

京都府
日本



A Spectral Clustering Approach to Optimally Combining Numerical Vectors with a Modular Network

Motoki Shiga, Ichigaku Takigawa,
Hiroshi Mamitsuka



Bioinformatics Center, ICR, Kyoto University, Japan

KDD 2007, San Jose, California, USA, August 12-15 2007

Table of Contents

1. Motivation

Clustering for heterogeneous data
(numerical + network)

2. Proposed method

Spectral clustering (numerical vectors + a network)

3. Experiments

Synthetic data and real data

4. Summary

Heterogeneous Data Clustering

Heterogeneous data : various information related to an interest

Ex. Gene analysis : gene expression, metabolic pathway, ..., etc.

Web page analysis : word frequency, hyperlink, ..., etc.

Gene expression

Gene

#experiments = S

Numerical
Vectors

S-th value

k-means
SOM, etc.

To improve clustering accuracy,
combine numerical vectors + network

metabolic
pathway

Network

Minimum edge cut
Ratio cut, etc.

Related work : semi-supervised clustering

- **Local property**

Neighborhood relation

-must-link edge, cannot-link edge

- **Hard constraint** (K. Wagstaff and C. Cardie, 2000.)
- **Soft constraint** (S. Basu etc., 2004.)
 - Probabilistic model (Hidden Markov random field)

Proposed method

- **Global property** (network modularity)
- **Soft constraint**
 - Spectral clustering

Table of Contents

1. Motivation

Clustering for heterogeneous data
(numerical + network)

2. Proposed method

Spectral clustering (numerical vectors + a network)

3. Experiments

Synthetic data and real data

4. Summary

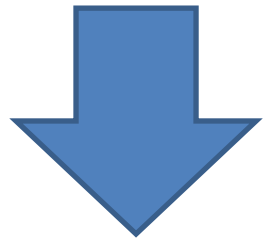
Spectral Clustering

L. Hagen, etc., IEEE TCAD, 1992., J. Shi and J. Malik, IEEE PAMI, 2000.

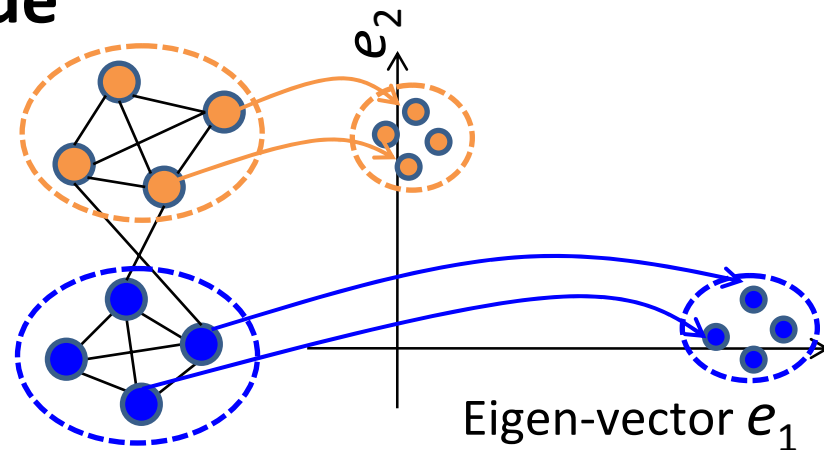
1. Compute affinity(dissimilarity) matrix \mathbf{M} from data
2. To optimize cost

$$J(\mathbf{Z}) = \text{tr}\{\mathbf{Z}^T \mathbf{M} \mathbf{Z}\} \text{ subject to } \mathbf{Z}^T \mathbf{Z} = \mathbf{I} \quad \textit{Trace optimization}$$

where $\mathbf{Z}(i,k)$ is 1 when node i belong to cluster k , otherwise 0,
compute **eigen-values and -vectors of matrix \mathbf{M}**
by relaxing $\mathbf{Z}(i,k)$ to a real value



Each node is by one or more
computed **eigenvectors**



3. Assign a cluster label to each node (by k-means)

Cost combining numerical vectors with a network

$$J = \text{tr}\{Z^T M Z\}$$
$$= (1 - \omega) J_{\text{num}}(Z) + \omega J_{\text{net}}(Z)$$

Cost of **numerical vector** **network**

cosine dissimilarity

$$J_{\text{num}}(Z) = \frac{1}{2} - \text{tr} \left(\frac{Z^T (2N)^{-1} Y Z}{Z^T Z} \right)$$

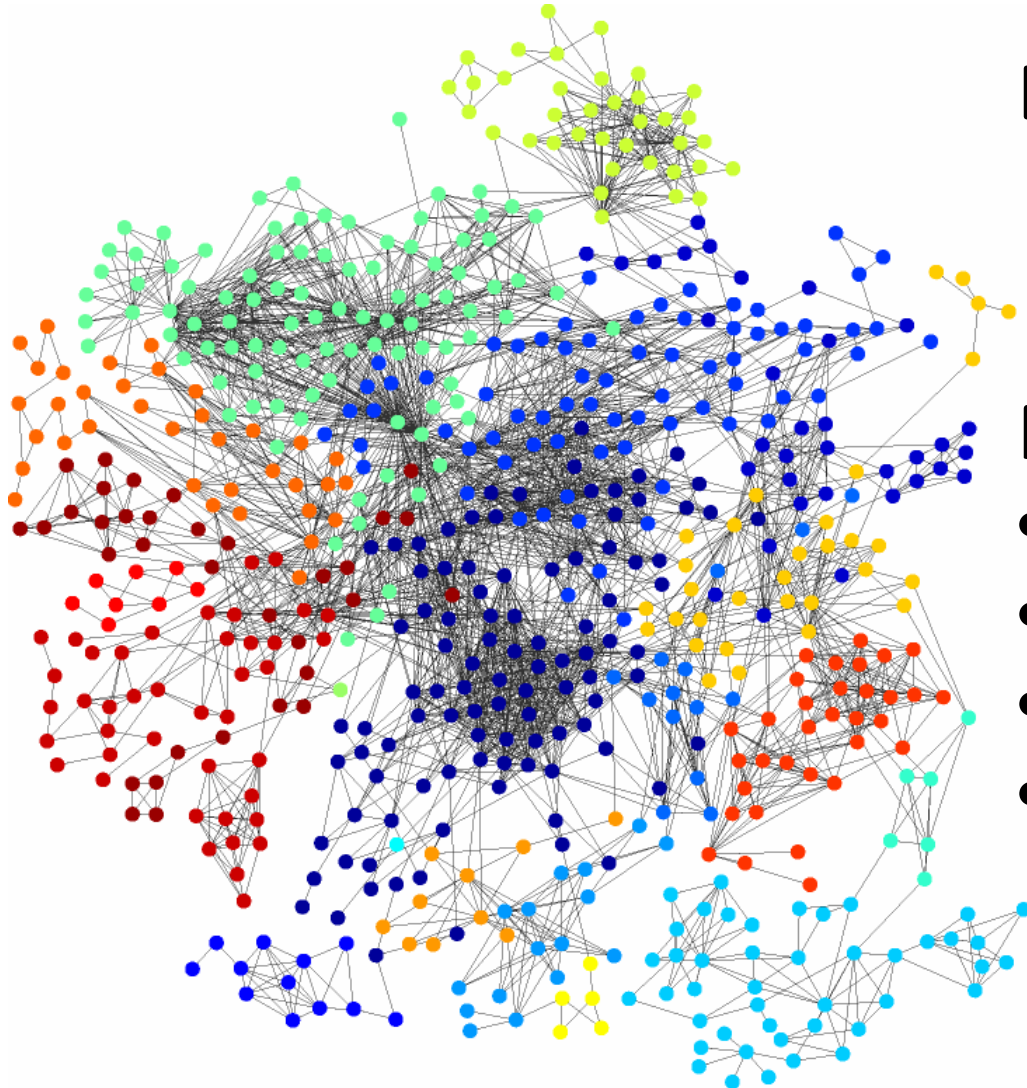
What cost?

N : #nodes,

Y : inner product of normalized numerical vectors

*To define a cost of a network,
use a property of complex networks*

Complex Networks



Ex. Gene networks,
WWW,
Social networks, ..., etc.

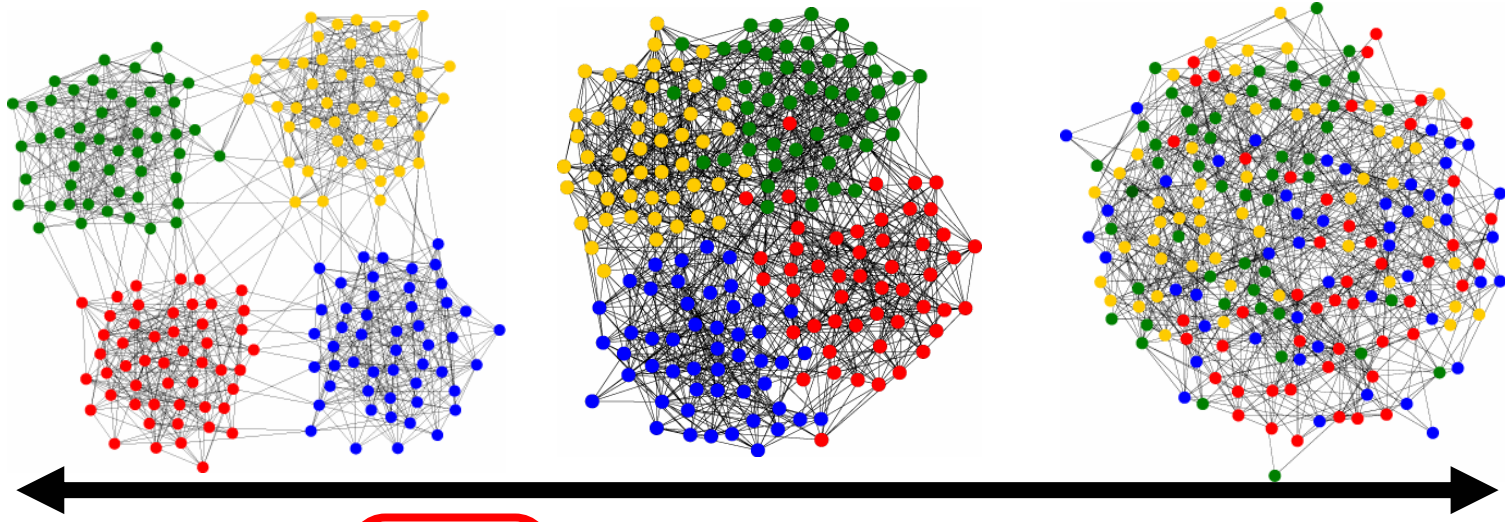
Property

- Small world phenomena
- Power law
- Hierarchical structure
- Network modularity

Ravasz, et al., Science, 2002.
Guimera, et al., Nature, 2005.

Network Modularity

= density of intra-cluster edges



$$Q(\mathcal{Z}) = \sum_{k=1}^K \left\{ \frac{L(\mathcal{Z}_k, \mathcal{Z}_k)}{L} - \left(\frac{L(\mathcal{Z}_k, \mathcal{Z})}{L} \right)^2 \right\}$$

intra-edges # total edges

normalize by cluster size

- \mathcal{Z} : set of whole nodes
- \mathcal{Z}_k : set of nodes in cluster k
- $L(A,B)$: #edges between A and B

Cost Combining Numerical Vectors with a Network

$$\begin{aligned}
 J &= \text{tr}\{Z^T M Z\} \\
 &= (1 - \omega) J_{\text{num}}(Z) + \omega J_{\text{net}}(Z)
 \end{aligned}$$

Cost of **numerical vector**

network

cosine dissimilarity

$$J_{\text{num}}(Z) = \frac{1}{2} - \text{tr} \left(\frac{Z^T (2N)^{-1} Y Z}{Z^T Z} \right)$$

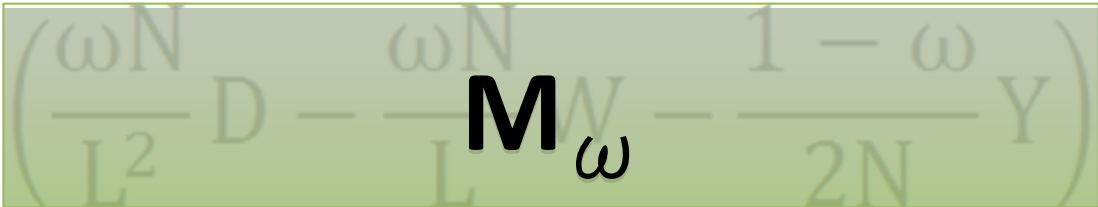
Normalized modularity
(Negative)

$$J_{\text{net}}(Z) = -\text{tr} \left(\frac{Z^T N \left(\frac{1}{L^2} D - \frac{1}{L} W \right) Z}{Z^T Z} \right)$$



$$\tilde{Z} = \frac{Z}{\sqrt{Z^T Z}}$$

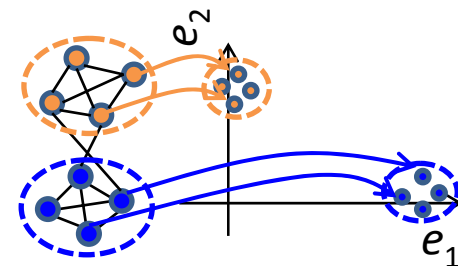
$$= \text{tr}\{\tilde{Z}^T \mathbf{M}_\omega \tilde{Z}\}$$



Our Proposed Spectral Clustering

for $\omega = 0 \dots 1$

1. Compute matrix $\mathbf{M}_\omega = \frac{\omega N}{L^2} \mathbf{D} - \frac{\omega N}{L} \mathbf{W} - \frac{1-\omega}{2N} \mathbf{Y}$
2. To optimize cost $J(\mathbf{Z}) = \text{tr}\{\mathbf{Z}^T \mathbf{M}_\omega \mathbf{Z}\}$ subject to $\mathbf{Z}^T \mathbf{Z} = \mathbf{I}$, compute eigen-values and -vectors of matrix \mathbf{M}_ω by relaxing elements of \mathbf{Z} to a real value



Each node is represented by $K-1$ eigen-vectors

3. Assign a cluster label to each node by k-means. (k-means outputs $\text{Cost}_{\text{spectral}}$ in spectral space.)

end

• **Optimize weight ω**

$$\omega^* \leftarrow \text{argmin}_{0 \leq \omega \leq 1} \text{Cost}_{\text{spectral}}$$

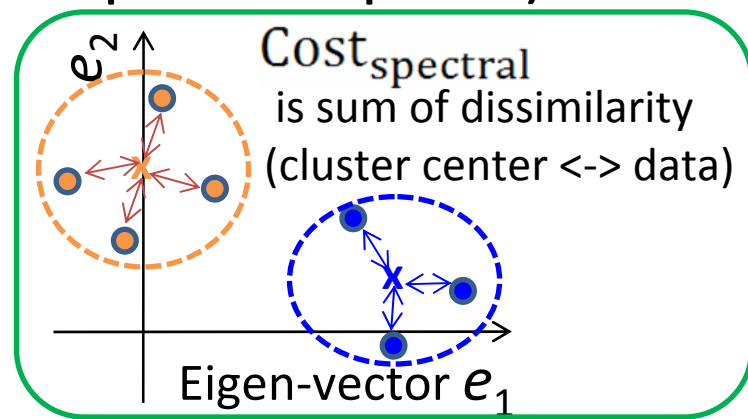


Table of Contents

1. Motivation

Clustering for heterogeneous data
(numerical + network)

2. Proposed method

Spectral clustering (numerical vectors + a network)

3. Experiments

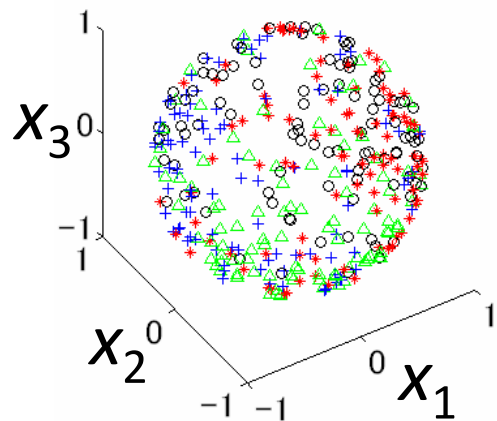
Synthetic data and real data

4. Summary

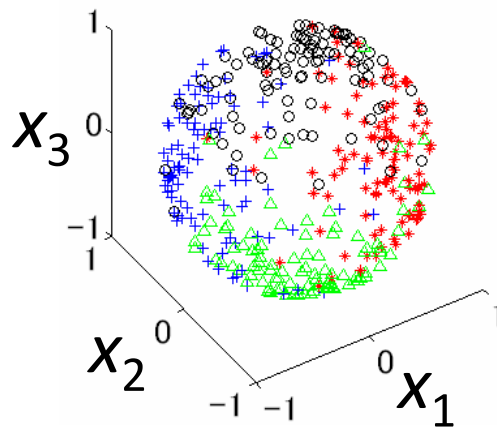
Synthetic Data

Numerical vectors (von Mises-Fisher distribution)

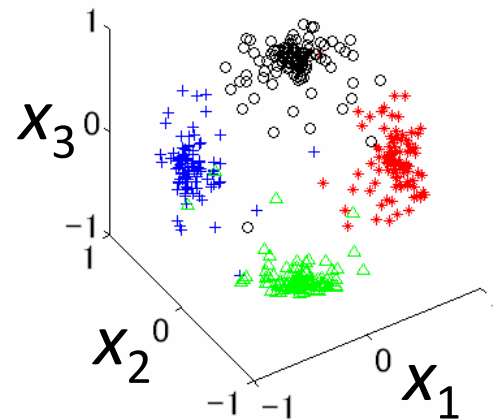
$\theta = 1$



5



50

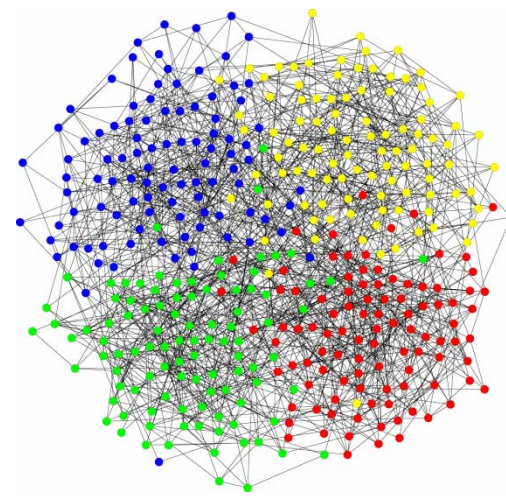
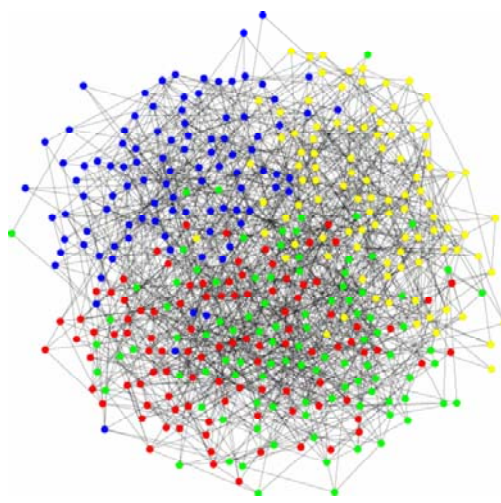
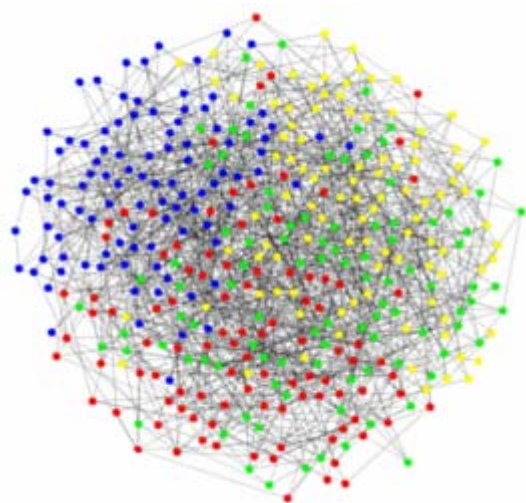


Network (Random graph) #nodes = 400, #edges = 1600

Modularity = 0.375

0.450

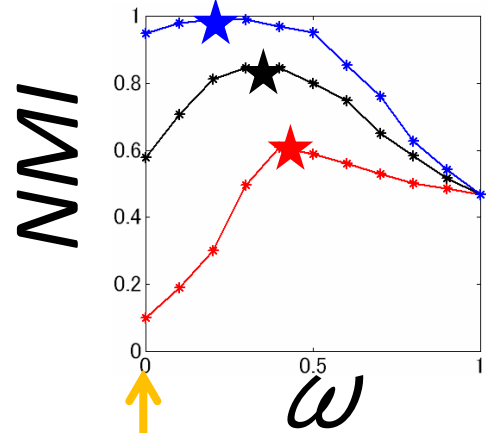
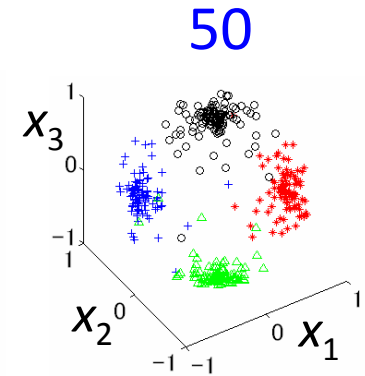
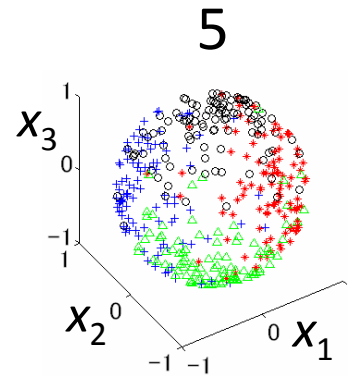
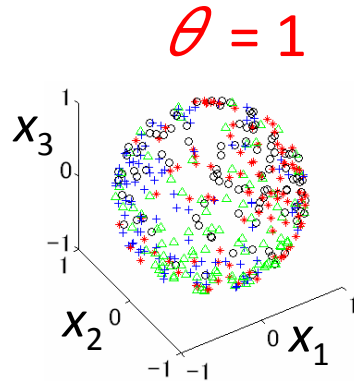
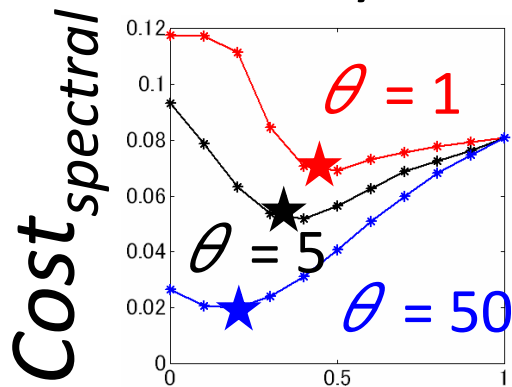
0.525



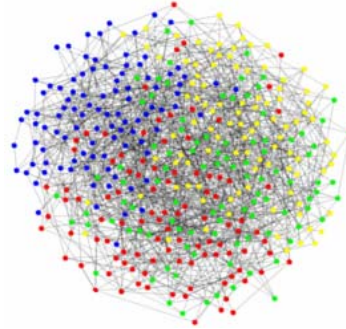
Results for Synthetic Data

Modularity = 0.375

Numerical vectors



Network



#nodes = 400, #edges = 1600
Modularity = 0.375

Numerical vectors only
(k-means)

Network only
(maximum modularity)

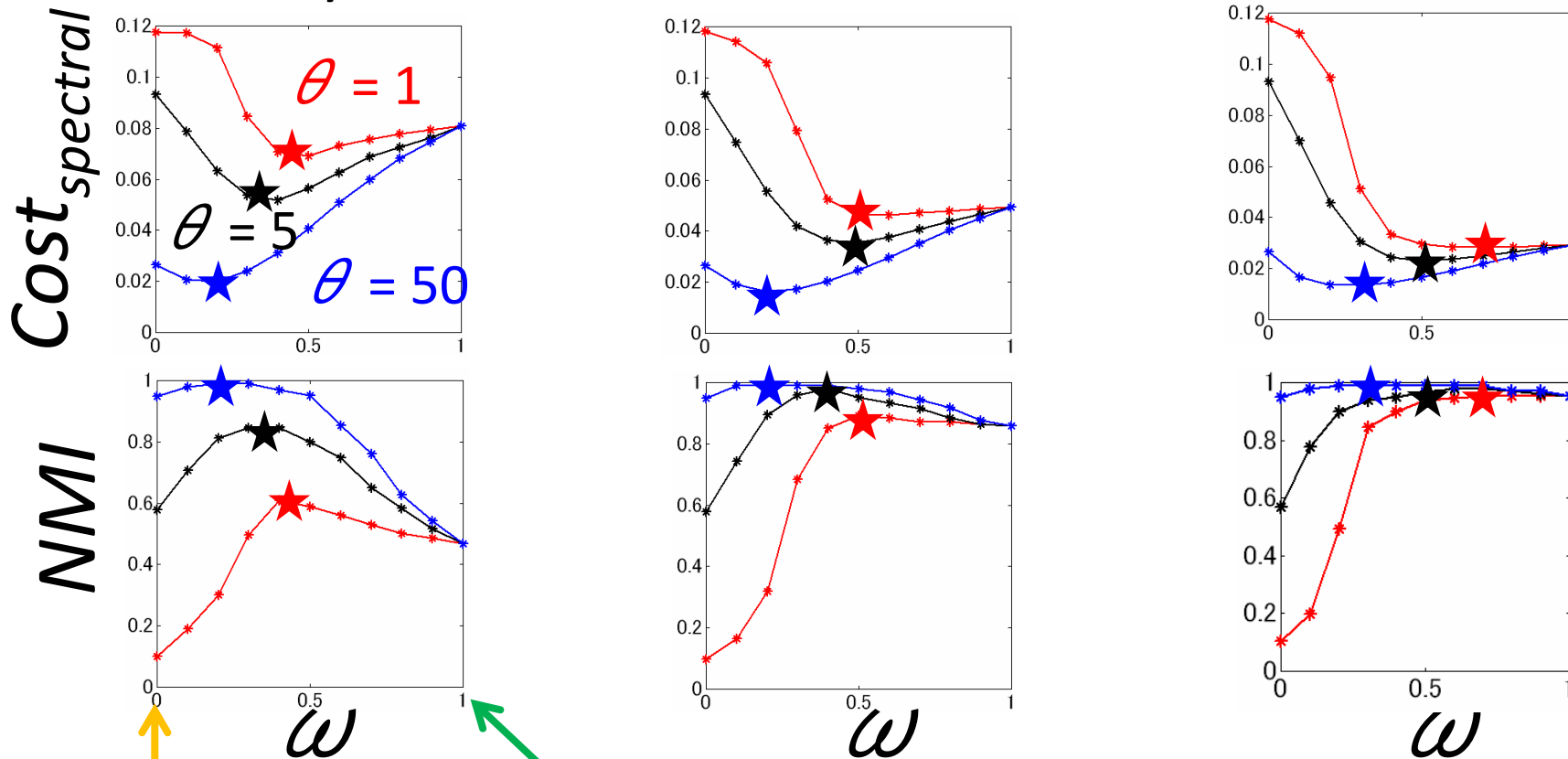
- Best NMI (Normalized Mutual Information) is in $0 < \omega < 1$
- Can be optimized using $Cost_{spectral}$

Results for Synthetic Data

Modularity = 0.375

0.450

0.525



Numerical vectors only
(k-means)

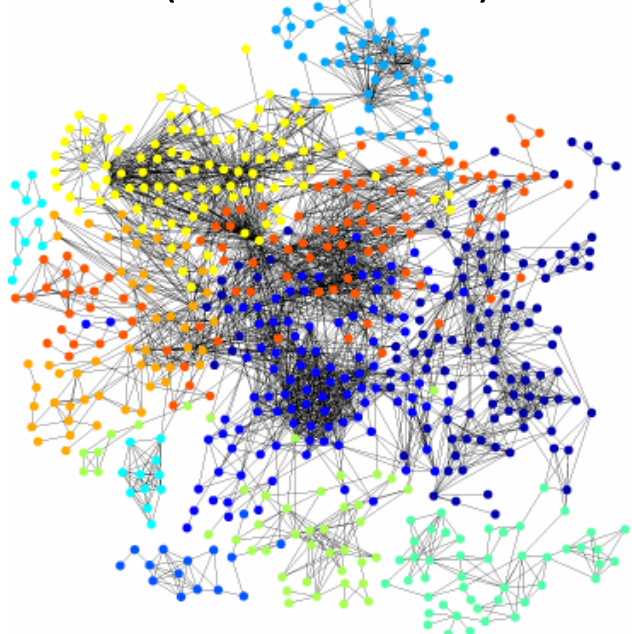
Network only
(maximum modularity)

- Best NMI (Normalized Mutual Information) is in $0 < \omega < 1$
- Can be optimized using $\text{Cost}_{\text{spectral}}$

Synthetic Data (Numerical Vector) + Real Data (Gene Network)

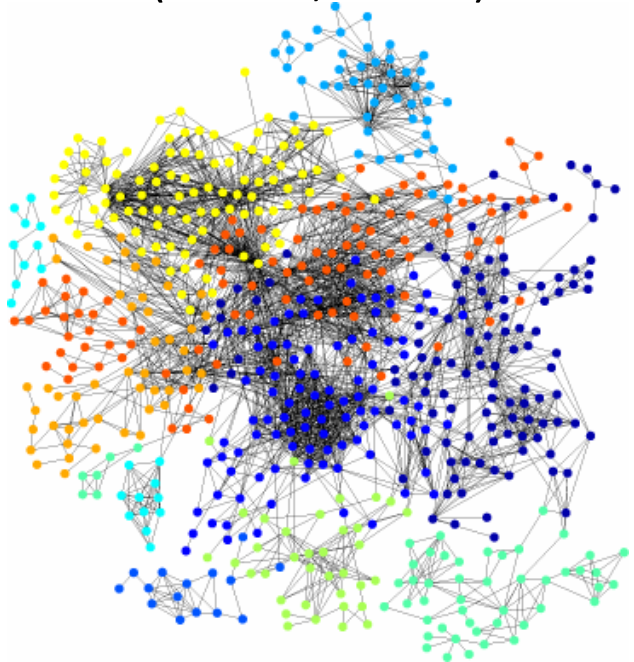
True cluster

(#clusters = 10)



Resultant cluster

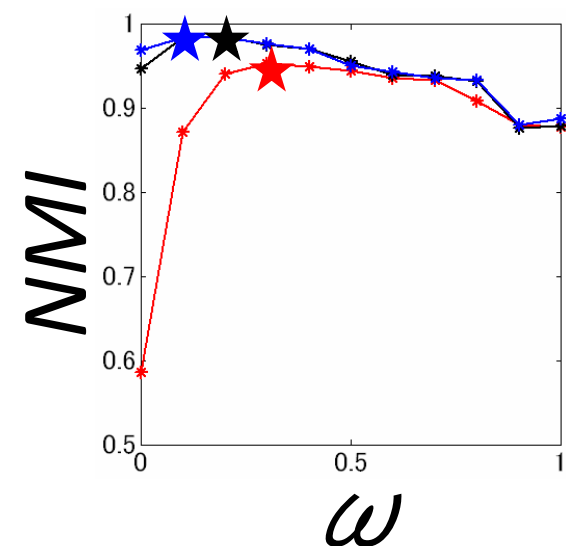
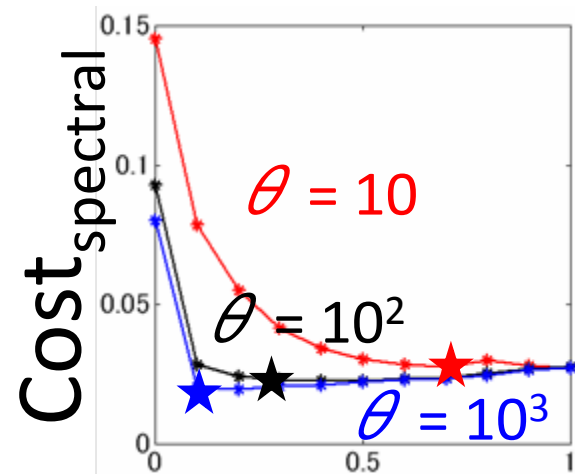
($\omega=0.5, \theta=10$)



Gene network
by KEGG metabolic pathway

- Best NMI is in $0 < \omega < 1$
- Can be optimized using

Cost



Summary

- **New spectral clustering method proposed**
combining numerical vectors with a network
 - **Global network property** (normalized network modularity)
 - Clustering can be optimized by the weight
- **Performance confirmed experimentally**
 - Better than numerical vectors only and a network only
 - **Optimizing the weight** with synthetic dataset and semi-real dataset

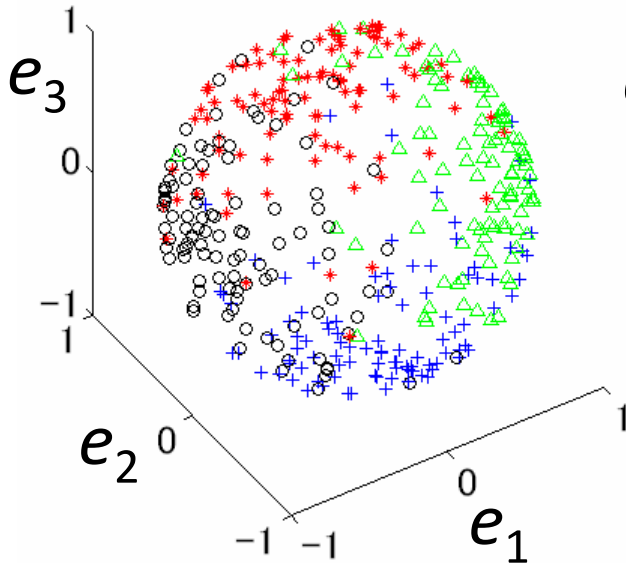
Thank you for your attention!

our poster **#16**

Spectral Representation of M_ω

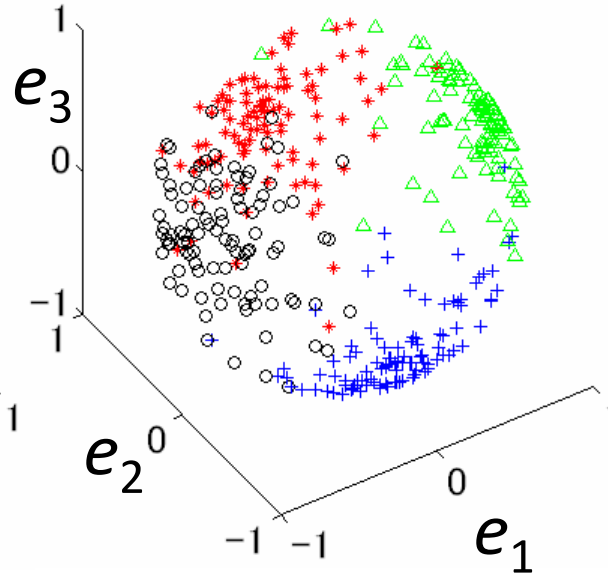
(concentration $\theta = 5$, Modularity = 0.375)

$\omega = 0$



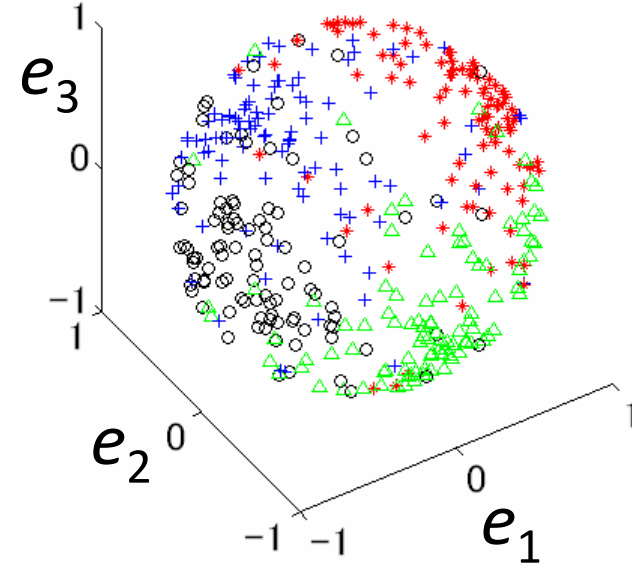
Cost $J_\omega = 0.0932$

$\omega = 0.3$



0.0538

$\omega = 1$



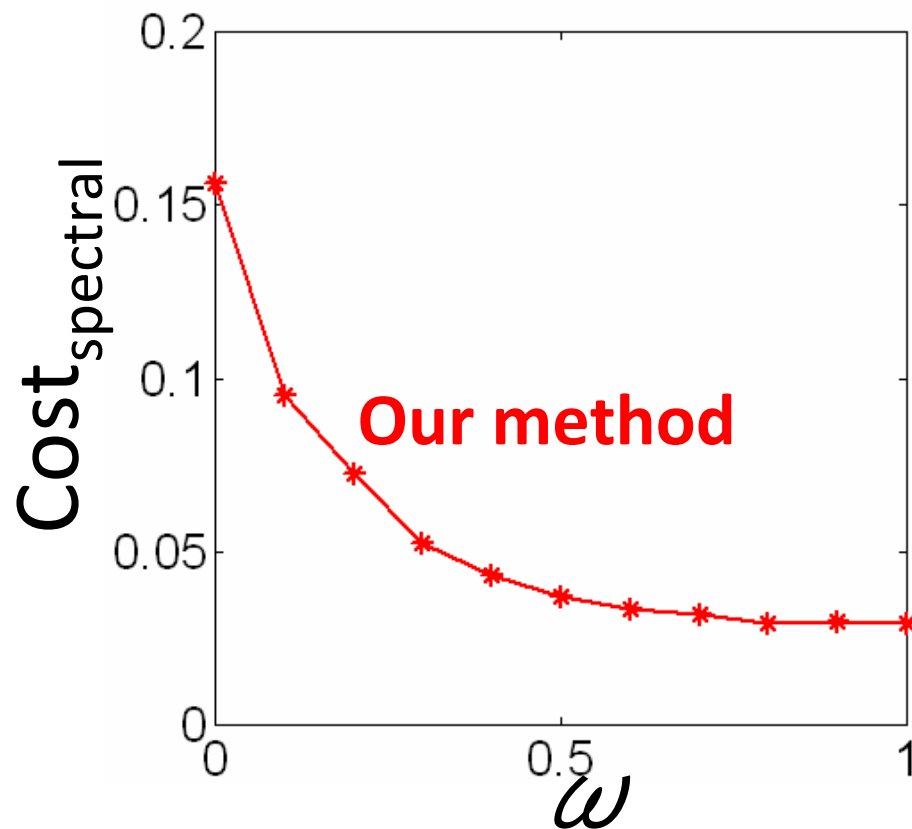
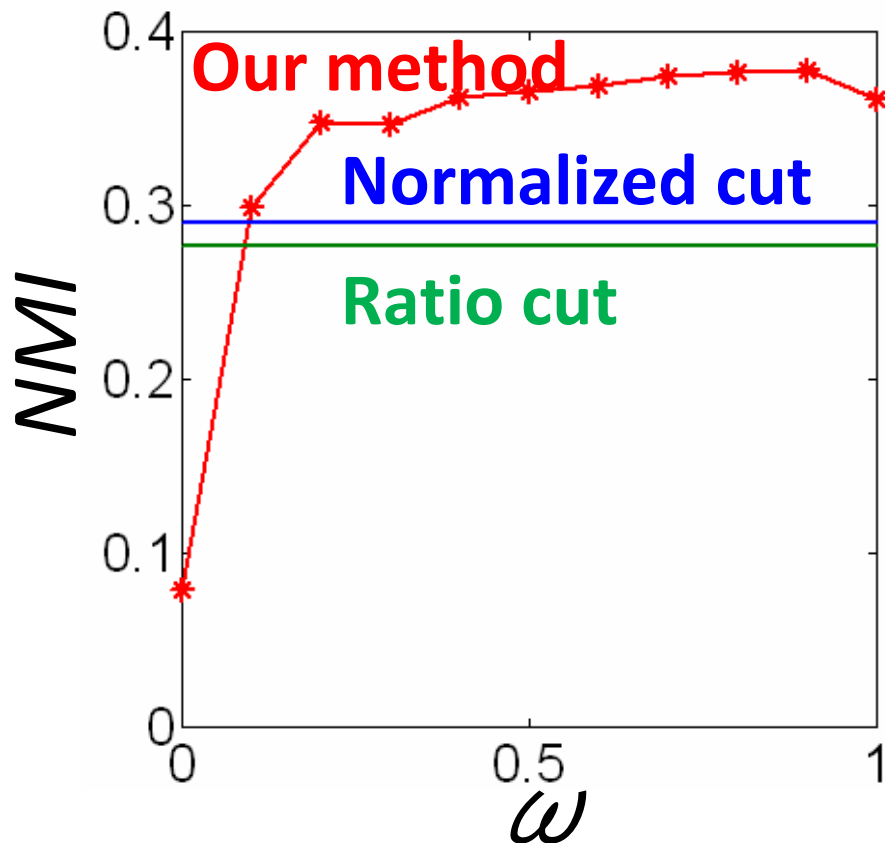
0.0809

Select ω by minimizing $\text{Cost}_{\text{spectral}}$
(clusters are divided most separately)

Result for Real Genomic Data

- Numerical vectors : Hughes' expression data (*Hughes, et al., cell, 2000*)
- Gene network : Constructed using KEGG metabolic pathways

(*M. Kanehisa, etc. NAR, 2006*)



Evaluation Measure

Normalized Mutual Information (NMI)

between estimated cluster and the standard cluster

$$NMI = \frac{H(C) + H(G) - H(C, G)}{\sqrt{H(C)}\sqrt{H(G)}}$$

$H(C)$: Entropy of probability variable C ,

C : Estimated clusters, G : Standard clusters

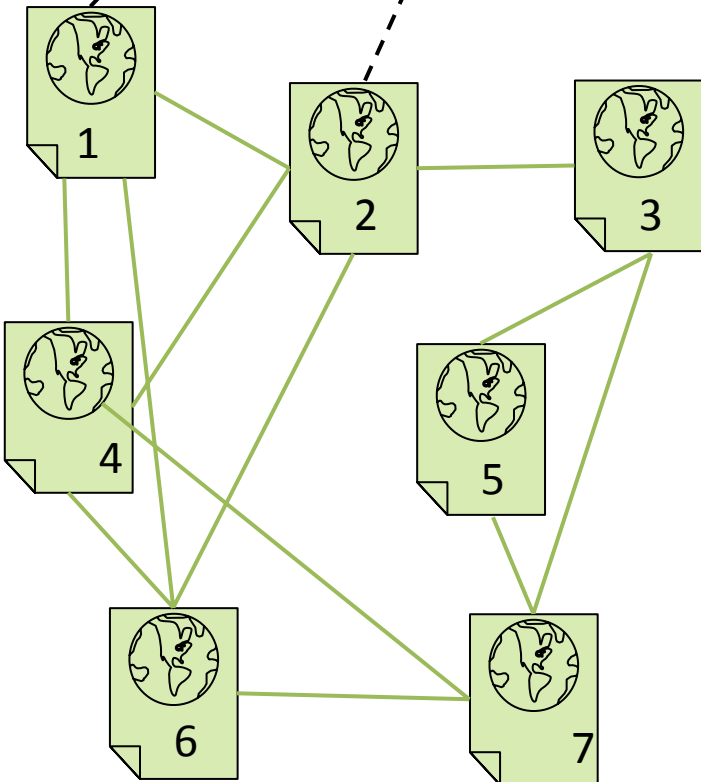
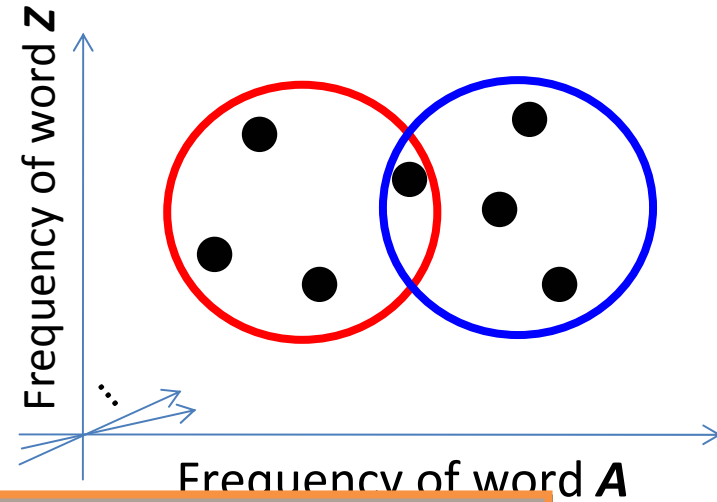
The more **similar** clusters

C and G are, **the larger the NMI.**

Web Page Clustering

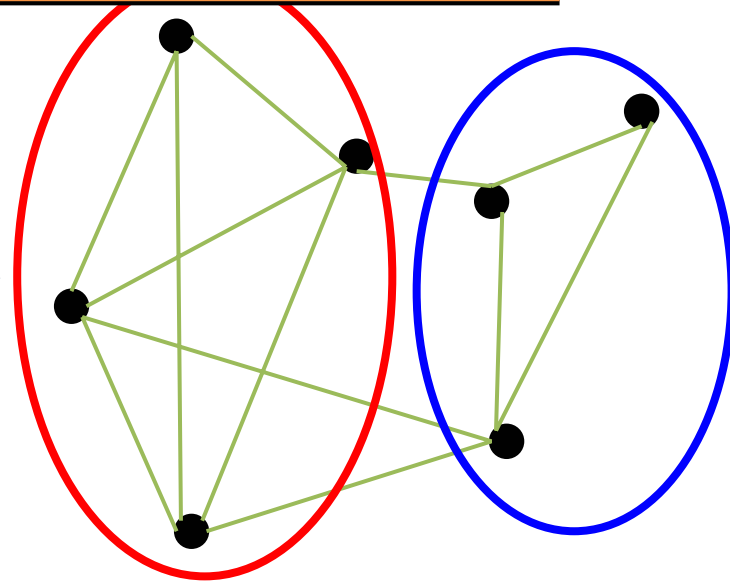
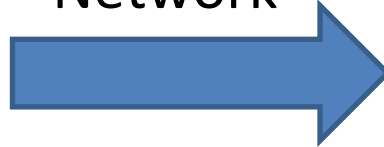
Word	Frequency	Frequency
A	$n_{(A,1)}$	
B	$n_{(B,1)}$	$n_{(A,2)}$
...	...	$n_{(B,2)}$

Numerical
Vector



To improve accuracy,
combine heterogynous data

Network



Spectral Clustering for Graph Partitioning

Ratio cut

$$\sum_k \frac{L(\mathcal{Z}_k, \mathcal{Z} \setminus \mathcal{Z}_k)}{|\mathcal{Z}_k|}$$



$$\text{tr}(\tilde{\mathbf{Z}}^T (\mathbf{D}_d - \mathbf{W}) \tilde{\mathbf{Z}})$$

Subject to $\tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} = \mathbf{I}_K$

Normalized cut

$$\sum_k \frac{L(\mathcal{Z}_k, \mathcal{Z} \setminus \mathcal{Z}_k)}{L(\mathcal{Z}_k, \mathcal{Z})}$$



$$\text{tr}(\tilde{\mathbf{Z}}^T (\mathbf{D}_d - \mathbf{W}) \tilde{\mathbf{Z}})$$

Subject to $\tilde{\mathbf{Z}}^T \mathbf{D}_d \tilde{\mathbf{Z}} = \mathbf{I}_K$