



NATIONAL RESEARCH
UNIVERSITY

Network communities

Network Science

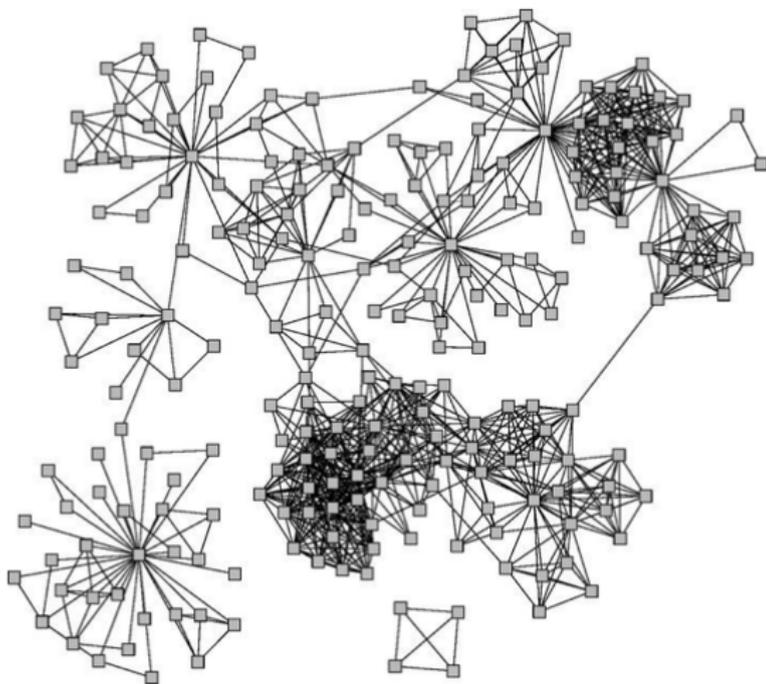
Lecture 8

Leonid Zhukov

lzhukov@hse.ru

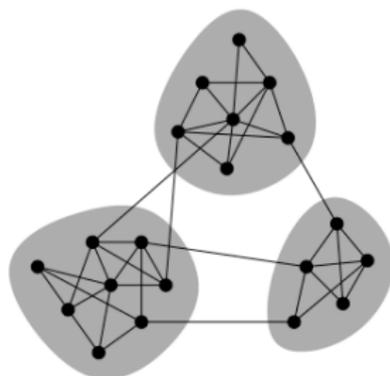
www.leonidzhukov.net/hse/2022/networks

National Research University Higher School of Economics
School of Data Analysis and Artificial Intelligence, Department of Computer Science



Definition

Network communities are groups of vertices such that vertices inside the group connected with many more edges than between groups.



- Community detection is an assignment of vertices to communities.

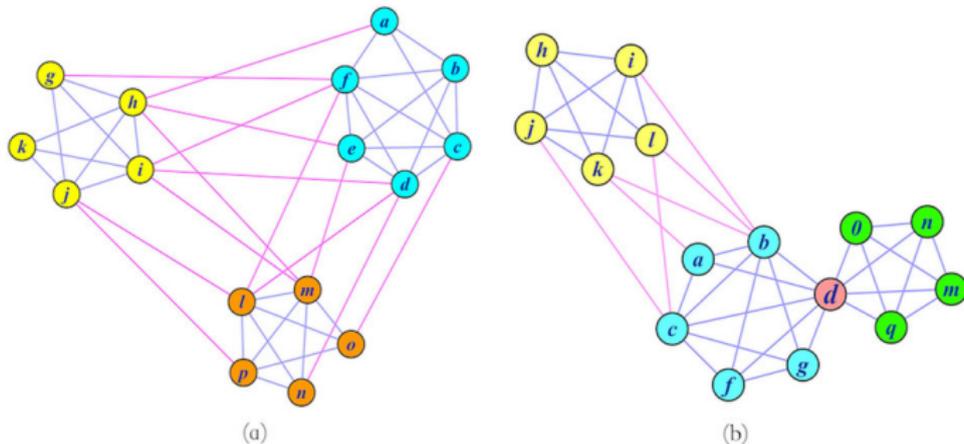


image from W. Liu, 2014

Community is a cohesive subgroup:

- Mutuality of ties. Almost everyone in the group has ties (edges) to one another
- Compactness. Closeness or reachability of group members in small number of steps, not necessarily adjacency
- Density of edges. High frequency of ties within the group
- Separation. Higher frequency of ties among group members compared to non-members

Wasserman and Faust

- Graph $G(V, E)$, $n = |V|$, $m = |E|$
- Community - set of nodes S
 n_s -number of nodes in S , m_s - number of edges in S
- Graph density

$$\rho = \frac{m}{n(n-1)/2}$$

- community internal density

$$\delta_{int} = \frac{m_s}{n_s(n_s-1)/2}$$

- external edges density

$$\delta_{ext} = \frac{m_{ext}}{n_s(n-n_s)}$$

- community (cluster): $\delta_{int} > \rho$, $\delta_{ext} < \rho$

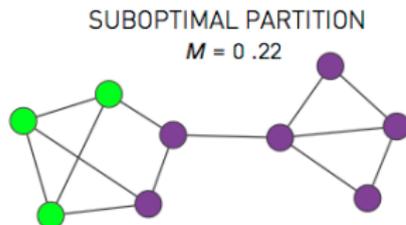
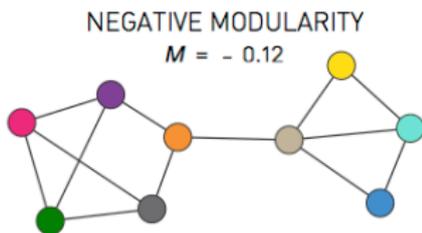
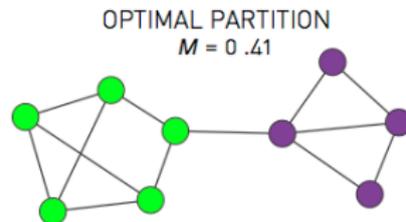
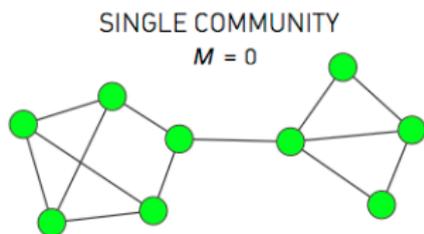
- Compare fraction of edges within the cluster to expected fraction in random graph with identical degree sequence
- Modularity score

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j), = \sum_u \left(\frac{m_u}{m} - \left(\frac{k_u}{2m} \right)^2 \right)$$

m_u - number of internal edges in a community u ,

k_u - sum of node degrees within a community

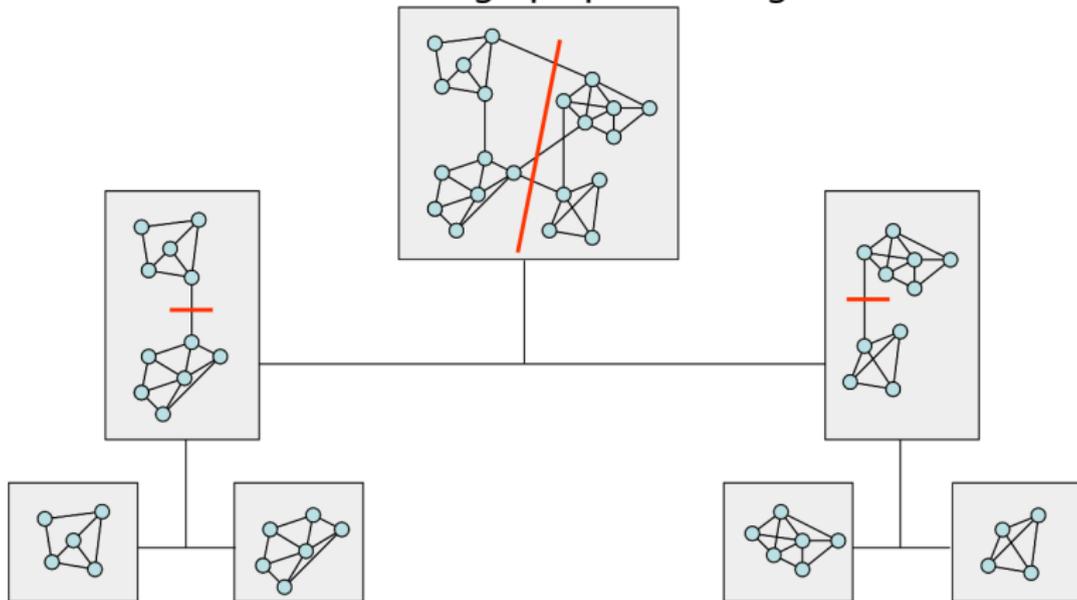
- Modularity score range $Q \in [-1/2, 1)$, single community
 $Q = 0$



- The higher the modularity score - the better are communities

from A.L. Barabasi 2016

Recursive graph partitioning



Finding optimal partition

Graph $G(E, V)$ partition: $V = V_1 + V_2$

- Graph cut

$$Q = \text{cut}(V_1, V_2) = \sum_{i \in V_1, j \in V_2} e_{ij}$$

- Ratio cut:

$$Q = \frac{\text{cut}(V_1, V_2)}{\|V_1\|} + \frac{\text{cut}(V_1, V_2)}{\|V_2\|}$$

- Normalized cut:

$$Q = \frac{\text{cut}(V_1, V_2)}{\text{Vol}(V_1)} + \frac{\text{cut}(V_1, V_2)}{\text{Vol}(V_2)}$$

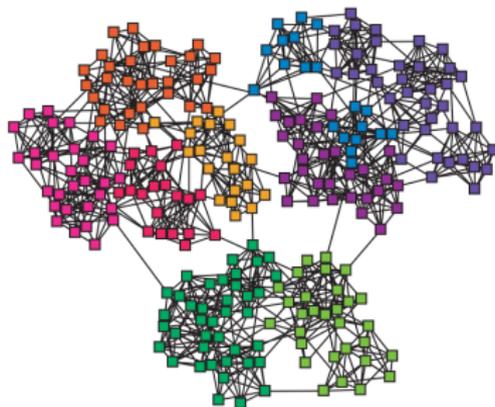
- Quotient cut (conductance):

$$Q = \frac{\text{cut}(V_1, V_2)}{\min(\text{Vol}(V_1), \text{Vol}(V_2))}$$

where: $\text{Vol}(V_1) = \sum_{i \in V_1, j \in V} e_{ij} = \sum_{i \in V_1} k_i$

Edge betweenness - number of shortest paths $\sigma_{st}(e)$ going through edge e

$$C_B(e) = \sum_{s \neq t} \frac{\sigma_{st}(e)}{\sigma_{st}}$$



Recover communities by progressively removing edges

Newman-Girvan, 2004

Algorithm: Edge Betweenness

Input: graph $G(V,E)$

Output: Dendrogram

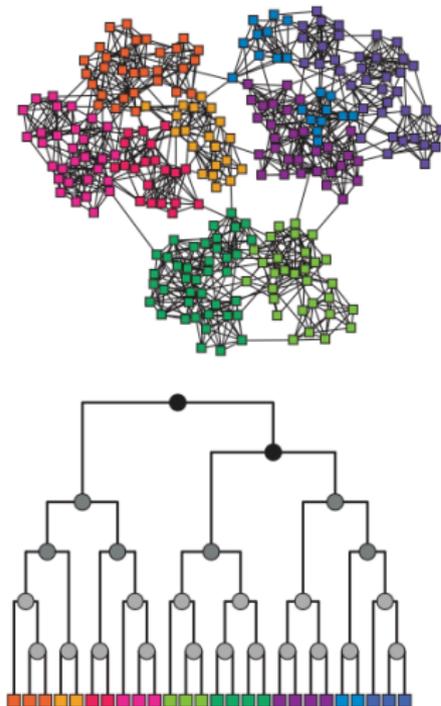
repeat

 For all $e \in E$ compute edge betweenness $C_B(e)$;
 remove edge e_i with largest $C_B(e_i)$;

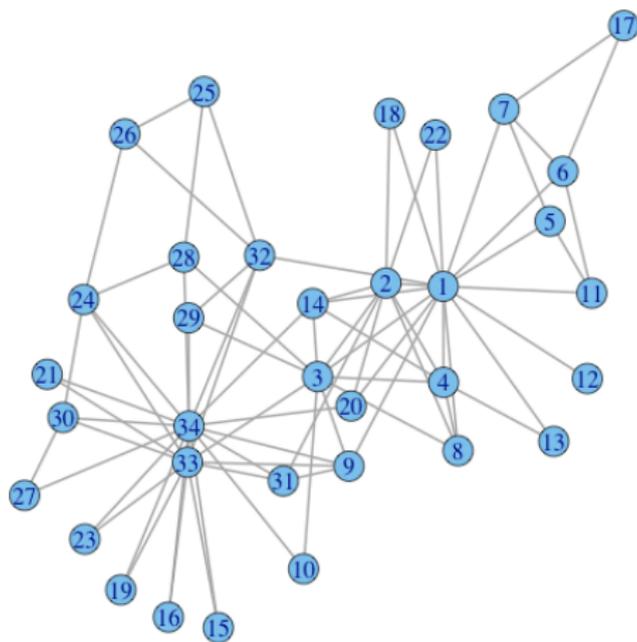
until *edges left*;

If bi-partition, then stop when graph splits in two components
(check for connectedness)

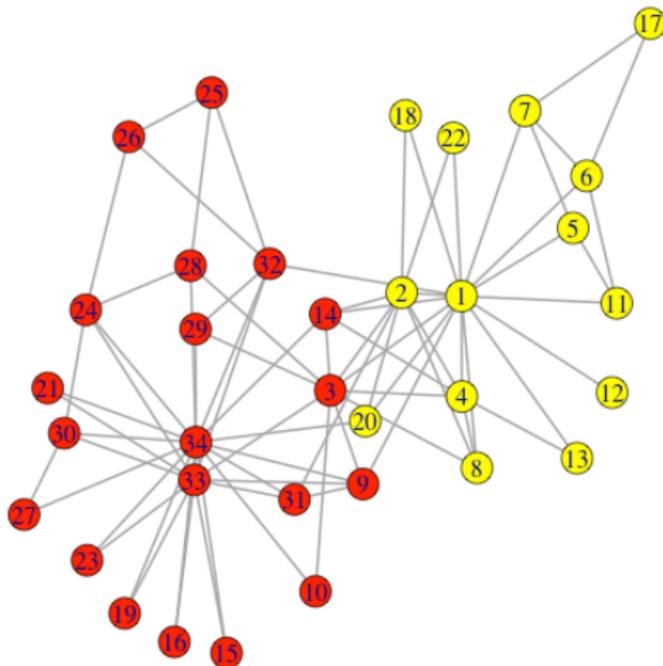
Hierarchical algorithm, dendrogram



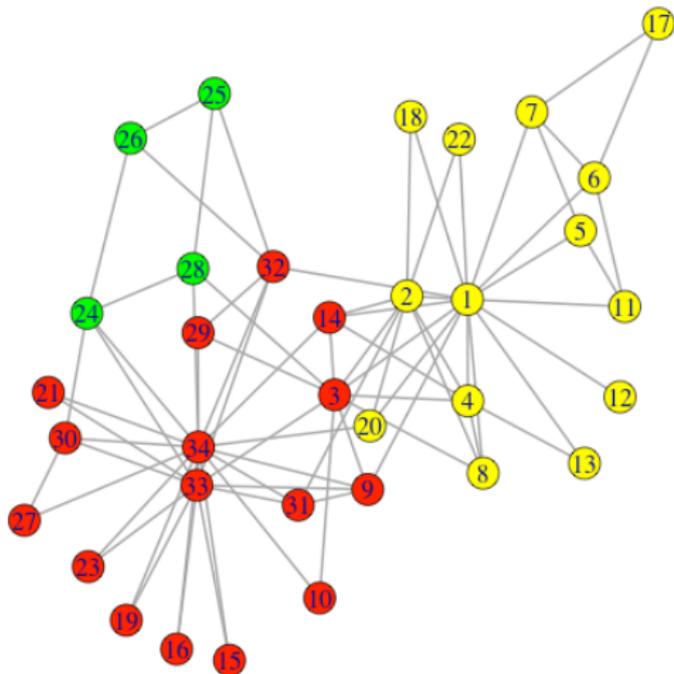
Zachary karate club

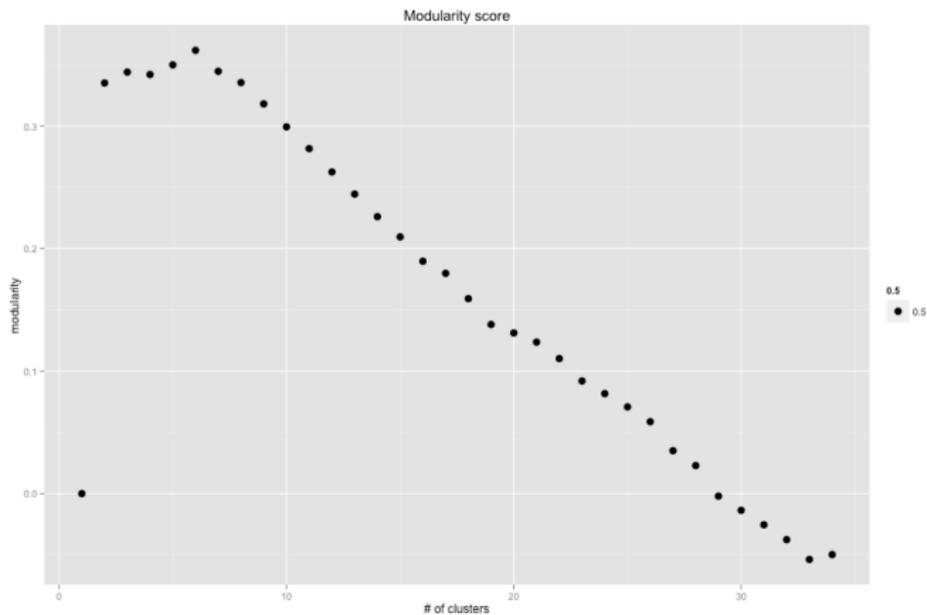


Zachary karate club

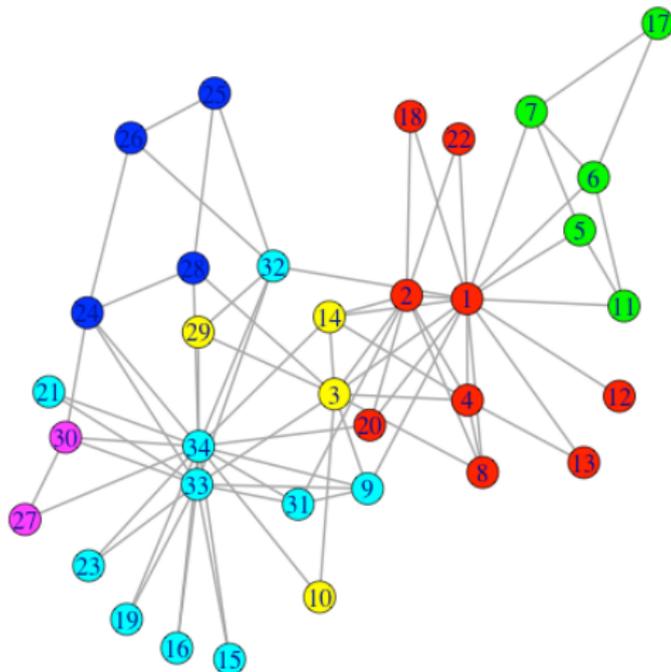


Zachary karate club

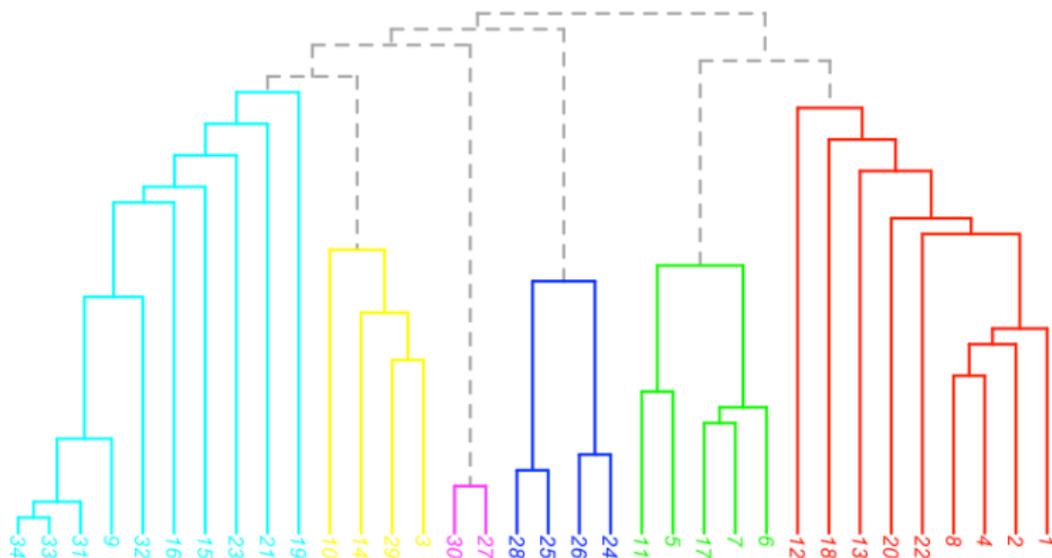




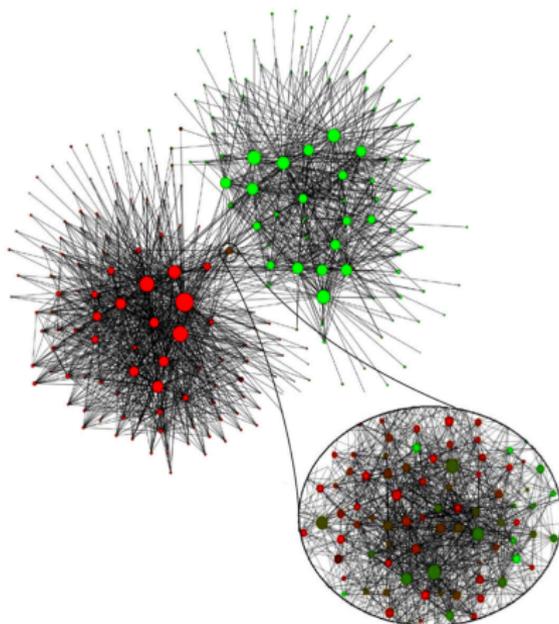
best: clusters = 6, modularity = 0.345



Zachary karate club



Multi-resolution scalable method



2 mln mobile phone network

V. Blondel et.al., 2008

“The Louvain method”

- Heuristic method for greedy modularity optimization
- Find partitions with high modularity
- Multi-level (multi-resolution) hierarchical scheme
- Scalable

Modularity:

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

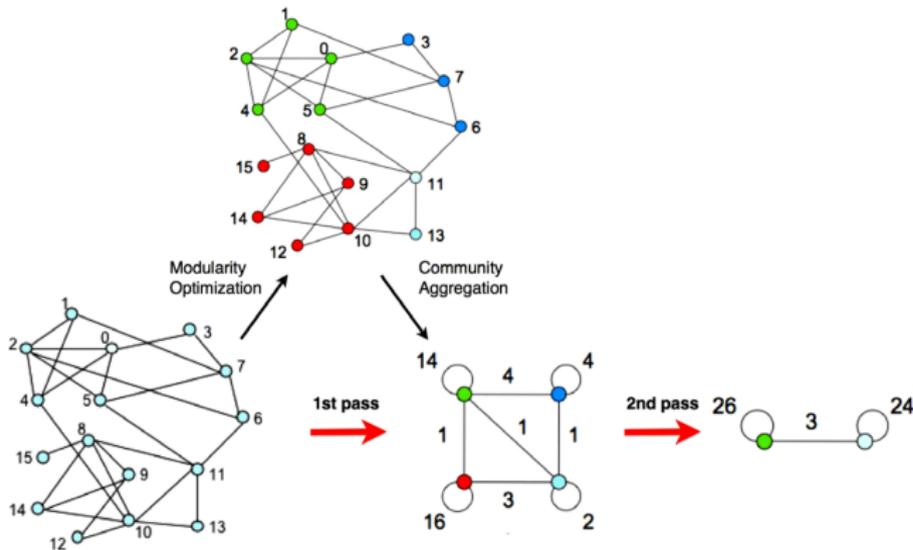
V. Blondel et.al., 2008

Algorithm

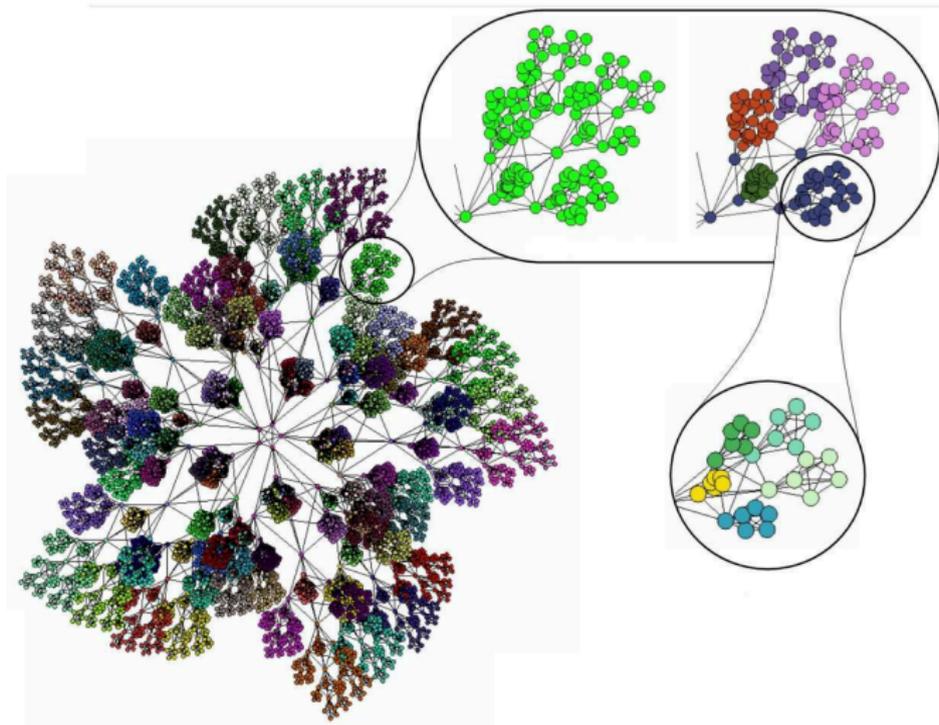
- Assign every node to its own community
- Phase I
 - For every node evaluate modularity gain from removing node from its community and placing it in the community of its neighbor
 - Place node in the community maximizing modularity gain
 - repeat until no more improvement (local max of modularity)
- Phase II
 - Nodes from communities merged into "super nodes"
 - Weight on the links added up
- Repeat until no more changes (max modularity)

V. Blondel et.al., 2008

Fast community unfolding



V. Blondel et.al., 2008

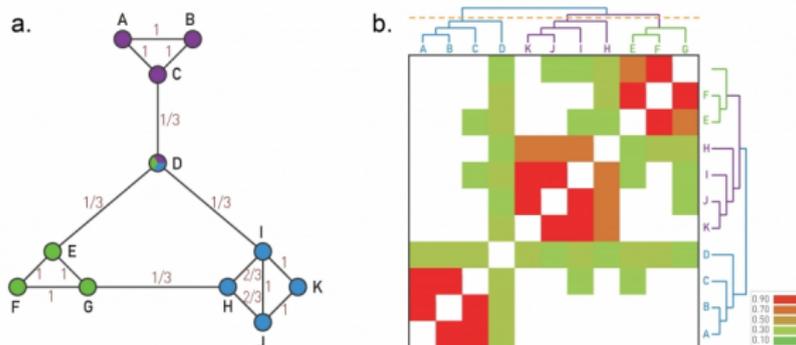


V. Blondel et al., 2008

Similarity matrix

$$x_{ij} = \frac{N(i,j)}{\min(k_i, k_j) + 1 - \theta(A_{i,j})}$$

- N_{ij} - number of common neighbors
- $\theta()$ - step function, 0 for $x \leq 0$ and 1 for $x > 0$

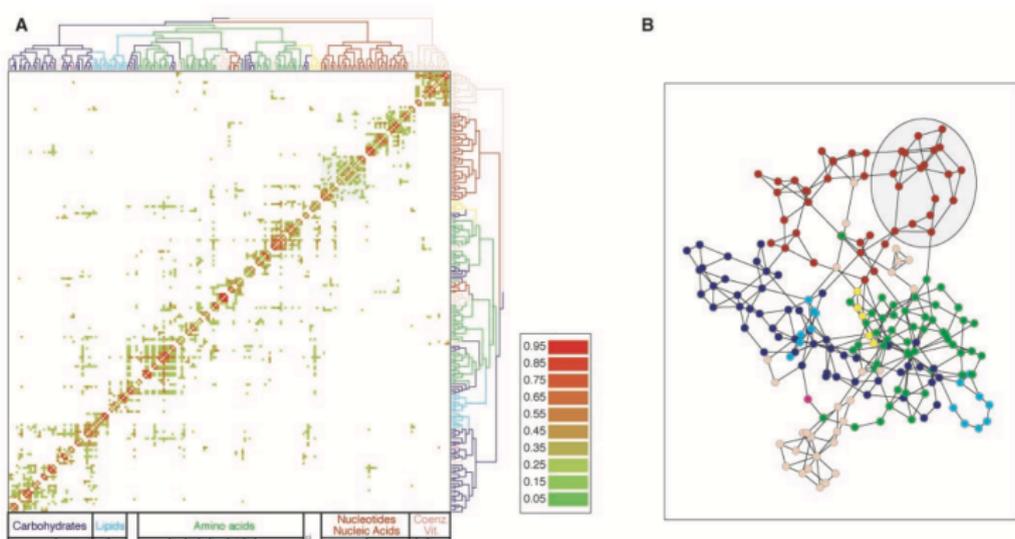




The Ravasz algorithm:

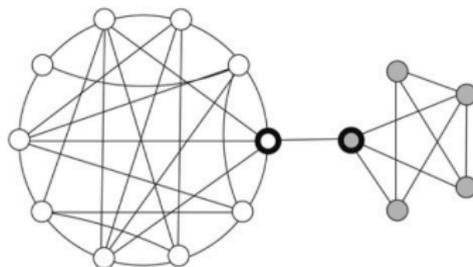
- Assign each node to a community of its own and evaluate x_{ij} for all node pairs.
- Find the node (community) pair with the highest similarity and merge them into a single community.
- Calculate the similarity between the new community and all other communities.
- Repeat Steps 2 and 3 until all nodes form a single community.
- Find the optimal cut of the dendrogram

E. Ravasz et.al., 2002



Reduced *E. coli* metabolic network

E. Ravasz et al., 2002



- Random walks on a graph tend to get trapped into densely connected parts corresponding to communities.

Walktrap

- Consider random walk on graph
- At each time step walk moves to NN uniformly at random
 $P_{ij} = \frac{A_{ij}}{d(i)}$, $P = D^{-1}A$, $D_{ii} = \text{diag}(d(i))$
- P_{ij}^t - probability to get from i to j in t steps, $t \ll t_{\text{mixing}}$
- Assumptions: for two i and j in the same community P_{ij}^t is high
- if i and j are in the same community, then $\forall k, P_{ik}^t \approx P_{jk}^t$
- Distance between nodes:

$$r_{ij}(t) = \sqrt{\sum_{k=1}^n \frac{(P_{ik}^t - P_{jk}^t)^2}{d(k)}} = \|D^{-1/2}P_i^t - D^{-1/2}P_j^t\|$$

Computing node distance r_{ij}

- Direct (exact) computation: $P_{ij}^t = (P^t)_{ij}$ or $P_i^t = P^t p_i^0, p_i^0(k) = \delta_{ik}$
- Approximate computation (simulation):
 - Compute K random walks of length t starting from node i
 - Approximate $P_{ik}^t \approx \frac{N_{ik}}{K}$, number of walks end up on k

Distance between communities:

$$P_{C_j}^t = \frac{1}{|C|} \sum_{i \in C} P_{ij}^t$$

$$r_{C_1 C_2}(t) = \sqrt{\sum_{k=1}^n \frac{(P_{C_1 k}^t - P_{C_2 k}^t)^2}{d(k)}} = \|D^{-1/2} P_{C_1}^t - D^{-1/2} P_{C_2}^t\|$$

Algorithm (hierarchical clustering)

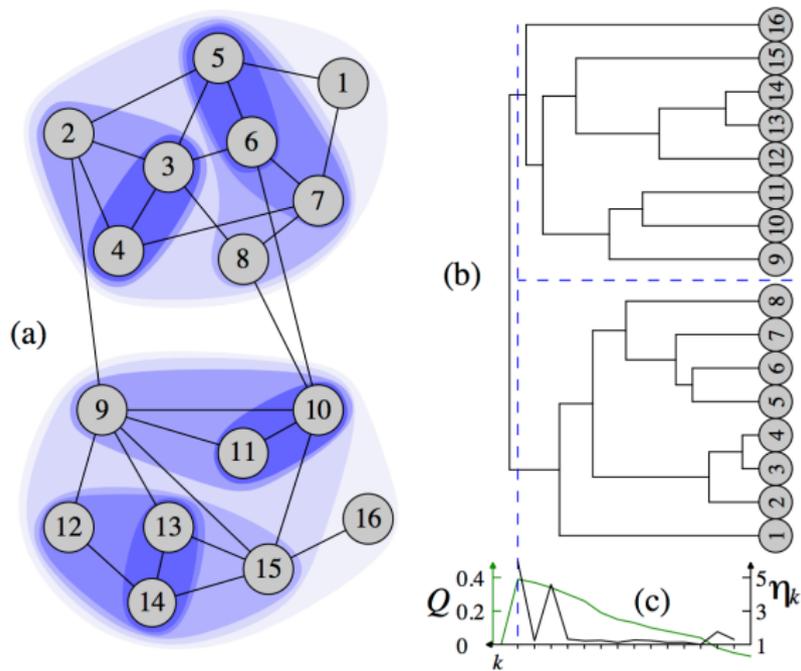
- Assign each vertex to its own community $S_1 = \{\{v\}, v \in V\}$
- Compute distance between all adjacent communities $r_{C_i C_j}$
- Choose two "closest" communities that minimizes (Ward's methods):

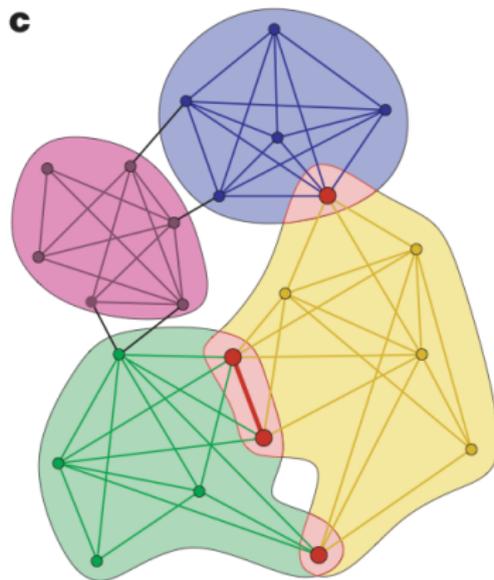
$$\Delta\sigma(C_1, C_2) = \frac{1}{n} \left(\sum_{i \in C_3} r_{iC_3}^2 - \sum_{i \in C_1} r_{iC_1}^2 - \sum_{i \in C_2} r_{iC_2}^2 \right)$$

and merge them $S_{k+1} = (S_k \setminus \{C_1, C_2\}) \cup C_3, C_3 = C_1 \cup C_2$

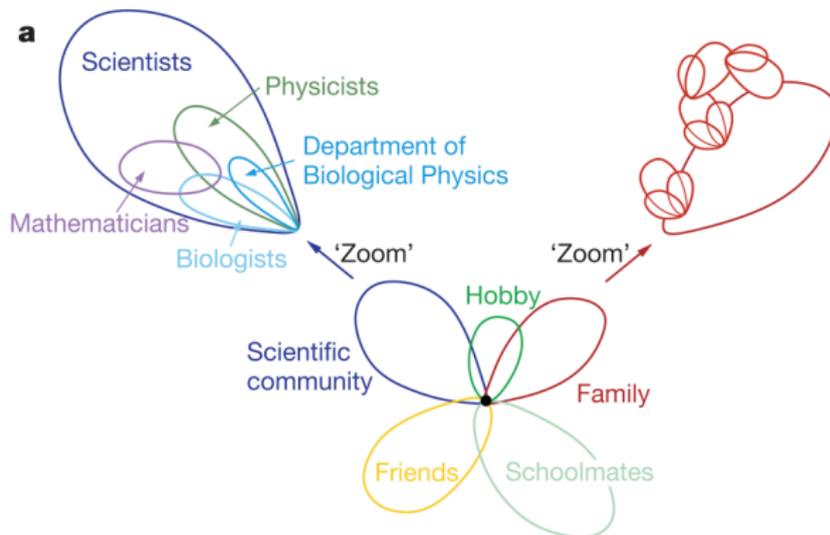
- update distance between communities

After $n - 1$ steps finish with one community $S_n = \{V\}$



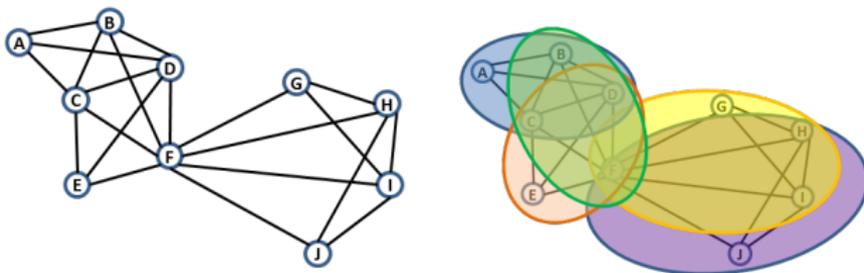


Palla, 2005

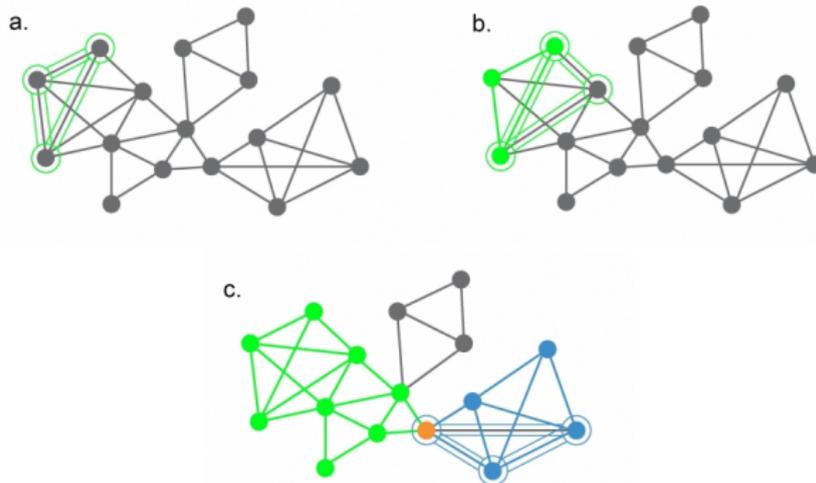


Palla, 2005

- k -clique is a clique (complete subgraph) with k nodes
- k -clique community a union of all k -cliques that can be reached from each other through a series of adjacent k -cliques
- two k -cliques are said to be adjacent if they share $k - 1$ nodes.



Adjacent 4-cliques

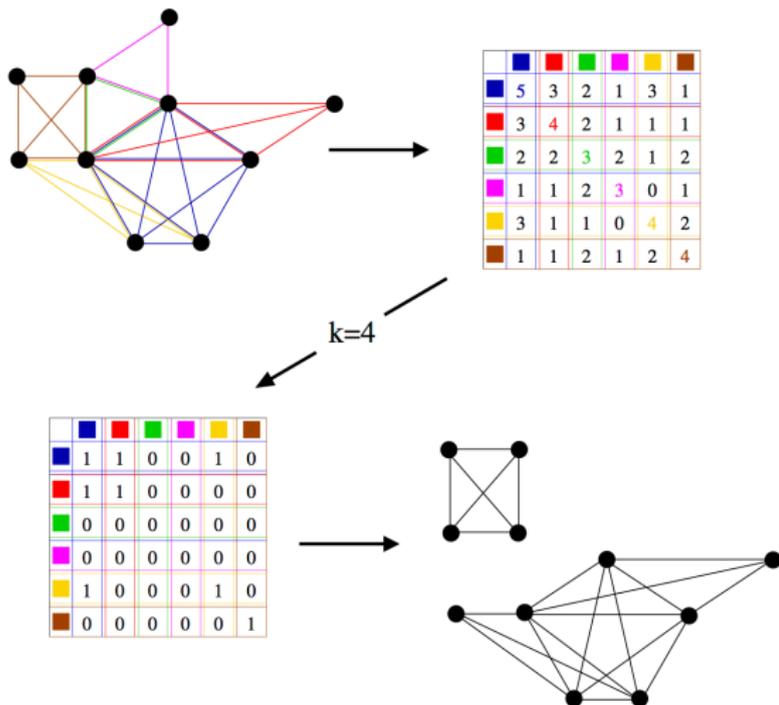


Palla, 2005

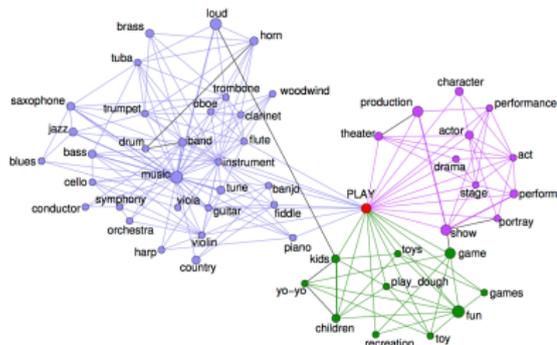


- Find all maximal cliques
- Create clique overlap matrix
- Threshold matrix at value $k - 1$
- Communities = connected components

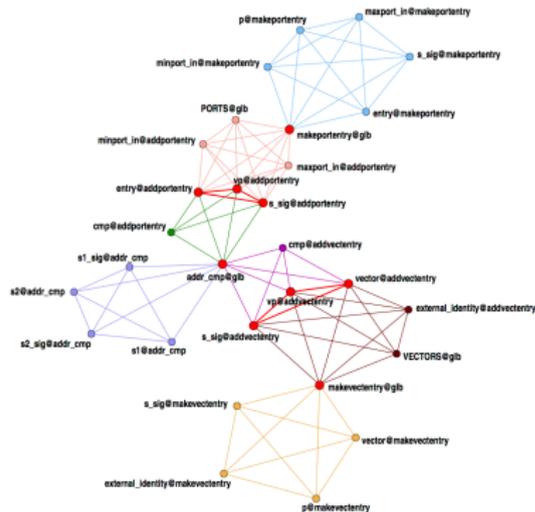
Palla, 2005



k-clique percolation

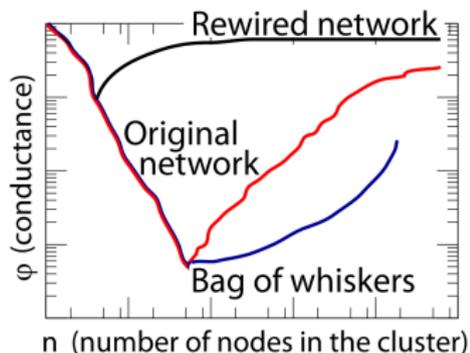


$k = 4$

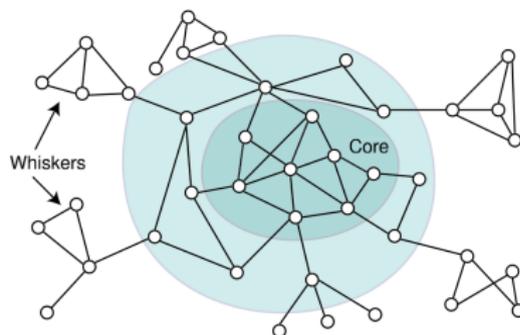


$k = 5$

Palla, 2005



(a) Typical NCP plot



(b) Caricature of network structure

Best conductance of a vertex set S of size k :

$$\Phi(k) = \min_{S \in V, |S|=k} \phi(S), \quad \phi(S) = \frac{\text{cut}(S, V \setminus S)}{\min(\text{vol}(S), \text{vol}(S \setminus V))}$$

where $\text{vol}(S) = \sum_{i \in S} k_i$ - sum of all node degrees in the set

Community detection algorithms



Author	Ref.	Label	Order
Eckmann & Moses	(Eckmann and Moses, 2002)	EM	$O(m(k^2))$
Zhou & Lipowsky	(Zhou and Lipowsky, 2004)	ZL	$O(n^3)$
Latapy & Pons	(Latapy and Pons, 2005)	LP	$O(n^3)$
Clauset et al.	(Clauset <i>et al.</i> , 2004)	NF	$O(n \log^2 n)$
Newman & Girvan	(Newman and Girvan, 2004)	NG	$O(nm^2)$
Girvan & Newman	(Girvan and Newman, 2002)	GN	$O(n^2 m)$
Guimerà et al.	(Guimerà and Amaral, 2005; Guimerà <i>et al.</i> , 2004)	SA	parameter dependent
Duch & Arenas	(Duch and Arenas, 2005)	DA	$O(n^2 \log n)$
Fortunato et al.	(Fortunato <i>et al.</i> , 2004)	FLM	$O(m^3 n)$
Radicchi et al.	(Radicchi <i>et al.</i> , 2004)	RCCLP	$O(m^4/n^2)$
Donetti & Muñoz	(Donetti and Muñoz, 2004, 2005)	DM/DMN	$O(n^3)$
Bagrow & Bollt	(Bagrow and Bollt, 2005)	BB	$O(n^2)$
Capocci et al.	(Capocci <i>et al.</i> , 2005)	CSCC	$O(n^2)$
Wu & Huberman	(Wu and Huberman, 2004)	WH	$O(n + m)$
Palla et al.	(Palla <i>et al.</i> , 2005)	PK	$O(\exp(n))$
Reichardt & Bornholdt	(Reichardt and Bornholdt, 2004)	RB	parameter dependent

Author	Ref.	Label	Order
Girvan & Newman	(Girvan and Newman, 2002; Newman and Girvan, 2004)	GN	$O(nm^2)$
Clauset et al.	(Clauset <i>et al.</i> , 2004)	Clauset et al.	$O(n \log^2 n)$
Blondel et al.	(Blondel <i>et al.</i> , 2008)	Blondel et al.	$O(m)$
Guimerà et al.	(Guimerà and Amaral, 2005; Guimerà <i>et al.</i> , 2004)	Sim. Ann.	parameter dependent
Radicchi et al.	(Radicchi <i>et al.</i> , 2004)	Radicchi et al.	$O(m^4/n^2)$
Palla et al.	(Palla <i>et al.</i> , 2005)	Cfinder	$O(\exp(n))$
Van Dongen	(Dongen, 2000a)	MCL	$O(nk^2)$, $k < n$ parameter
Rosvall & Bergstrom	(Rosvall and Bergstrom, 2007)	Infomod	parameter dependent
Rosvall & Bergstrom	(Rosvall and Bergstrom, 2008)	Infomap	$O(m)$
Donetti & Muñoz	(Donetti and Muñoz, 2004, 2005)	DM	$O(n^3)$
Newman & Leicht	(Newman and Leicht, 2007)	EM	parameter dependent
Ronhovde & Nussinov	(Ronhovde and Nussinov, 2009)	RN	$O(m^\beta \log n)$, $\beta \sim 1.3$



- M.A Porter, J-P Onella, P.J. Mucha. Communities in Networks, Notices of the American Mathematical Society, Vol. 56, No. 9, 2009
- Finding and evaluating community structure in networks, M.E.J. Newman, M. Girvan, Phys. Rev E, 69, 2004
- Modularity and community structure in networks, M.E.J. Newman, PNAS, vol 103, no 26, pp 8577-8582, 2006
- S. E. Schaeffer. Graph clustering. Computer Science Review, 1(1):27–64, 2007.
- S. Fortunato. Community detection in graphs, Physics Reports, Vol. 486, Iss. 3–5, pp 75-174, 2010



- G. Palla, I. Derenyi, I. Farkas, T. Vicsek, Uncovering the overlapping community structure of complex networks in nature and society, *Nature* 435 (2005) 814-818.
- P. Pons and M. Latapy, Computing communities in large networks using random walks, *Journal of Graph Algorithms and Applications*, 10 (2006), 191-218.
- V.D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities in large networks, *J. Stat. Mech.* P10008 (2008).
- J. Leskovec, K.J. Lang, A. Dasgupta, and M.W. Mahoney. Statistical properties of community structure in large social and information networks. In *WWW 08: Procs. of the 17th Int. Conf. on World Wide Web*, pages 695-704, 2008.