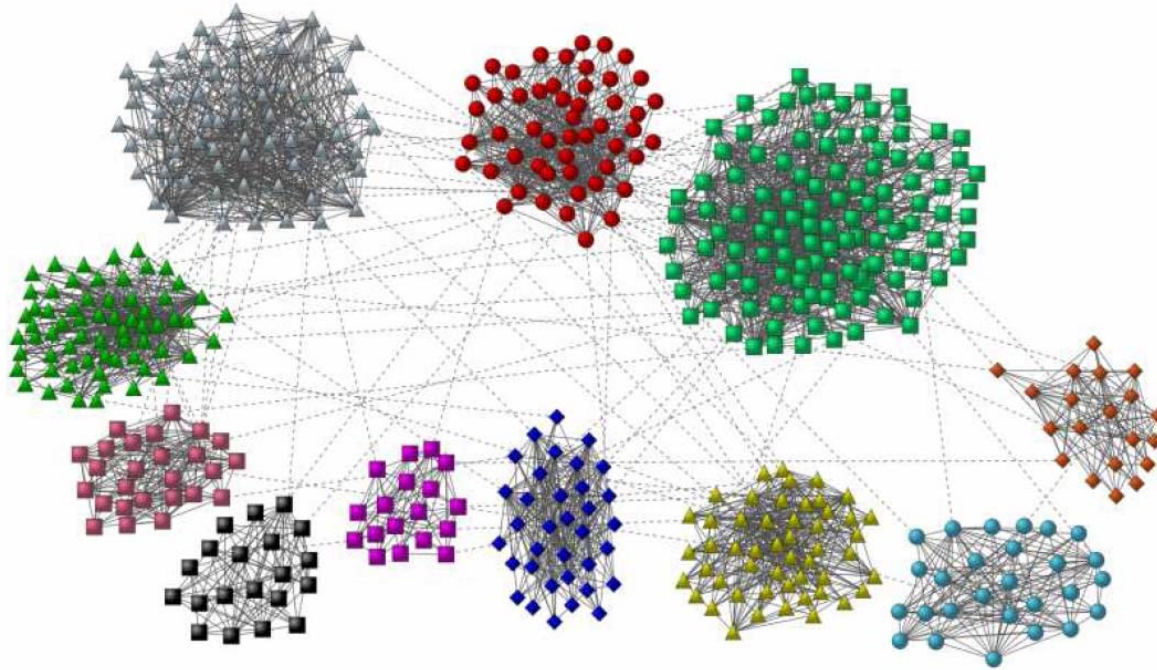# Significance of Community Structure

胡延庆

北京师范大学管理学院

2009年12月21日

# Definition of Community



In communities, the link density is comparatively high and among communities the link density is comparatively low.

Many complex systems can be represented as networks and separating a network into communities could simplify the functional analysis considerably.
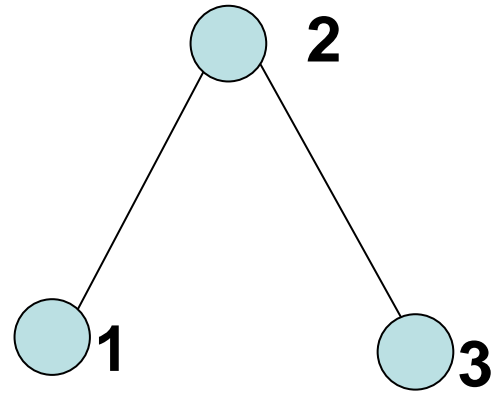
# Why We Pay Attention to Significance of Community Structure

1. Some networks contain error links.

2. Many detecting algorithms have random factors.

So we should evaluate the sensitivity of community structure

# Quantitive Definition of Multi-Communities Structure

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



$$0 = \lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \cdots \leq \lambda_n$$
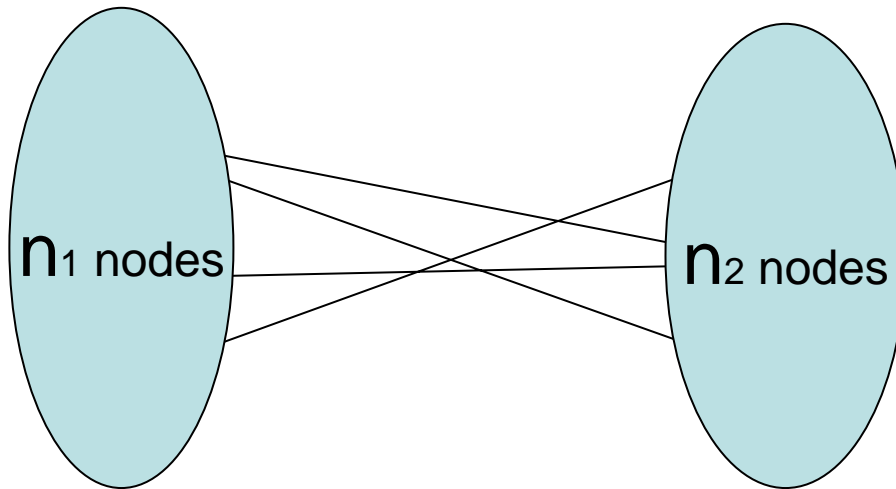
$v_i$ is the corresponding orthogonal and normalized eigenvector

$$L = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix}$$

$$span\{v_1, v_2, \cdots, v_n\} = R^n$$

$$\forall s \in R^n, s = a_1 v_1 + \cdots + a_n v_n$$

$$s = [1, -1, -1, 1, \cdots]$$

## Bi-community

denotes the bi-community structure



n$_1$ nodes          n$_2$ nodes

**Minimum number of links between the two communities.**

**Objective Function**

$$MinZ = s^T L s = \sum a_i^2 \lambda_i$$

$$st. \quad a^2{}_1 + a_2{}^2 + \cdots + a_n{}^2 = n$$

$$\text{where } a_i = v_i^T s$$

$$MinZ \approx Max\hat{Z} \, a_2^2 \lambda_2 \quad \text{simple}$$

For multi-communities structure, we define $S_j$ as community $j's$ community structure vector. If node $i$ belongs to community $j$ we let $S_{ji} = 1$ otherwise $S_{ji} = -1$
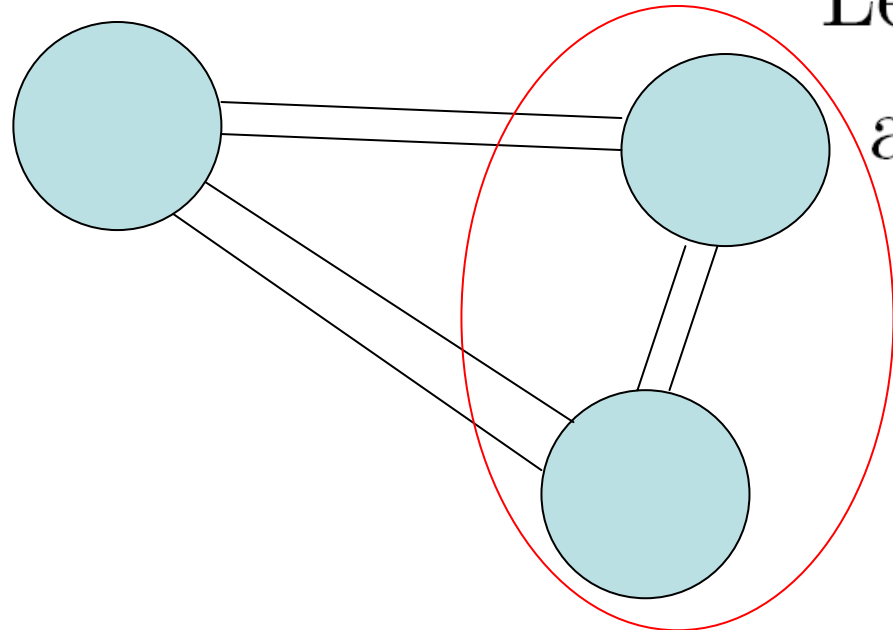
The optimal partition is defined as:

$$MinZ = \sum_{i=1}^{c} S_i^T L S_i$$

Let $S = (S_1^T, S_2^T, \cdots, S_c^T)^T$ and $\hat{L} = diag(L, L, \cdots, L)$

thus we have

$$MinZ = S^T \hat{L} S$$

We can obtain all orthogonal and normalized eigenvectors $u_q$ and the corresponding eigenvalues $\tau_q$ of $\hat{L}$, where $q = 1, 2, \cdots, n \times c$. Obviously, each eigenvalue of $L$ will be $\hat{L}$'s eigenvalue and each them will repeat $c$ time. Without losing any generality, we can let $\tau_{ci-c+j} = \lambda_i, j = 1, 2, \cdots, c$. Let $SU$ as the eigenvectors set of the eigenvalues of $\lambda_2, \lambda_3, \cdots, \lambda_c$ of matrix $\hat{L}$. $SU$ can be written as $SU = \{(v_2^T, 0 \cdots, 0), \cdots, (v_c^T, 0, \cdots, 0), \cdots, (0, 0, \cdots, v_c^T)\}$, where each $0$ denote a $n$-dimensional vector and $SU$ has $c \times (c-1)$ elements. We can expand $SU$ as a space $SSU$ in which each point is the liner combination of the elements in set $SU$. As the bi-partition problem, the multi-partition problem can be written as:

$$MinZ = \sum_{q=1}^{n \times c} b_q^2 \tau_q \approx Max\hat{Z} = \sum_{u_q \in SSU} b_q^2 \tau_q \approx \bar{\lambda} \sum_{u_q \in SSU} b_q^2$$

where $b_q = S^T u_q$ and $\overline{\lambda}$ is the average value of $\lambda_2$ to $\lambda_c$. $\sum_{u_q \in SSU} b_q^2$ denotes the length of vector $S$ projection in space $SSU$. Obviously, the longer of the projection, the $S$ is more optimal.

# Robustness of Space SSU

The space $SSU$ is expanded by the simple combination of $v_2, v_3, \cdots, v_c$, therefore, the robustness of space $SSU$ is the robustness of the eigenvalues $\lambda_2, \lambda_3, \cdots, \lambda_c$ and eigenvectors $v_2, v_3, \cdots, v_c$.

$$(\delta L + L)(\delta v_i + v_i) = (\delta \lambda_i + \lambda_i)(\delta v_i + v_i)$$

deleting the second-order small quantities we have

$$\delta L v_i + L \delta v_i = \lambda_i \delta v_i + \delta \lambda_i v_i$$

after some deductions we obtain:

$$\delta \lambda_i = \frac{v_i^T \delta L v_i}{v_i^T v_i}$$

$$\delta v_i = \sum_{j=1}^{n} h_{ij} v_j$$

$$h_{ij} = \frac{v_j^T \delta L v_i}{v_j^T v_j (\lambda_i - \lambda_j)}, (i \neq j)$$

$$|\delta\lambda_i| \leq \|\delta L\|$$

Which implies that eigenvalues are always not related to robustness of community structure for unweighted network.

Without losing any generality, for any $i \neq 1$ we can let $a_{i1} = a_{ii} = 0$.
Then the comparative error of $v_i$ can be denoted as

$$\frac{|\delta v_i|}{|v_i|} \leq \| \delta L \| \sum_{j \neq i, j=2}^{n} \frac{1}{|\lambda_i - \lambda_j|}$$

$\| \delta L \|$ is the perturbation strength

$\sum_{j \neq i, j=2}^{n} \frac{1}{|\lambda_i - \lambda_j|}$ is the robustness

Integrating the robustness of $\lambda_2$ to $\lambda_c$

we define $R$ as the robustness of space $SSU$

$$R = \sum_{j=c+1}^{n} \frac{1}{|\bar{\lambda} - \lambda_j|}$$

# Index of Significance

the most significant community structure

$$MinR = \sum_{i=c+1}^{n} \frac{1}{\lambda_{c+1} - \bar{\lambda}}$$

$$\sum_{i=1}^{n} \lambda_i = nk$$

Employ lagrange multiplier method

$R$ will achieve it's global minimum

$$R = \frac{(n-c)^2}{nk} \approx \frac{n}{k}$$

when $\bar{\lambda} = 0, \lambda_{c+1}, \lambda_{c+2}, \cdots, \lambda_n = \frac{nk}{n-c}$,

$\bar{\lambda} = 0$ implies that there are no any connections among communities.

R  holds

$$R = h\frac{n}{k}$$

$$R \propto n, \frac{1}{R} \propto k$$

Define significance index as

$$H = \frac{1}{h} \qquad H \in (0,1)$$

$k_{in}=20, k_{out}=5$

$R$

$n$

- c=2
- c=3
- c=4
- c=5

$R$

$n$

- $\mu=0.1, \gamma=2.5, \beta=1$
- $\mu=0.1, \gamma=2.5, \beta=2$
- $\mu=0.3, \gamma=2.5, \beta=1$
- $\mu=0.3, \gamma=2.5, \beta=1$

$\frac{1}{R}$

$k$

- ER,c=2,$k_{out}$=5
- ER,c=2,$k_{out}$=50
- ER,c=3,$k_{out}$=5
- ER,c=3,$k_{out}$=50

$k$

- LFR,$\mu=0.1, \beta=1$
- LFR,$\mu=0.1, \beta=2$
- LFR,$\mu=0.3, \beta=1$
- LFR,$\mu=0.3, \beta=2$

FIG. 3: The performance of $H$ index in both GN-benchmark and LFR-benchmark. In GN-benchmark, we can see that $H$ decrease with increasing of $k_{out}$. When the community structure is very clear $H$ close to 1 very much, and the network close to no community structure network $H$ close to 0.3 which implies that for a given network when $H$ is less than 0.3 it is not safe to say there exit significant community structure. In LFR-benchmark, the average degree $k = 20$, maximum degree is 50 and $p(k) \propto k^{\gamma}$. Maximum and minimum community sizes are 50 and 20 respectively, more over $p(m) \propto m^{\beta}$ where, $m$ denotes the community size. We can see that with the increase of mix parameter $\mu$, the $H$ index decrease. When $\mu \geq 0.5$ (no significant community) $H$ is near 0.3 which is similar with GN-benchmark.
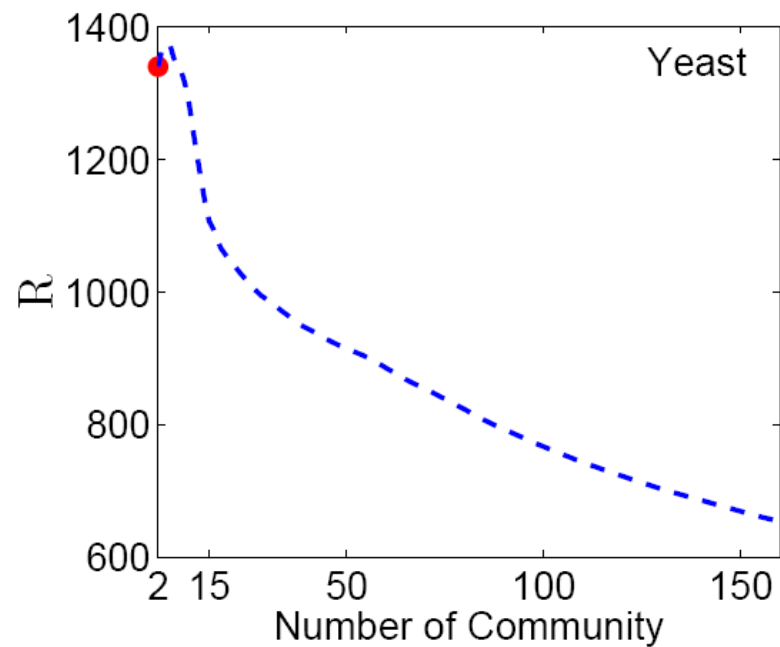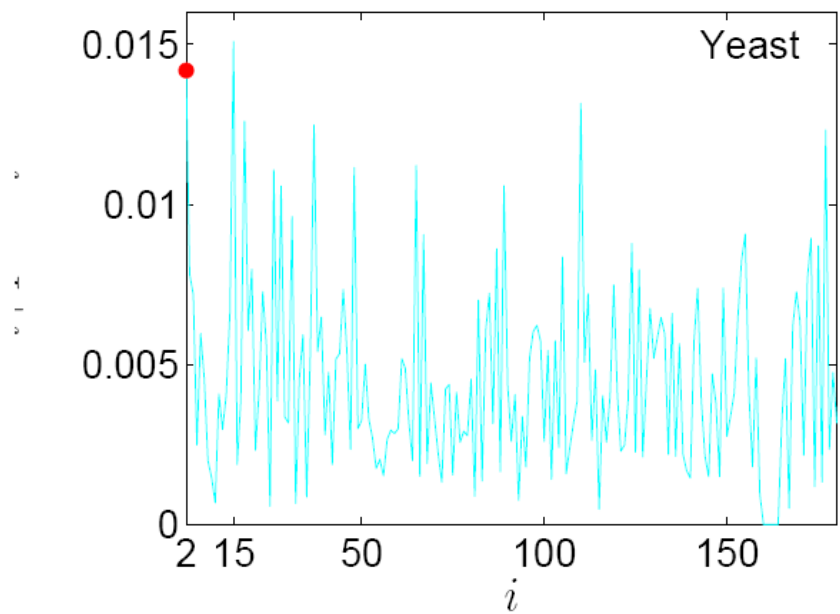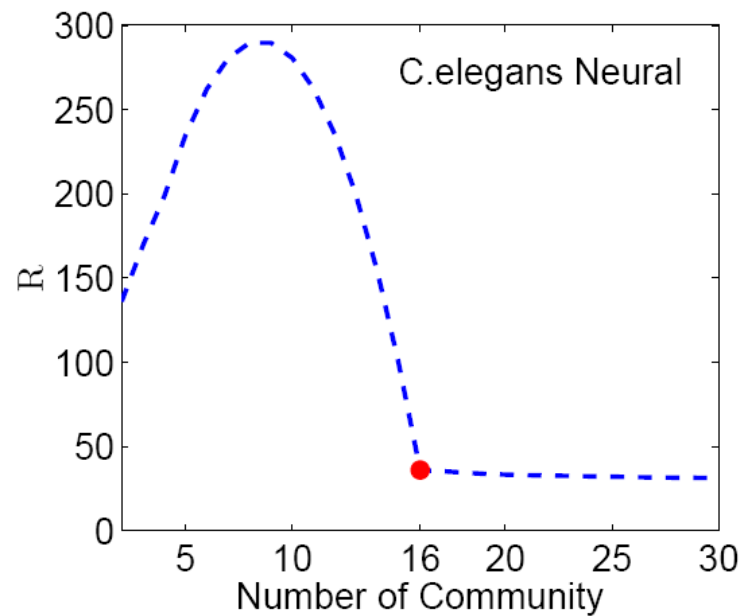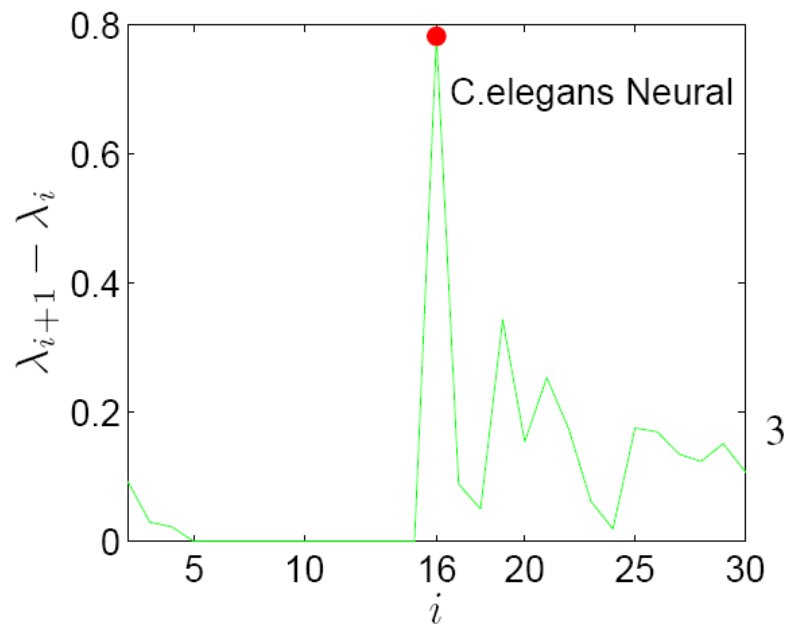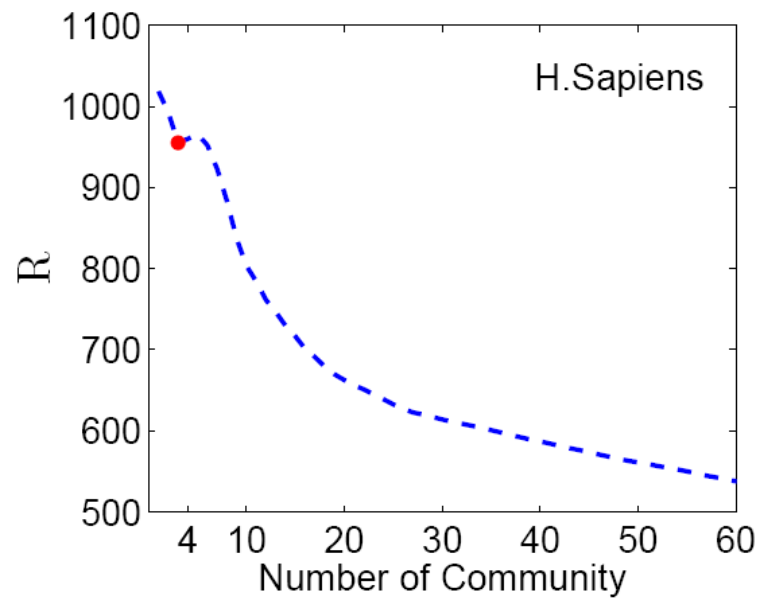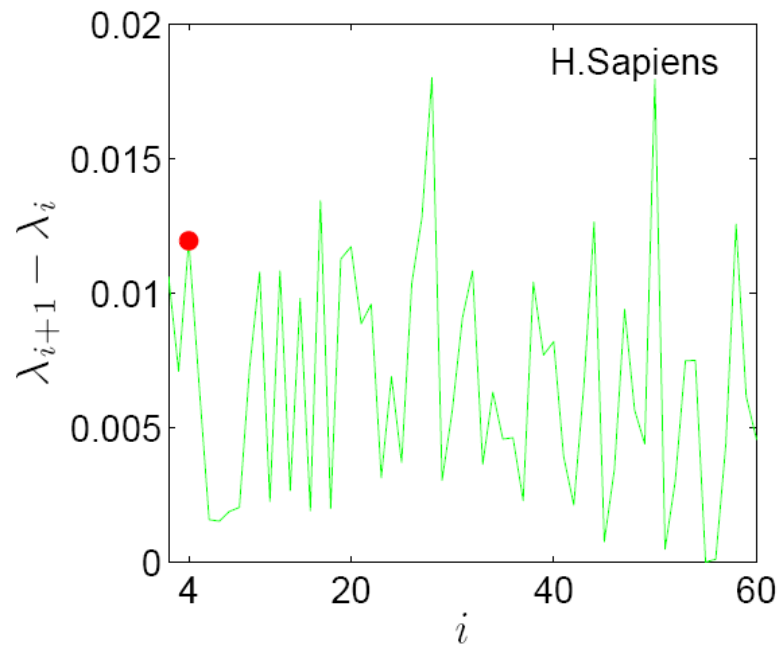
# How to obtain the optimal community number c

| network | n | edge number | $\hat{R}$ | $H$ | type |
|---|---|---|---|---|---|
| E.coli | 1442 | 5873 | 0.14 | 0.14 | protein |
| Yeast | 1458 | 1993 | 0.14 | 0.40 | |
| H.Sapiens | 693 | 982 | 0.21 | 0.21 | |
| Celegans metabolic | 453 | 4596 | 0.19 | 0.62 | metabolic |
| Aquifex aeolicus | 1437 | 3272 | 0.19 | 0.36 | |
| Helicobacter pylori | 1341 | 3087 | 0.19 | 0.36 | |
| Yersinia pestis | 1922 | 4383 | 0.18 | 0.36 | |
| Celegans neural | 297 | 2148 | 0.24 | 0.52 | neural |
| Santa Fe scientists | 260 | 612 | 0.14 | 0.22 | social |
| Zachary karate | 34 | 78 | 0.27 | 0.46 | |
| Dolphin | 62 | 159 | 0.27 | 0.42 | |
| College football | 115 | 613 | 0.38 | 0.79 | |
| Jazz | 198 | 2742 | 0.42 | 0.47 | |
| Email | 1133 | 5452 | 0.22 | 0.42 | |
| Political blogs | 1222 | 16716 | 0.29 | 0.22 | |
| Political books | 105 | 441 | 0.34 | 0.56 | |

# Thank you!