

Network communities

Leonid E. Zhukov

School of Applied Mathematics and Information Science
National Research University Higher School of Economics

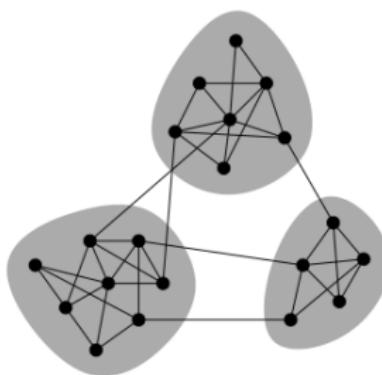
24.02.2014



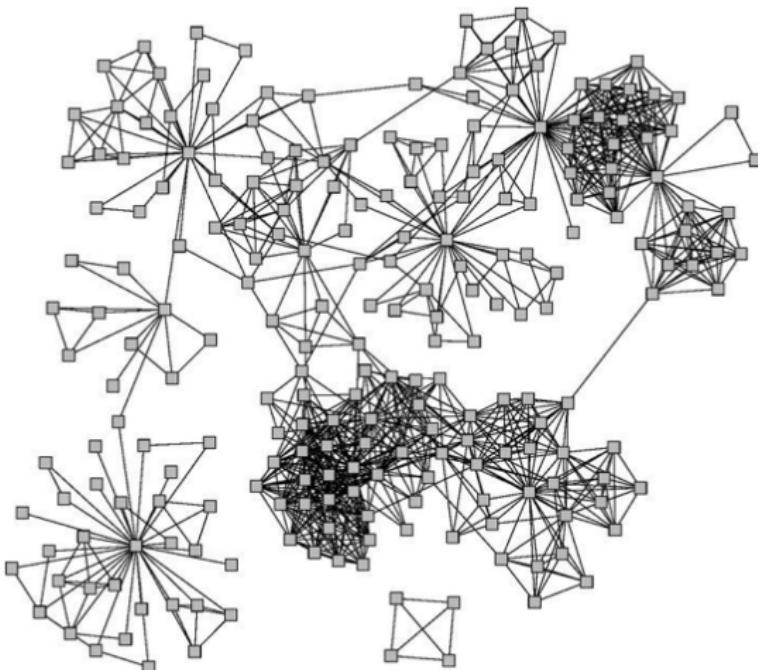
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

Definition

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



Network Communities



Network Communities

- graph density

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

- community (cluster) density

$$\delta_{int}(C) = \frac{m_c}{n_c(n_c-1)/2}$$

- external edges density

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

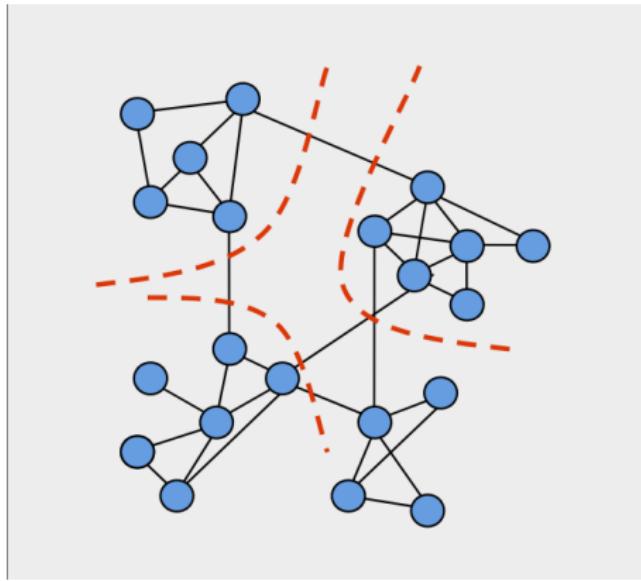
- community (cluster): $\delta_{int} > \rho, \delta_{ext} < \rho$
- cluster detection

$$\max(\delta_{int} - \delta_{ext})$$

Community detection

- Consider only sparse graphs $m \ll n^2$
- Each community should be connected
- Combinatorial optimization problem:
 - optimization criterion
 - optimization method
- Exact solution NP-hard
(bi-partition: $n = n_1 + n_2$, $n!/(n_1!n_2!)$ combinations)
- Solved by greedy, approximate algorithms or heuristics
- Recursive top-down 2-way partition, multiway partition
- Balanced class partition vs communities

Optimization criterion: graph cut



Optimization criterion : graph cut

Graph: $G(E, V)$

$$||E|| = m, ||V|| = n, V = V_1 + V_2$$

- min cut

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

- quotient cut

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

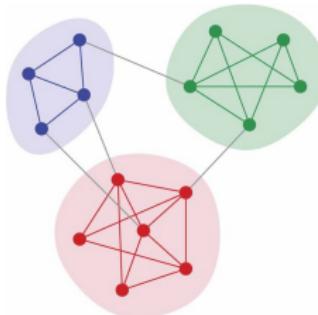
- normalized cut

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

[Minimization!]

Optimization criterion: modularity

- Let n_c - number of classes, c_i - class label per node



- Modularity:

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

$$\delta(c_i, c_j) = \begin{cases} 1 & \text{if } c_i = c_j \\ 0 & \text{if } c_i \neq c_j \end{cases}$$
 - kronecker delta

[Maximization!]

Modularity

- Random network
 - consider an edge e attached to node i , degree k_i
 - probability that it is attached to node j , degree k_j is $k_j/2m$
 - expected number of edges (average) between i and j is $k_i k_j / 2m$
- expected number of edges within the same class, c_i - class, $\delta(c_i, c_j)$ - kronecker delta

$$\langle m_c \rangle = \frac{1}{2} \sum_{ij} \frac{k_i k_j}{2m} \delta(c_i, c_j)$$

- Modularity:

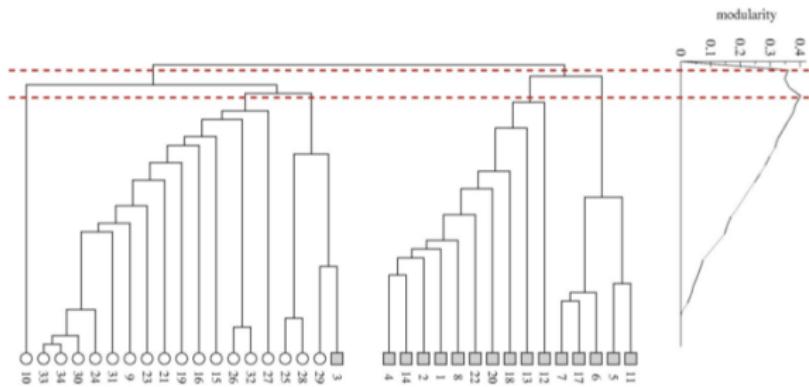
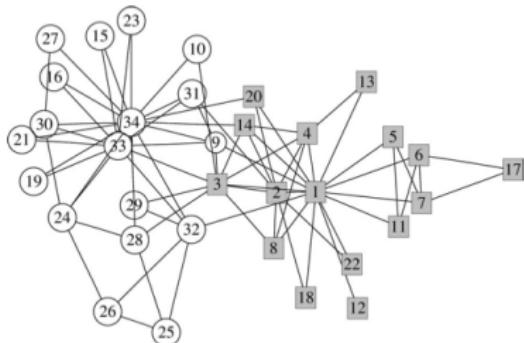
$$Q = \frac{m_c - \langle m_c \rangle}{m} = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

- Modularity matrix:

$$B_{ij} = A_{ij} - \frac{k_i k_j}{2m}$$

- Single class, $\delta(c_i, c_j) = 1$, $Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) = 0$

Dendrogram and modularity score



Optimization methods

- "Graph cut algorithms":
 - Kernighan-Lin
 - Spectral normalized cuts
 - s-t flow
 - Multilevel graph partitioning
- "Modularity algorithms":
 - Greedy
 - Spectral modularity maximization
- "Heuristics algorithms":
 - Edge betweenness
 - Random walks

Spectral Modularity Maximization

- Direct modularity maximization (bi-partitioning), [Newman, 2006]
- For two classes C_1, C_2 indicator variable $s = \pm 1$

$$\delta(c_i, c_j) = \frac{1}{2}(s_i s_j + 1) = \begin{cases} 1 : & i, j \in C_1 \text{ or } i, j \in C_2 \\ 0 : & i \in C_1, j \in C_2 \text{ or } i \in C_2, j \in C_1 \end{cases}$$

- Modularity

$$Q = \frac{1}{4m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) (s_i s_j + 1) = \frac{1}{4m} \sum_{i,j} B_{ij} s_i s_j$$

where

$$B_{ij} = A_{ij} - \frac{k_i k_j}{2m}$$

Spectral Modularity Maximization

- Quadratic form:

$$Q(\mathbf{s}) = \frac{1}{4m} \mathbf{s}^T \mathbf{B} \mathbf{s}$$

- Integer optimization - NP, relaxation $s \rightarrow x, x \in R$
- Keep norm $\|x\|^2 = \sum_i x_i^2 = \mathbf{x}^T \mathbf{x} = n$
- Quadratic optimization

$$Q'(\mathbf{x}) = \frac{1}{4m} \mathbf{x}^T \mathbf{B} \mathbf{x} - \lambda (\mathbf{x}^T \mathbf{x} - n)$$

- Eigenvector problem

$$\mathbf{B} \mathbf{x}_i = \lambda'_i \mathbf{x}_i$$

- Approximate modularity

$$Q'(\mathbf{x}_i) = \frac{n}{4m} \lambda_i$$

- Maximization - maximal λ

Spectral Modularity Maximization

- Can't choose $\mathbf{s} = \mathbf{x}_k$, can select optimal \mathbf{s}
- Decompose in the basis: $\mathbf{s} = \sum_j a_j \mathbf{x}_j$, where $a_j = \mathbf{x}_j^T \mathbf{s}$
- Modularity

$$Q(\mathbf{s}) = \frac{1}{4m} \mathbf{s}^T \mathbf{B} \mathbf{s} = \frac{1}{4m} \sum_i (\mathbf{x}_i^T \mathbf{s})^2 \lambda_i$$

- $\max Q(\mathbf{s})$ reached when $\lambda_1 = \lambda_{\max}$ and $\max \mathbf{x}_1^T \mathbf{s} = \sum_j x_{1j} s_j$
- Choose $\mathbf{s} \parallel \mathbf{x}_1$, $\mathbf{s} = sign(\mathbf{x}_1)$

Modularity maximization

Algorithm: Spectral modularity maximization: two-way partition

Input: adjacency matrix \mathbf{A}

Output: class indicator vector \mathbf{s}

compute $\mathbf{k} = \deg(\mathbf{A})$;

compute $\mathbf{B} = \mathbf{A} - \frac{1}{2m}\mathbf{k}\mathbf{k}^T$;

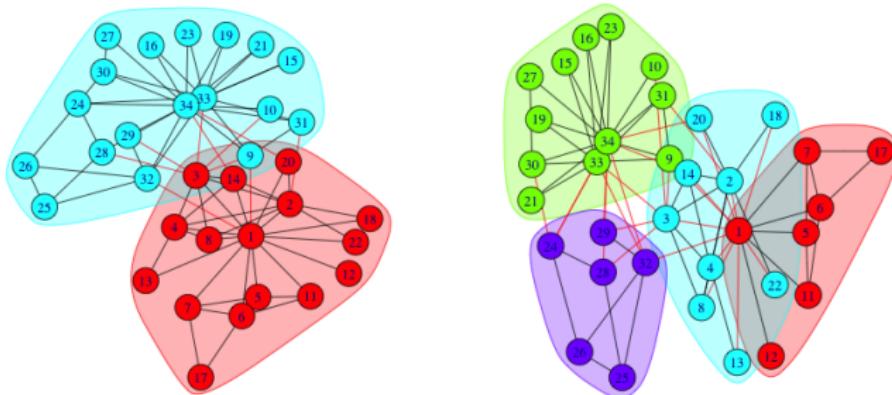
solve for maximal eigenvector $\mathbf{Bx} = \lambda\mathbf{x}$;

set $\mathbf{s} = sign(\mathbf{x}_1)$

Recursive bisection

Spectral modularity maximization

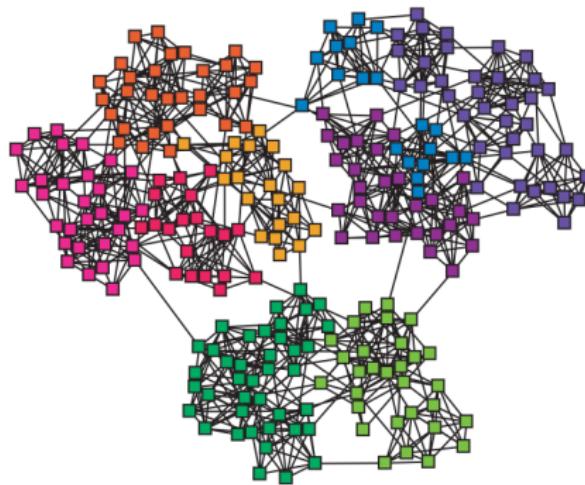
- Compute modularity matrix $B_{ij} = A_{ij} - \frac{k_i k_j}{2m}$
- Solve for maximal eigenvalue $\mathbf{Bx} = \lambda \mathbf{x}$
- set $\mathbf{s} = sign(\mathbf{x}_1)$
- recurse on each partition



Edge betweenness

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

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



Edge betweenness

Newman-Girvan, 2004

Algorithm: Edge Betweenness

Input: graph $G(V,E)$

Output: Dendrogram

repeat

 For all $e \in E$ compute edge betweenness $C_B(e)$;

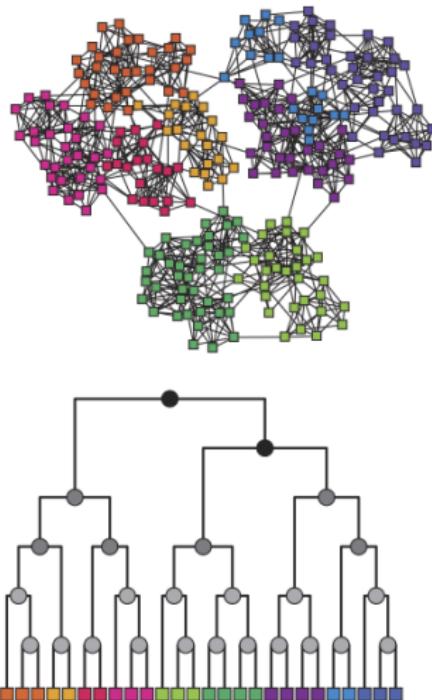
 remove edge e_i with largest $C_B(e_i)$;

until edges left;

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

Edge betweenness

Dendrogram



References

- 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