI *ET. AL* GENERATIVE MODEL

PARAMETER	VALUE
n	{1000}
β	{1,2}
γ	{2,3}
$\langle k \rangle$	{5,15,30}
k_{max}	{15,45,90}
μ	[0.05,0.95]

2.3 NORMALIZED MUTUAL INFORMATION (NMI)

- Quality assessment for a dataset partition
 - m : confusion matrix
 - n : dataset size

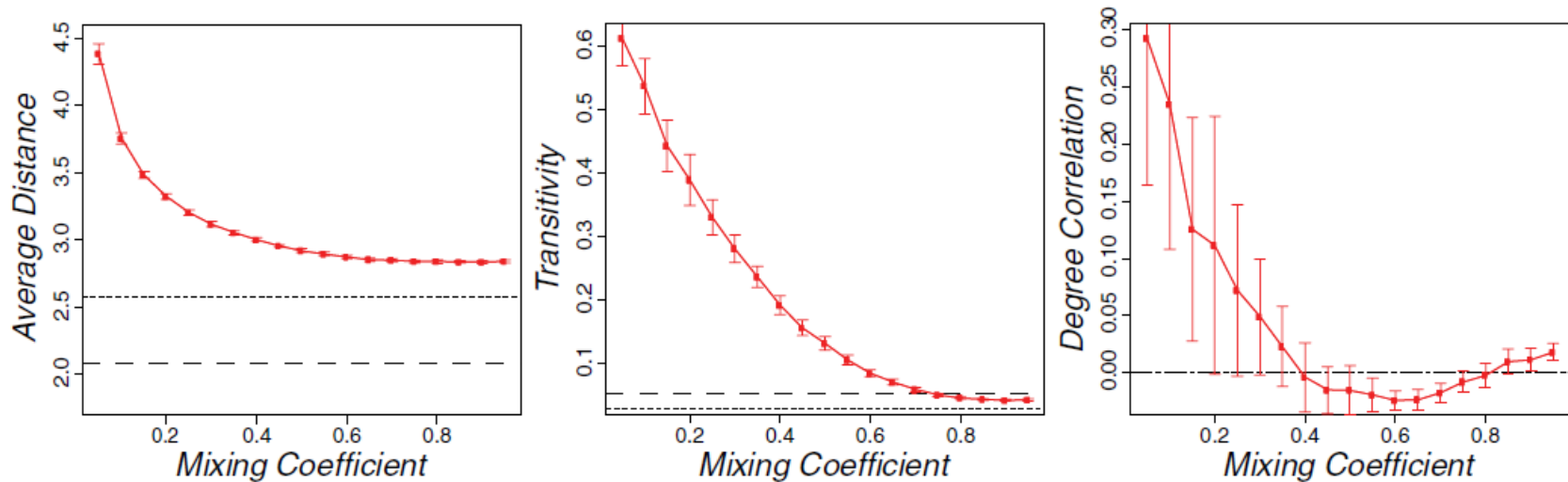
$$I = \frac{-2 \sum_i \sum_j m_{ij} \log(nm_{ij}/m_{i+}m_{+j})}{\sum_i m_{i+} \log(m_{i+}/n) + \sum_j m_{+j} \log(m_{+j}/n)}$$

- Values:
 - 0: random partition
 - 1: perfect partition

3. RESULTS & DISCUSSIONS

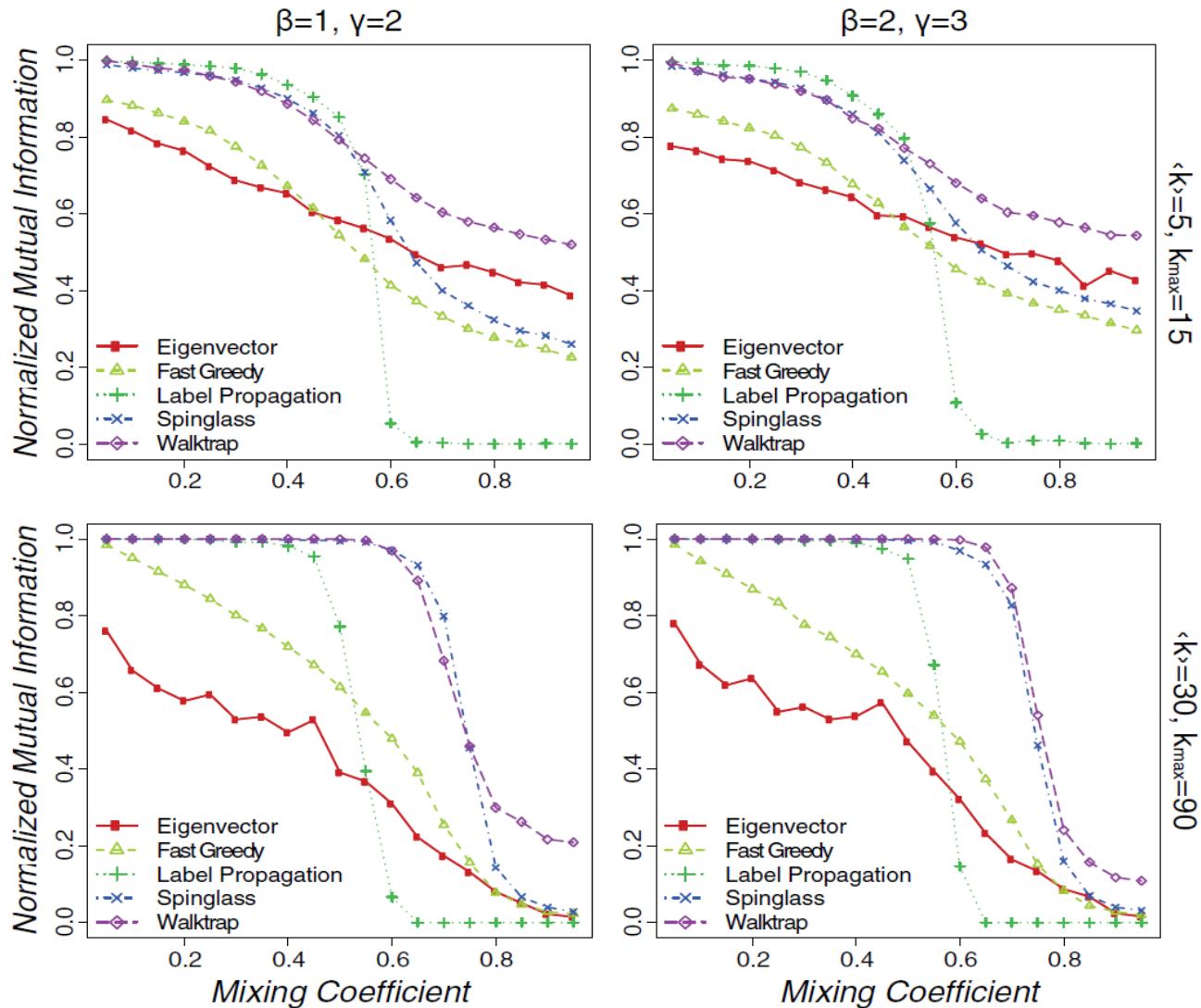
NETWORKS' UNCONTROLLED PROPERTIES

6/27/2009

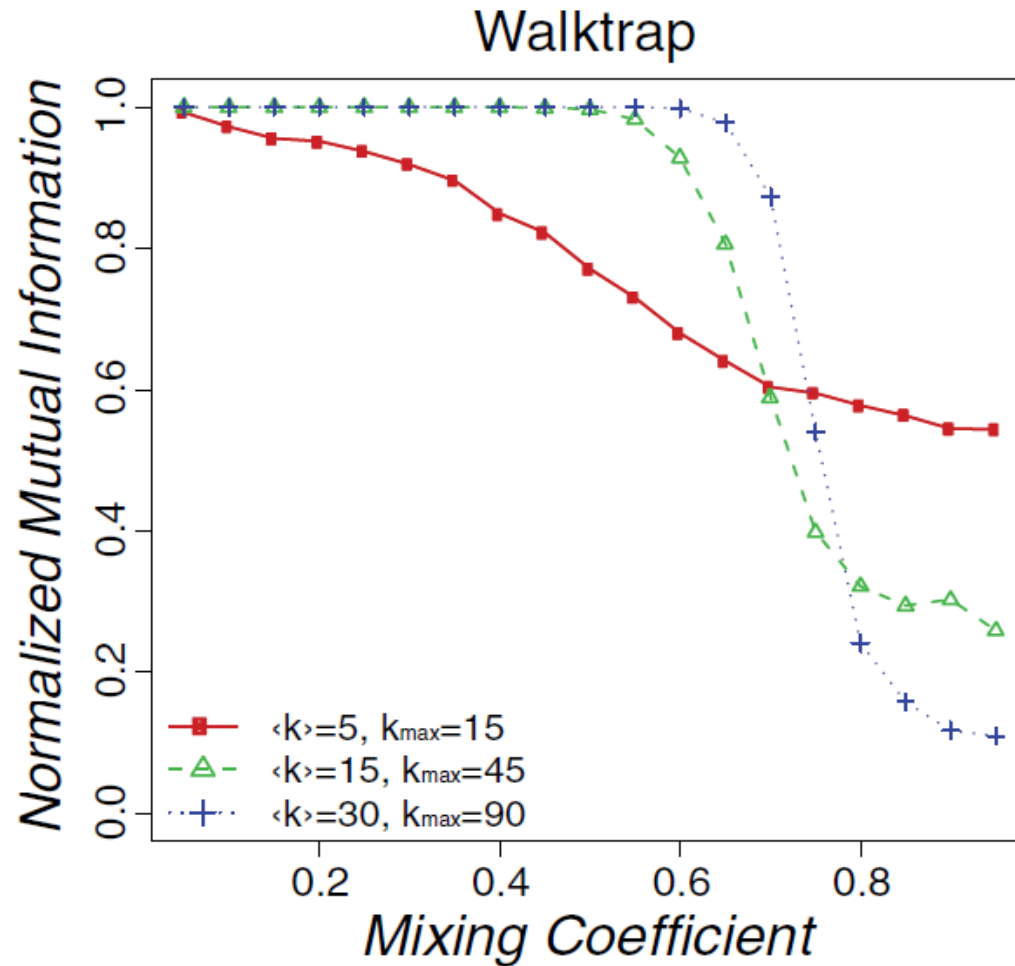


$$\langle k \rangle = 30 \quad k_{max} = 90$$

3. RESULTS & DISCUSSIONS



3. RESULTS & DISCUSSIONS



3. RESULTS & DISCUSSIONS

- Parameters effects:
 - β and γ : almost no effect ($p < 0.06$)
 - $\langle k \rangle$: higher average degree improves performance
- Algorithms:
 - Partition quality:
 - WT and SG performs better
 - LP is not robust
 - Speed:
 - SG is slow,
 - LP, FG and WT are fast
 - EV lie somewhere in between

4. CONCLUSION

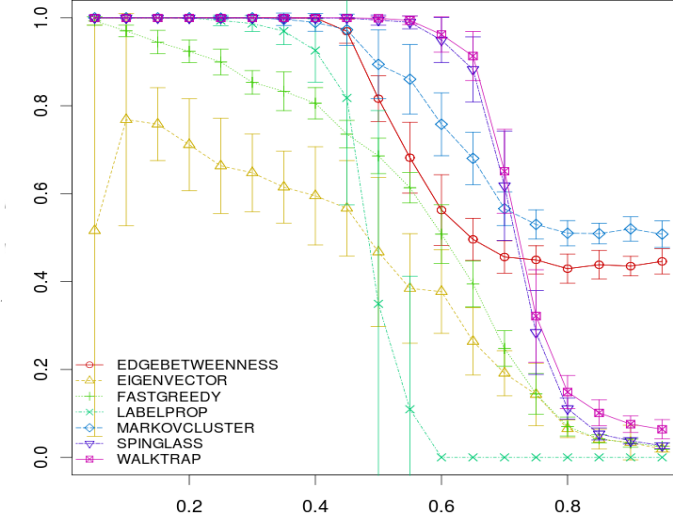
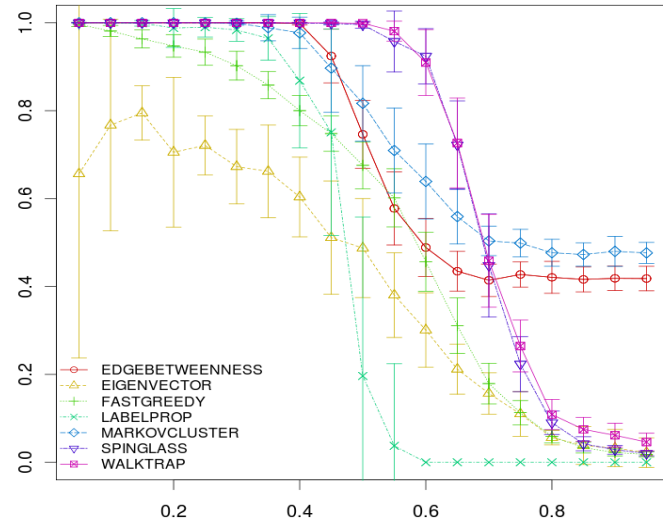
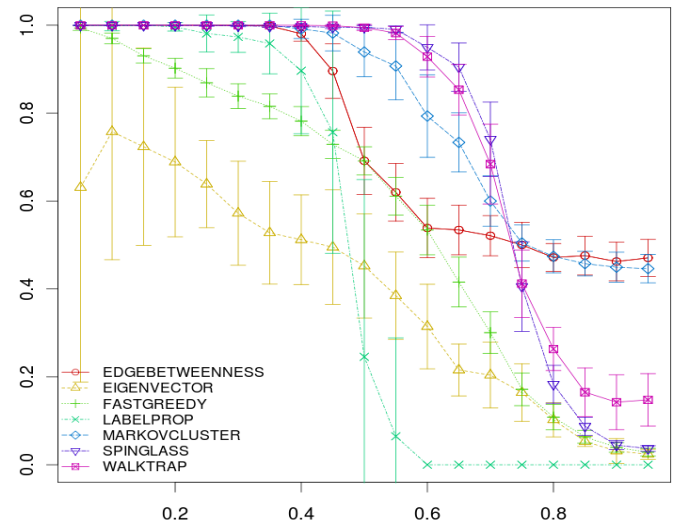
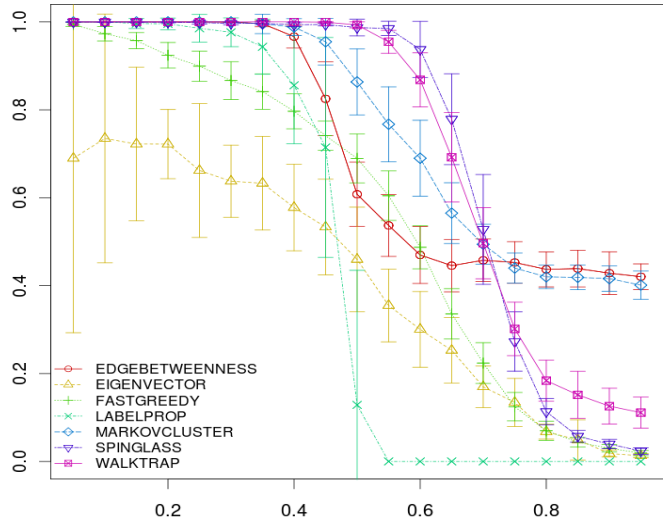
Further experiments

- on larger networks
- with more realistic networks
- with more algorithms
- with different performance measures



QUESTIONS

6/27/2009



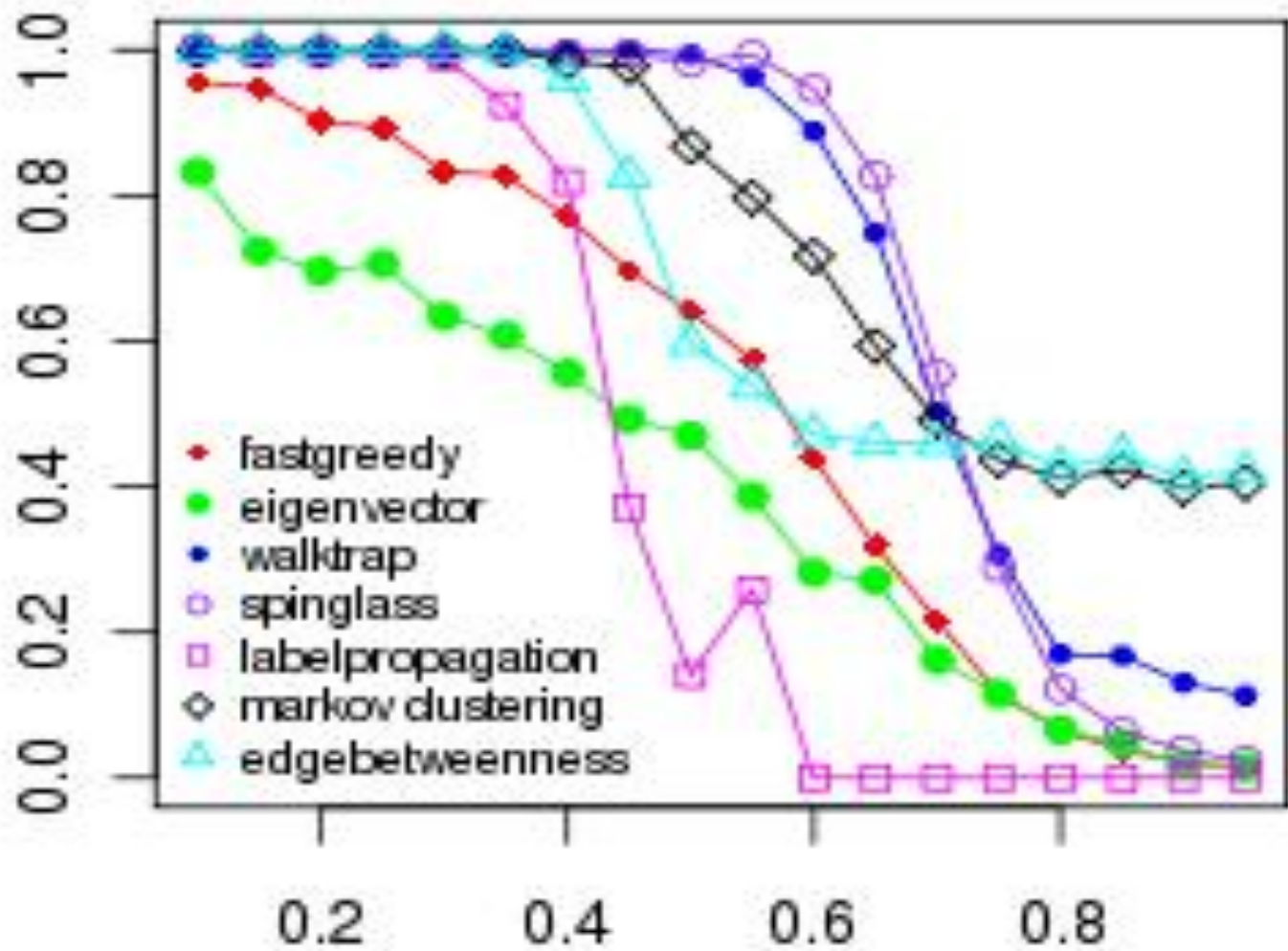
ALGORITHMS COMPLEXITIES

6/27/2009

LP	FG	EV	WT	MC	SG	EB
	$O(m+n)$	$O(n \log^2 n)$	$O(n^2)$	$O(n^2 \log n)$	$O(nk^2n)$	$O(m^2n)$

n: node number; **m:** link number ; **k:** number of ressource allocated

gamma 2 beta 1



degree 30 gamma 2 beta 1

