

Case Study 2: Document Retrieval

Spectral Clustering

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Emily Fox

February 12th, 2013

Document Retrieval

- **Goal:** Retrieve documents of interest



Task 1: Find Similar Documents

■ Setup

- **Input:** Query article **X**
- **Output:** Set of k similar articles



X



k-Nearest Neighbor

- Articles $X = \{x^1, \dots, x^N\}$, $x^i \in \mathbb{R}^d$

- Query: $x \in \mathbb{R}^d$

- k-NN

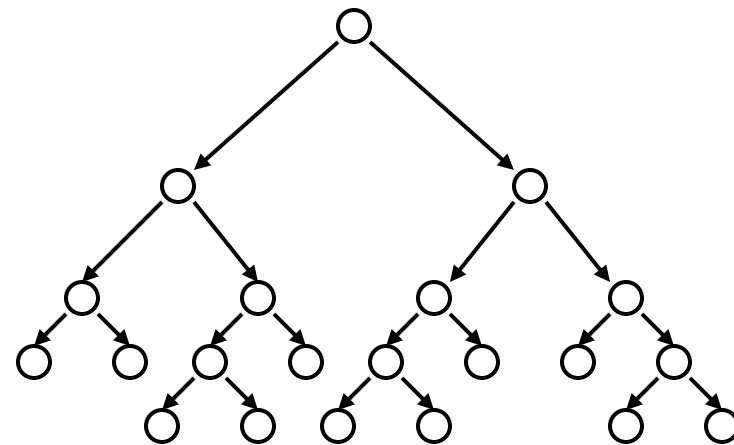
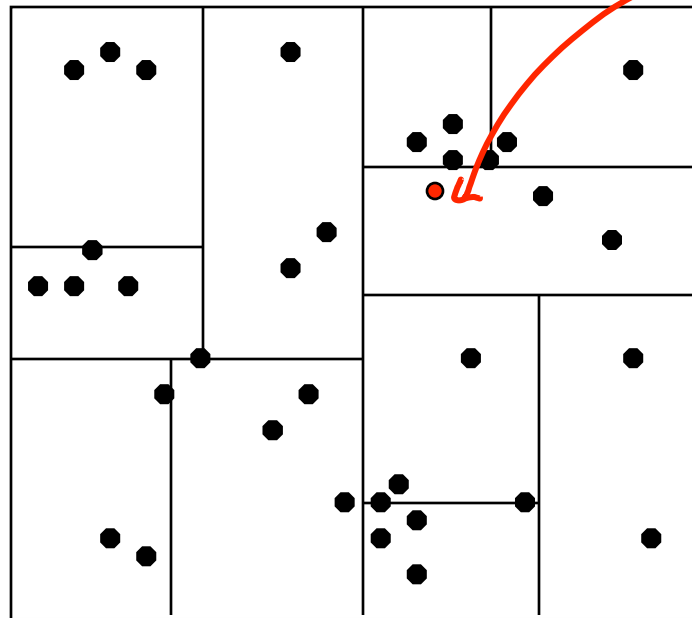
- Goal:

Find k articles in X closest x

- Formulation:

$$X^{NN} = \{x^{NN_1}, \dots, x^{NN_k}\} \subseteq X$$
$$\text{s.t. } \forall x^i \in X \setminus X^{NN}$$
$$d(x^i, x) \geq \max_{x^{NN_i} \in X^{NN}} d(x^{NN_i}, x)$$

Nearest Neighbor with KD Trees



- Traverse the tree looking for the nearest neighbor of the query point.

Task 2: Cluster Documents

■ Setup

- **Input:** Corpus of documents
- **Output:** Topic assignment per document

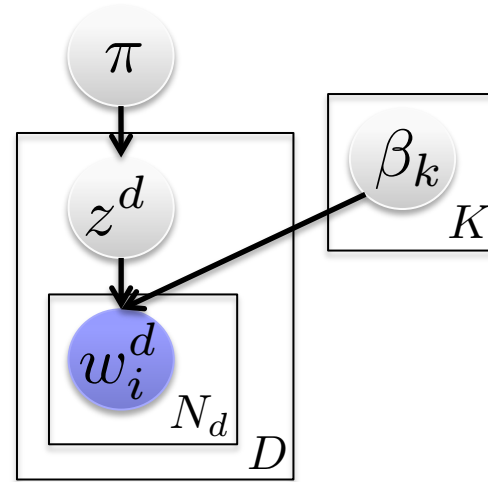
Sports

world news



A Generative Model

- Documents: x^1, \dots, x^D
- Associated topics: z^1, \dots, z^D
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:



$$z^d \sim \pi \quad d=1, \dots, D$$

$$w_i^d | z^d \sim \beta_{z^d} \quad i=1, \dots, N$$

↑ word prob. for cluster/topic z^d

Bayesian approach:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$

$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V)$$

← size of vocab.
 β_k is a V -dim pmf

Inference

■ Two tasks

□ Point estimation:

$\hat{\theta}^{ML}$, or $\hat{\theta}^{MAP}$ ← $p(\theta)$ prior

- Expectation-Maximization (EM)

□ Characterize posterior:

- Gibbs sampling
- Variational methods
- Stochastic variational inference

EM Algorithm

- Initial guess: $\hat{\theta}^{(0)}$
- Estimate at iteration t : $\hat{\theta}^{(t)}$

- E-Step

Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y | \theta) | x, \hat{\theta}^{(t)}]$

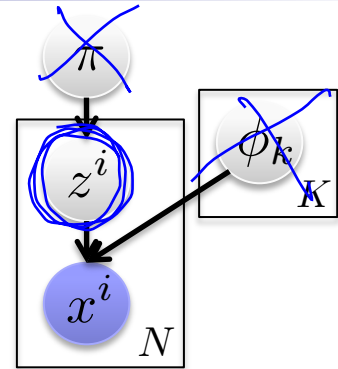
- M-Step

Compute $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)}) + \log p(\theta)$

Collapsed Gibbs Sampling

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$

$$\{\mu_k, \Sigma_k\} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

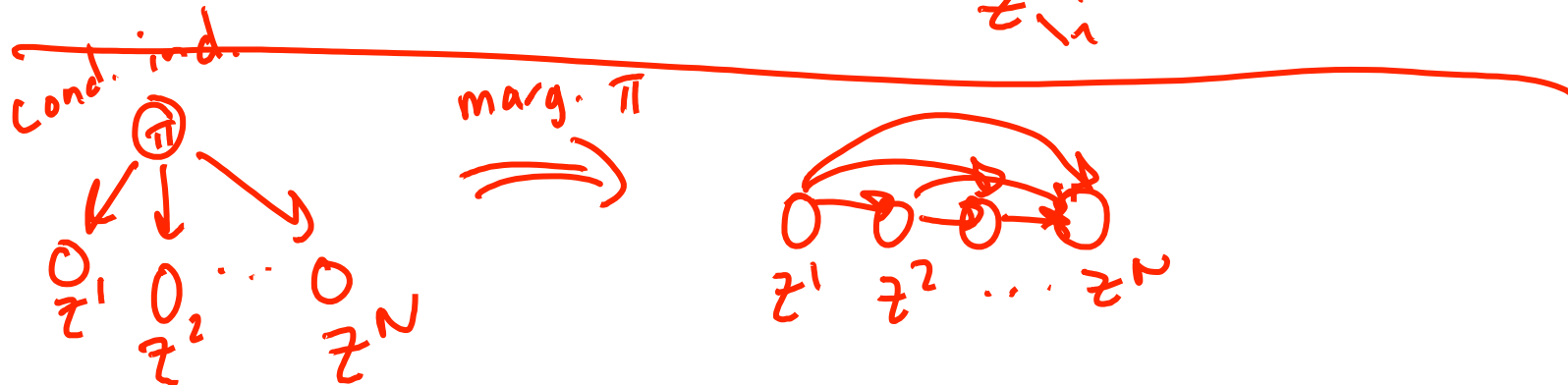


■ Collapsed sampler

For $i=1, \dots, N$

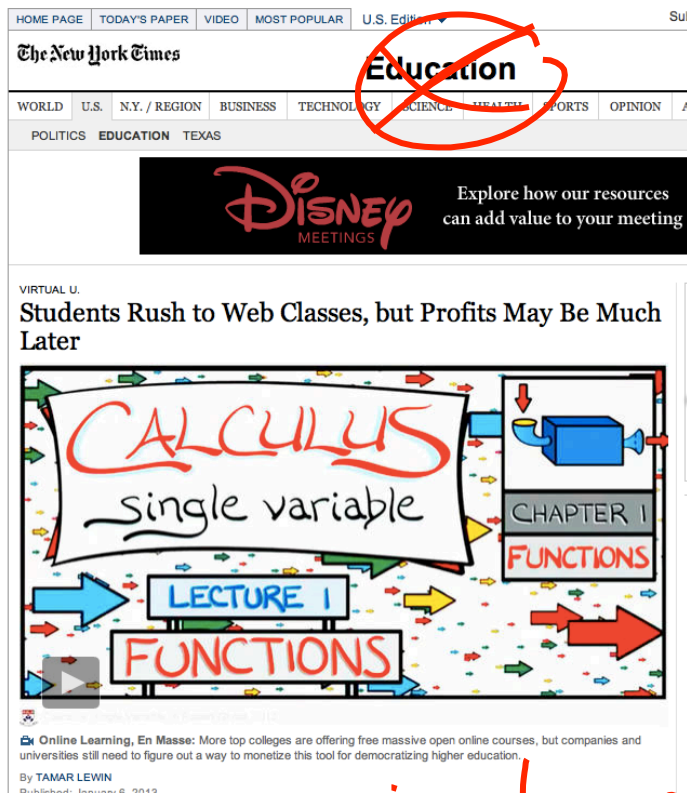
$$z^i(t) \sim p(z^i | z^{1(t)}, \dots, z^{i-1(t)}, z^{i+1(t)}, \dots, z^{N(t)}, X_{1:N}, \alpha, \phi)$$

$z_{-i}^{(t)}$



Task 3: Mixed Membership Model

- **Setup:** Document may belong to multiple clusters



EDUCATION

FINANCE

TECHNOLOGY

mixed membership

Latent Dirichlet Allocation (LDA)

each topic as a dist. over words { β_k }

Topics	
gene	0.04
dna	0.02
genetic	0.01
...	
life	0.02
evolve	0.01
organism	0.01
...	
brain	0.04
neuron	0.02
nerve	0.01
...	
data	0.02
number	0.02
computer	0.01
...	

Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson at Uppsala University in Sweden. "You arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

each doc is a mixture of these corpus-wide topics

Topic proportions and assignments

every word is assigned to a topic

each doc has its own prevalence of topics in doc

Variational Methods

- Recall task: Characterize the posterior $p(\theta, z | x)$
params \uparrow *latent vars* \uparrow *obs*
- Turn posterior inference into an optimization task
- Introduce a “tractable” family of distributions over parameters and latent variables
 - Family is indexed by a set of “free parameters”
 - Find member of the family closest to: $p(\theta, z | x)$
call the family Q and want $q \in Q$ that is closest to $p(\theta, z | x)$
- Questions:
 - How do we measure “closeness”?
 - If the posterior is intractable, how can we approximate something we do not have to begin with?

Variational Methods

- Similarity measure:

$$D(q(z, \theta) \parallel p(z, \theta | x)) = E_q[\log q(z, \theta)] - E_q[\log p(z, \theta | x)]$$

$$= E_q[\log q(z, \theta)] - E_q[\log p(z, \theta, x)]$$

$\neq \log p(x)$

- Evidence lower bound (ELBO)

$$\underbrace{\log p(x)}_{\text{const.}} = D(q(z, \theta) \parallel p(z, \theta | x)) + \underbrace{\mathcal{L}(q)}_{\text{add to a const}} \geq \underline{\underline{\mathcal{L}(q)}}$$

$-\mathcal{L}(q)$

- Therefore, minimizing KL is equivalent to maximizing a lower bound on the marginal likelihood:

□ Max $\mathcal{L}(q) = \min D(q \parallel p) = \max$ lower bound of $\log p(x)$

$$\mathcal{L}(q) = E_q[\log p(\theta, z, x)] \neq E_q[\log q(\theta, z)]$$

\leftarrow entropy of q

Task 2: Cluster Documents

■ Setup

- **Input:** Corpus of documents
- **Output:** Topic assignment per document

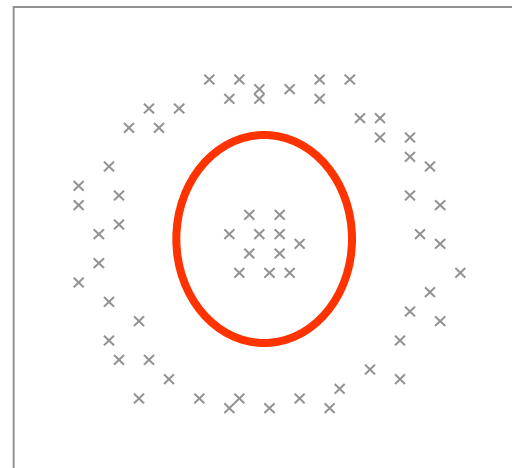
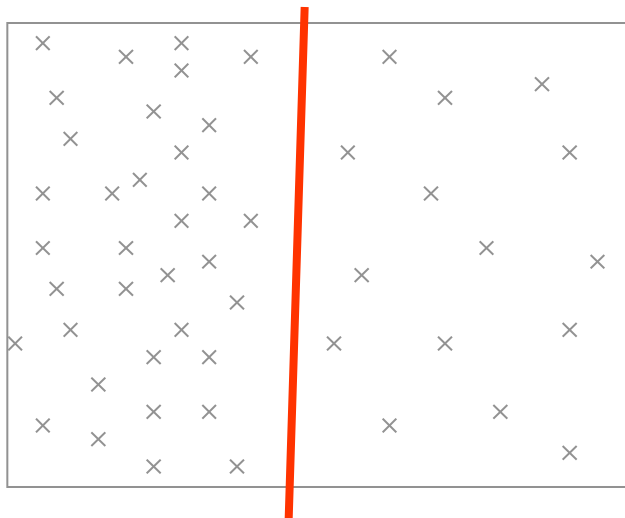
Sports

world news



New Approach: Spectral Clustering

- **Goal:** Cluster observations
- **Method:**
 - Use similarity metric between observations
 - Form a similarity graph
 - Use standard linear algebra and optimization techniques to cut graph into connected components (clusters)



Setup

- Data: x^1, \dots, x^N

$x^i = \text{doc } i \dots$ maybe use tf-idf

- Similarity metric:

S_{ij} bet x^i and x^j (eg cosine similarity $x^i \in \mathbb{R}^V$)

$$S_{ij} = S_{ji}$$

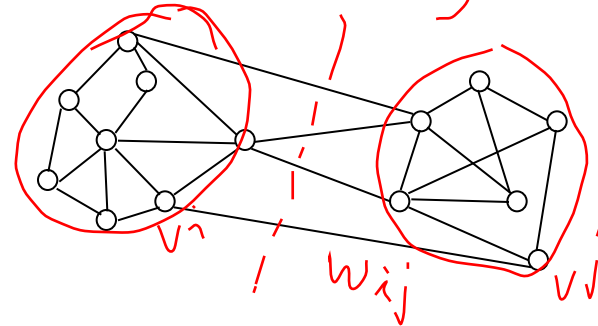
- Similarity graph

□ Nodes

v^i for each x^i

□ Edge weights

$$w_{ij} = \text{fcn of } S_{ij}$$



$$G = \{V, E\}$$

- Problem: Want to partition graph such that edges between groups have low weights

Types of Graphs

■ ϵ -neighborhood:

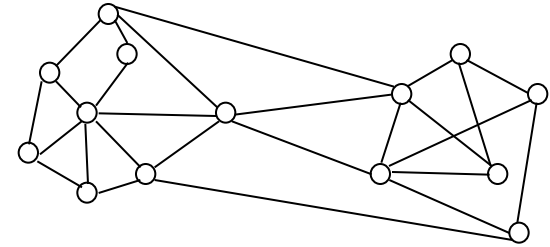
- Only include edges with distances $< \epsilon$
- Treat as unweighted $w_{ij} = \epsilon$

■ k-NN:

- Connect v_i and v_j if v_j is a k-NN of v_i
- Weighted by similarity $s_{ij} = w_{ij}$
- Directed \rightarrow undirected

■ Mutual k-NN:

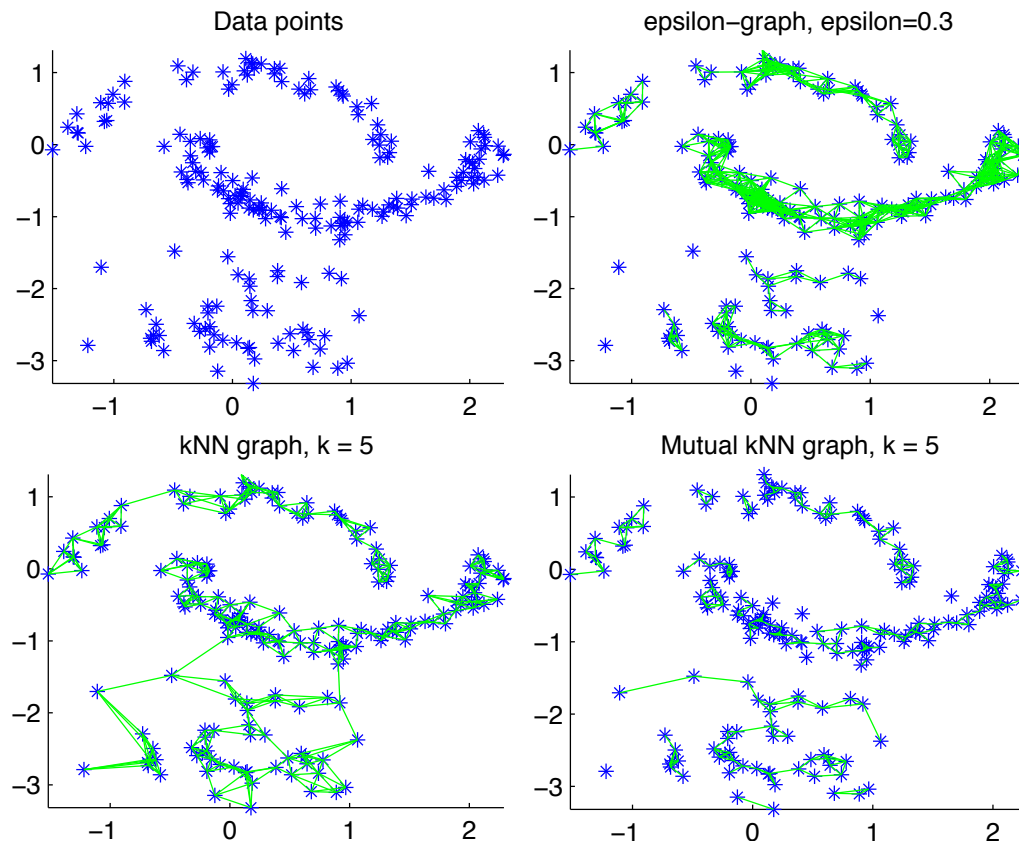
- Same as k-NN, but only include mutual k-NN



Issues with Choosing Graph

- Choosing graph construction techniques and parameters is non-trivial

Choice matters



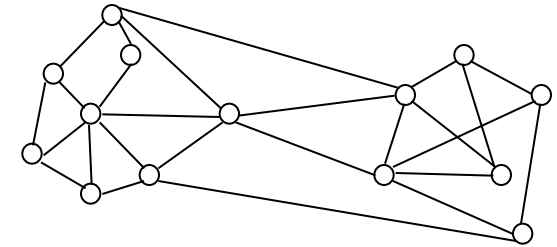
From
von Luxburg
2007

Graph Terminology I

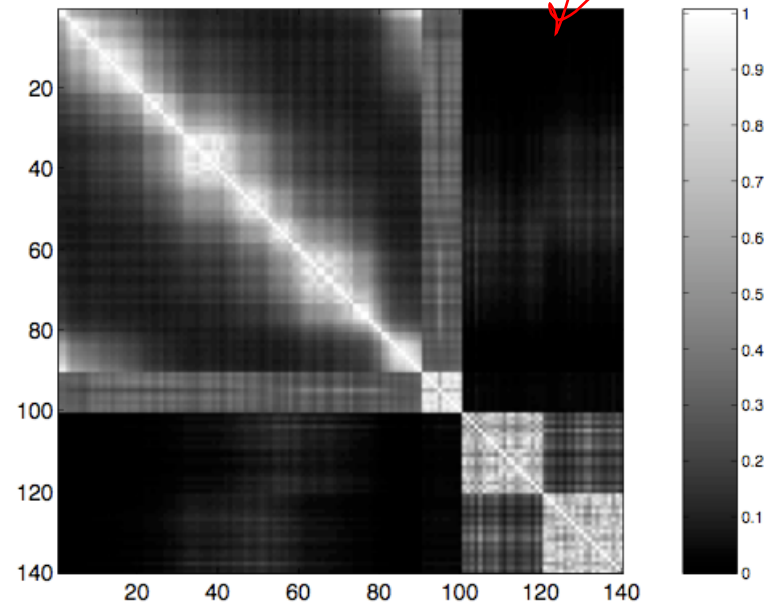
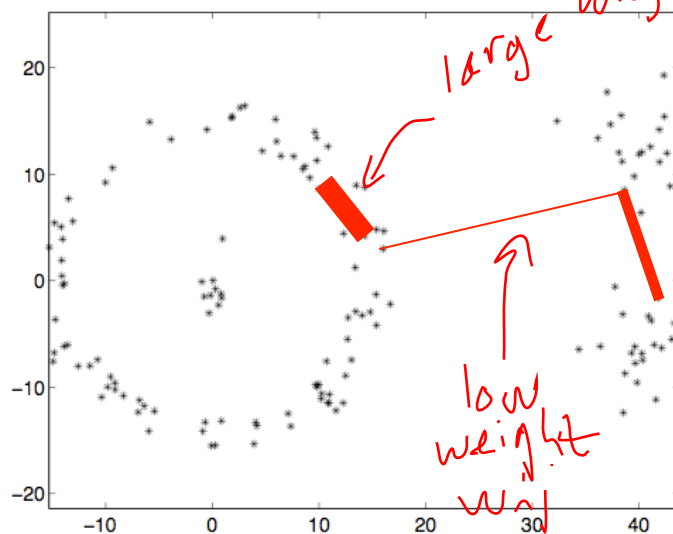
- Weighted adjacency matrix

$$W = (w_{ij})_{i,j=1,\dots,N}$$

$w_{ij} = 0 \Rightarrow$ no edge bt v_i and v_j
 $w_{ij} \geq 0$



W sparse



Graph Cuts

- **Problem:** Partition graph such that edges between groups have low weights

- Define: $\underline{W(A, B)} = \sum_{i \in A, j \in B} w_{ij}$

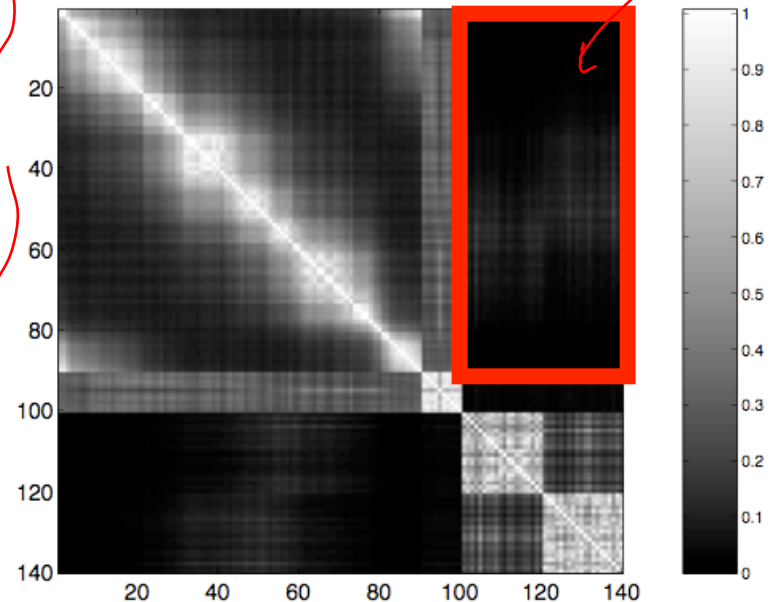
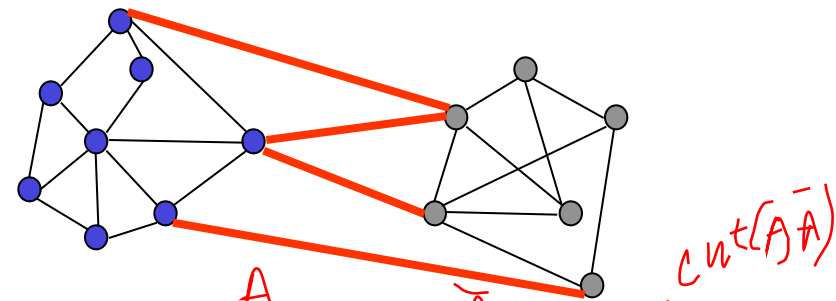
- MinCut problem:

$$\text{Cut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i)$$

Choose

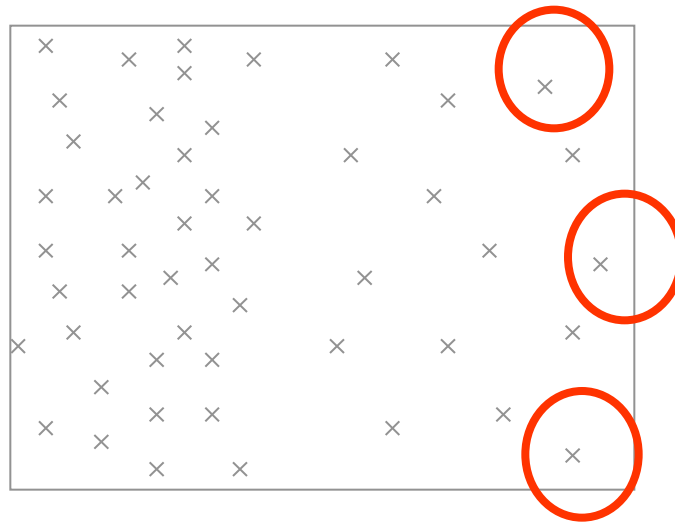
$$A_1, \dots, A_k = \underset{A_1, \dots, A_k \text{ disjoint } \subset V}{\text{argmin}} \text{Cut}(A_1, \dots, A_k)$$

- Trivial to solve for $k=2$



Issues with MinCut

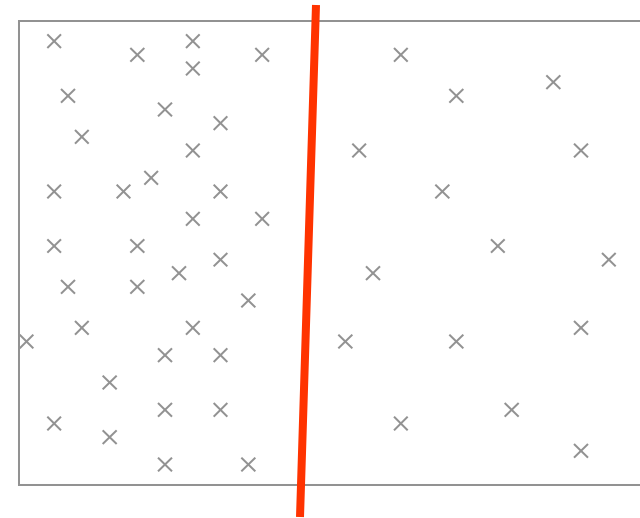
- MinCut favors isolated clusters



nothing working against this

Cuts Accounting for Size

- Ratio cuts (RatioCut)
- Normalized cuts (Ncut)
- Lead to “balanced” clusters



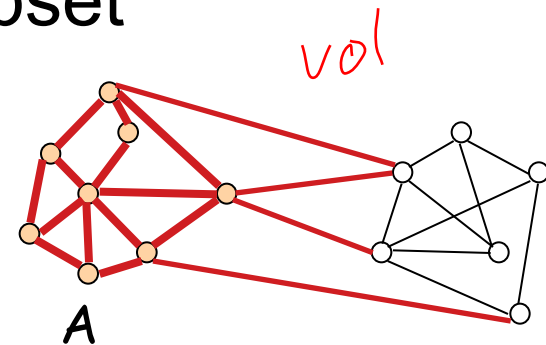
- First need more graph terminology...
to measure "size" of the clusters

Graph Terminology II

- Two measures of size of a subset

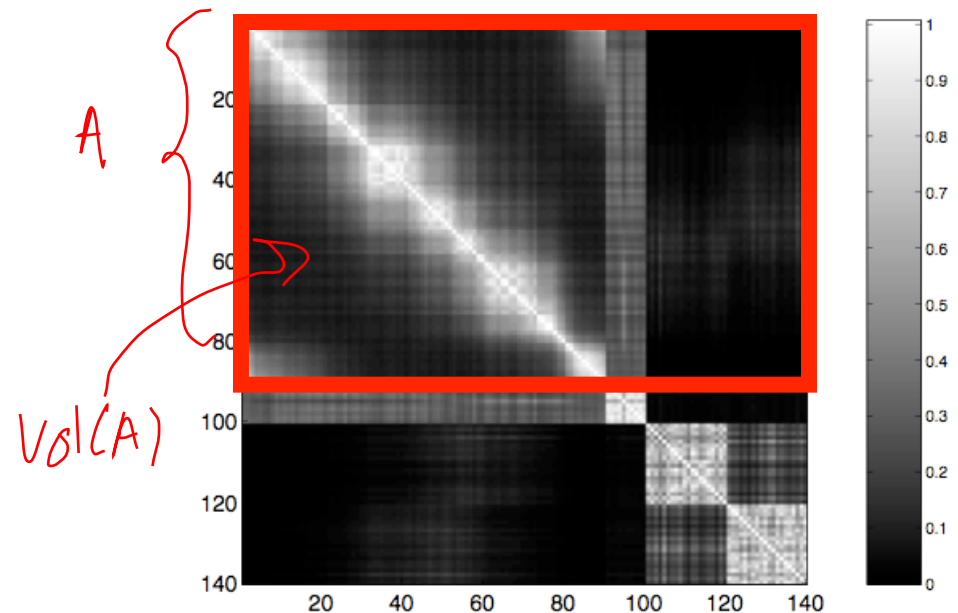
- Cardinality:

$$|A| = \# \text{ of vertices in } A$$



- Volume:

$$\text{vol}(A) = \sum_{i \in A} \sum_{j=1}^N w_{ij}$$



Cuts Accounting for Size

- Ratio cuts (RatioCut)

- $k=2$ $\text{RatioCut}(A, \bar{A}) = \underbrace{\text{cut}(A, \bar{A})}_{\text{min when } |A| \text{ and } |\bar{A}| \text{ coincide}} \left(\frac{1}{|A|} + \frac{1}{|\bar{A}|} \right)$

- General k

$$\text{RatioCut}(A_1, \dots, A_k) = \frac{1}{2} \sum_i \frac{W(A_i, \bar{A}_i)}{|A_i|}$$

min when $|A|$ and $|\bar{A}|$ coincide

- Normalized cuts (Ncut)

- $k=2$ $\text{Ncut}(A, \bar{A}) = \text{cut}(A, \bar{A}) \left(\frac{1}{\text{vol}(A)} + \frac{1}{\text{vol}(\bar{A})} \right)$

- General k

$$\text{Ncut}(A_1, \dots, A_k) = \frac{1}{2} \sum_i \frac{W(A_i, \bar{A}_i)}{\text{vol}(A_i)}$$

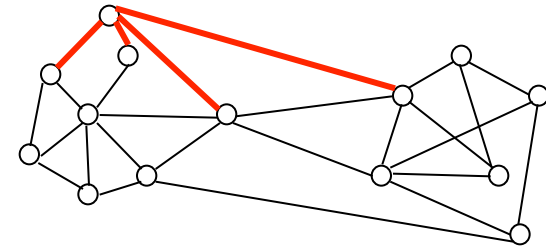
- Problem is NP-hard! Look at relaxation.

Graph Terminology III

- Degree

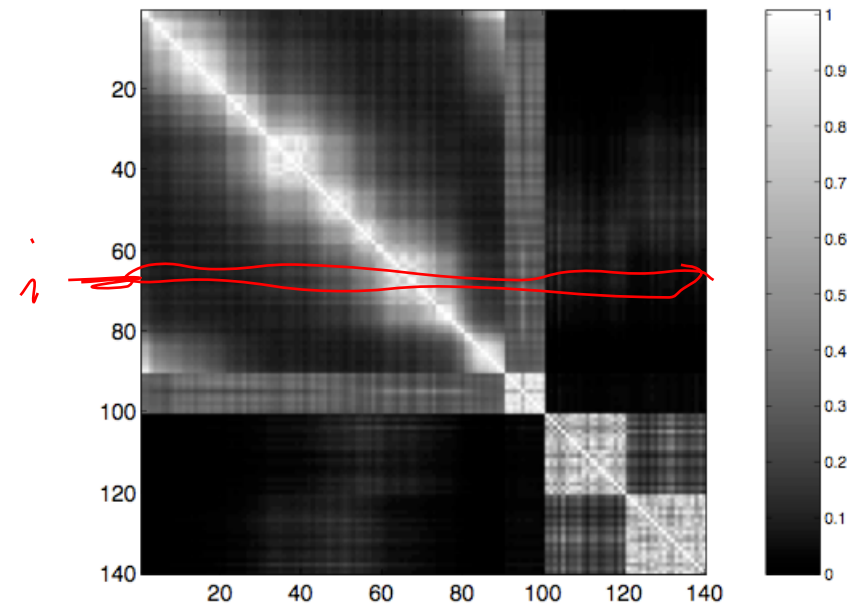
$$d_i = \sum_{j=1}^N w_{ij}$$

only counts neighbors

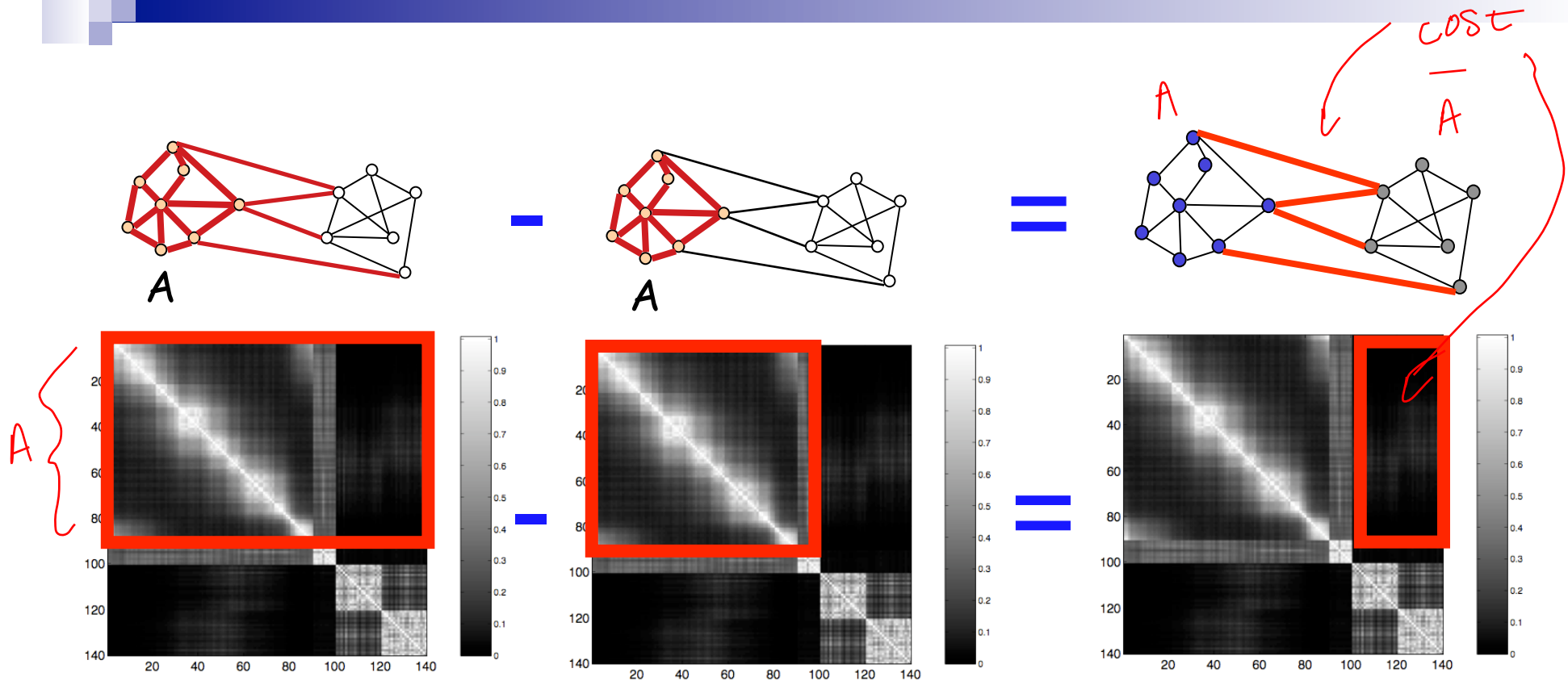


- Degree matrix

$$D = \begin{bmatrix} d_1 & & & 0 \\ & \ddots & & \\ 0 & & & d_N \end{bmatrix}$$



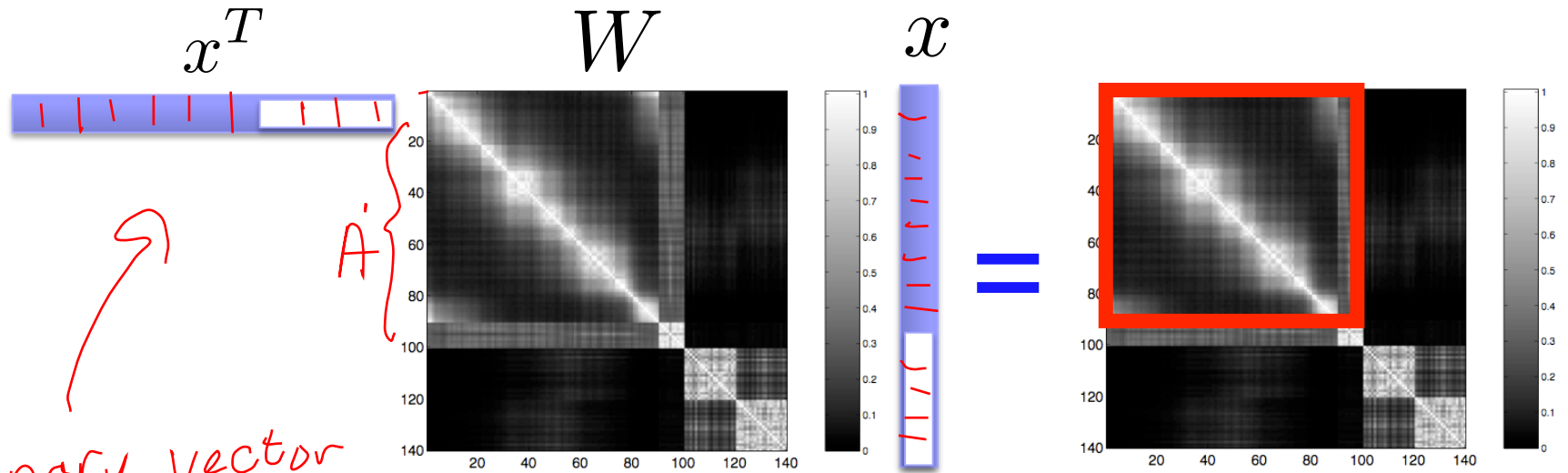
Restating Cut Metric



$$\text{Volume}(A) - \text{"association}(A)\text{"} = \text{cut cost}$$

Restating Cut Metric

- Assoc.



binary vector
that's an
indicator on set A

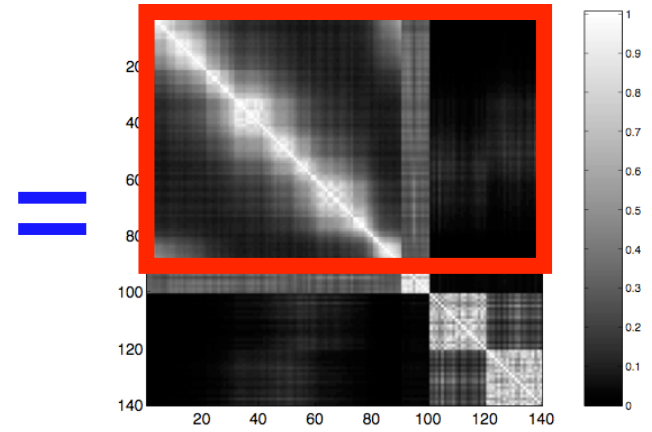
$\mathbb{1}_A$

Restating Cut Metric

- Volume

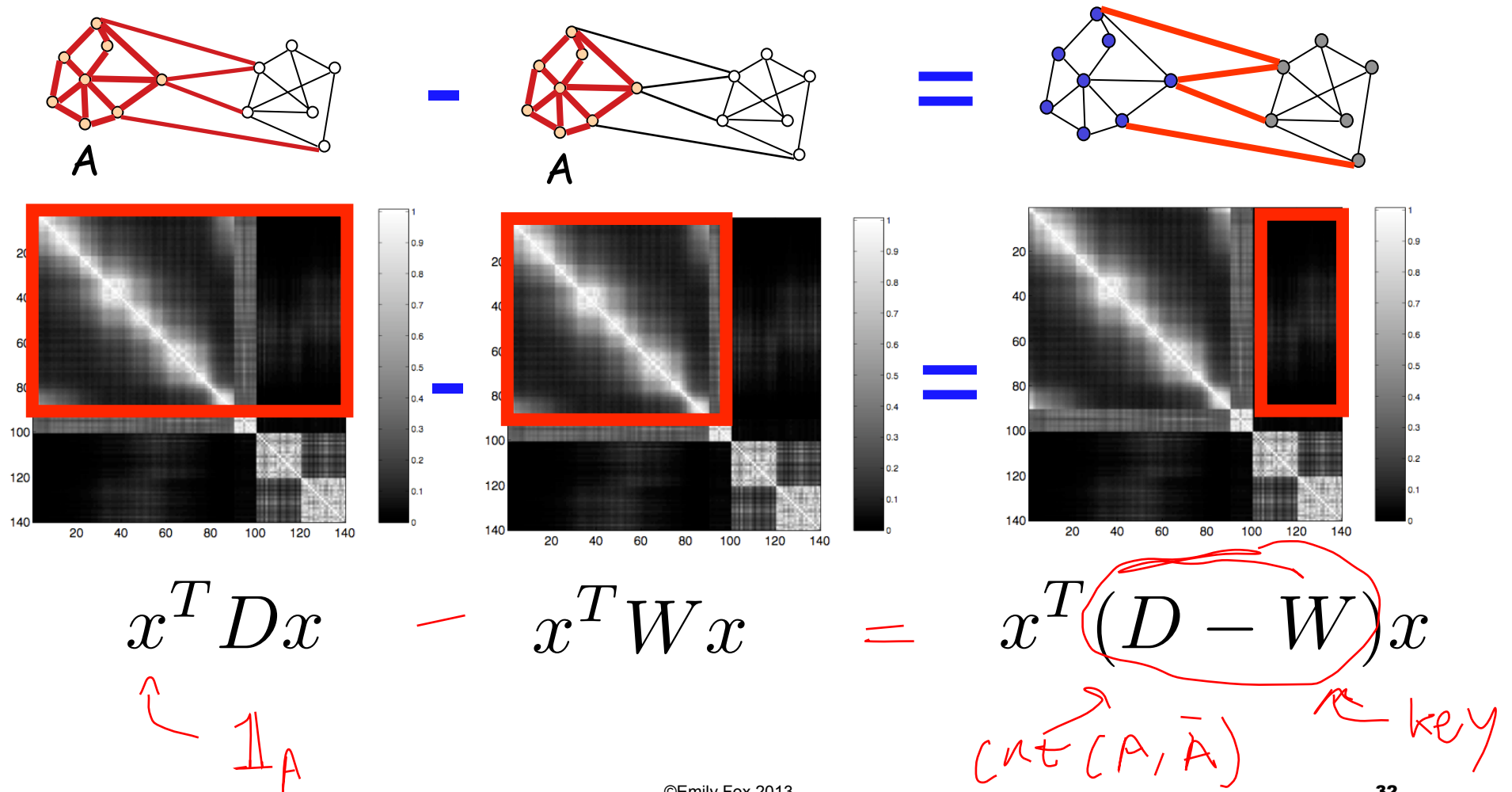
$$x^T \begin{bmatrix} d_1 & 0 & 0 & 0 \\ 0 & d_2 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & d_N \end{bmatrix} x$$

Sum of weights in row 2



Vol(A)

Restating Cut Metric



Graph Laplacian

- Definition: $L = D - W$

- Facts:

- Symmetric, positive semi-definite
- Eigenvalues

$$0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$$

eigvec $u_1 = \mathbb{1}$

- Invariance to self-edges

$$\left. \begin{aligned} L_{ii} &= d_i - w_{ii} \\ L_{ij} &= -w_{ij} \end{aligned} \right\} \text{ don't depend on } w_{ii}$$

- Inner ~~product~~ ^{norm} in L space

$$\forall f \in \mathbb{R}^N \quad f^T L f = \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2 \quad \text{useful later}$$

Relationship to Identifying Connected Components

- Proposition:

- The multiplicity k of eigenvalue 0 of L is equal to the number of connected components

Furthermore, $u_1, \dots, u_k = \mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_k}$

- Proof: Assume graph is connected ($k=1$)

$$0 = u_1^T L u_1 = \sum_{i,j} w_{ij} (u_{1,i} - u_{1,j})^2$$

$$\text{If } w_{ij} > 0 \Rightarrow u_{1,i} = u_{1,j}$$

Since \exists a path bt all i, j , then

$$u_1 = \text{constant} = \mathbb{1}$$

Relationship to Identifying Connected Components

- Proposition:

- The multiplicity k of eigenvalue 0 of L is equal to the number of connected components

- Proof: Assume k connected components A_1, \dots, A_k
 Assume WLOG that they're ordered

$$L = \begin{pmatrix} L_1 & & 0 \\ & \ddots & \\ 0 & & L_k \end{pmatrix}$$

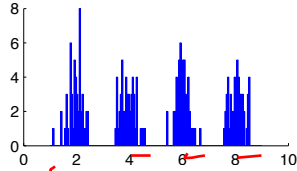
graph Laplacian for subgraph A_1

$\text{eigval}(L) = \cup \text{eigval}(L_i)$ ← each has 1 eigval equal to 0
 + corr. eigvec $\mathbb{1}_{A_i}$
 ⇒ eigvecs are indicators on the partition



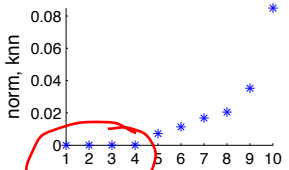
Example – Mixture of Gaussians

Histogram of the sample



← MoG w/ 4 clusters

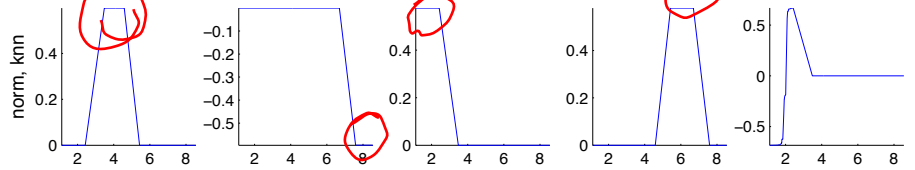
Eigenvalues



4 eigenval = 0

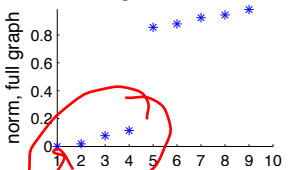
4 disc. comp. in graph

Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



indicators on the partition

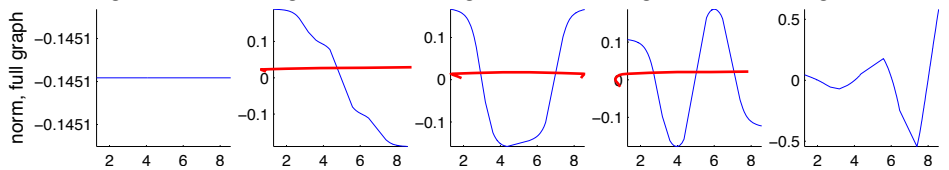
Eigenvalues



1 eigenval = 0

connected graph

Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



still info here

From von Luxburg 2007

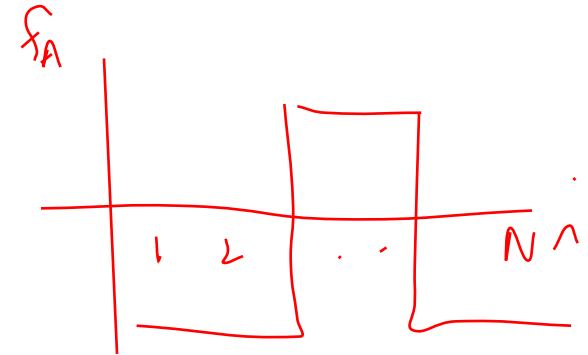
Graph Laplacians and Ratio Cuts

- Ratio cuts for $k=2$

$$V = (A, \bar{A})$$

- Define cluster indicator variables:

$$f_{Ai} = \begin{cases} \sqrt{|\bar{A}|} / |A| & v_i \in A \\ -\sqrt{|A|} / |\bar{A}| & v_i \in \bar{A} \end{cases}$$



- Properties:

$$\sum f_{Ai} = |A| \sqrt{|\bar{A}|} / |A| - |\bar{A}| \sqrt{|A|} / |\bar{A}| = 0$$

$$\|f_A\|^2 = N$$

- RatioCut

$$\text{RatioCut}(A, \bar{A}) = \frac{f_A' L f_A}{|V|} \quad \text{for } f_A \text{ as}$$

- Reformulating RatioCut problem

$$\min_{A \subset V} f_A' L f_A \quad \text{s.t. } f_A \text{ defined as above, } f_A \perp \mathbf{1}, \|f_A\| = \sqrt{N}$$

f_{Ai} are in a discrete set

Relaxation to Formulation

- Let f be arbitrary continuous vector

$$\min_{f \in \mathbb{R}^N} f' L f \quad \text{s.t.} \quad f \perp \mathbb{1} \quad \|f\| = \sqrt{N}$$

\uparrow graph Laplacian \uparrow 1st eigvec of L \uparrow const

- Rayleigh-Ritz Theorem

- Which vector maximizes objective subject to constraint that the vector is orthogonal to the first eigenvector and has bounded norm?

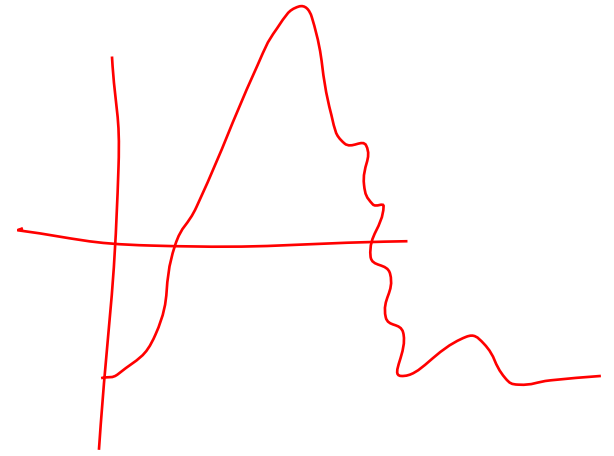
$$f = u_2(L) = \text{eigvec assoc. w/ 2nd smallest eigenval}$$

Mapping Back to Partition

- To obtain partition, transform continuous f to a discrete indicator

- Cluster coordinates

$f_i \in \mathbb{R}$ into C, \bar{C}
using k -means



- Return

$$\begin{cases} v_i \in A & \text{if } f_i \in C \\ v_i \in \bar{A} & \text{if } f_i \in \bar{C} \end{cases}$$

Ratio Cuts for General k

- Define cluster indicator variables:

$$F_{ij} = \begin{cases} 1/\sqrt{|A_j|} \\ 0 \end{cases} \quad \begin{array}{l} \text{v.e. } A_j \\ \text{ow} \end{array} \quad F_A \in \mathbb{R}^{N \times k} \quad \underline{F'_A F_A = I}$$

- RatioCut

$$\text{RatioCut}(A_1, \dots, A_k) = \sum_{i=1}^k f'_{A_i} L f_{A_i} = \underline{\text{Tr}(F'_A L F_A)}$$

- Reformulating RatioCut problem

$$\min_{A_1, \dots, A_k} \text{Tr}(F'_A L F_A) \quad \text{and } F_A \text{ w/ } F'_A F_A = I$$

- Relaxation

$$\min_{F \in \mathbb{R}^{N \times k}} \text{Tr}(F' L F) \quad \text{s.t. } F' F = \underline{I}$$

Ratio Cuts for General k

- Relaxation:

$$\min_{F \in \mathbb{R}^{N \times k}} \text{Tr}(F'LF) \quad \text{s.t.} \quad F'F = I$$

- Solution:

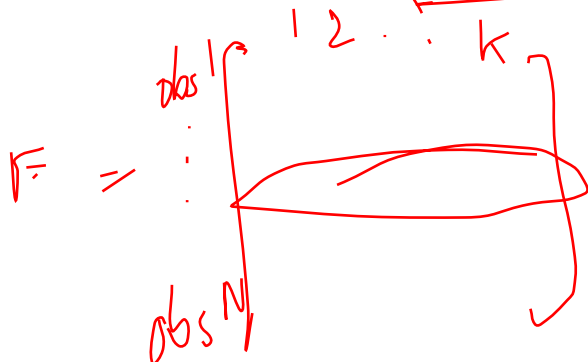
standard trace min problem

⇒ choose F containing first k eigvec(L)

$$\begin{bmatrix} | & & | \\ u_1 & \dots & u_k \\ | & & | \end{bmatrix}$$

- To obtain partition:

Cluster rows of F using k-means



if $obs\ i$ is in cluster w/ $obs\ j$
then the rows are the same

Graph Laplacians and Norm. Cuts

- Normalized cuts for $k=2$
- Define cluster indicator variables:

$$f_{A\bar{A}} = \begin{cases} \sqrt{\text{vol}(\bar{A}) / \text{vol}(A)} & v_i \in A \\ -\sqrt{\text{vol}(A) / \text{vol}(\bar{A})} & v_i \in \bar{A} \end{cases}$$

- Properties:

$$(Df_A)' \mathbb{1} = 0 \quad \text{and} \quad f_A' Df_A = \text{vol}(V)$$

- Ncut

$$\text{Ncut}(A, \bar{A}) = \frac{f_A' L f_A}{\text{vol}(V)}$$

- Reformulating Ncut problem

$$\min_{A \subset V} f_A' L f_A \quad \text{s.t.} \quad Df_A \mathbb{1} = \mathbb{1} \quad \text{and} \quad f_A' Df_A = \text{vol}(V)$$

Relaxation to Formulation

- Let f be arbitrary continuous vector

$$\min_{f \in \mathbb{R}^N} f' L f \quad \text{s.t.} \quad Df \perp \mathbb{1} \quad (f' D f = \text{vol}(V))$$

$$\iff f = D^{-1/2} g$$

$$\min_{g \in \mathbb{R}^n} g' \underbrace{D^{-1/2} L D^{-1/2}}_{\triangleq L_{\text{sym}}} g \quad \text{s.t.} \quad g \perp \underbrace{D^{1/2} \mathbb{1}}_{\substack{\uparrow \\ \text{1st eigvec of } L_{\text{sym}}}} \quad \|g\|^2 = \text{vol}(V) \quad \uparrow \text{const}$$

- Rayleigh-Ritz Theorem

$$g = u_2(L_{\text{sym}})$$

$$\Rightarrow f = D^{-1/2} u_2(L_{\text{sym}}) = u_2(L_{\text{rw}})$$

$$\boxed{\text{equiv to } f \text{ soln of } Lu = \lambda D u} \quad \leftarrow I - D^{-1} W$$

Normalized Cuts for General k

- Define cluster indicator variables:

$$F_{ij} = \begin{cases} 1/\sqrt{\text{vol}(A_j)} & v_i \in A_j \\ 0 & \text{ow} \end{cases} \quad \begin{aligned} F'_{\mathcal{A}} F_{\mathcal{A}} &= I \\ F'_{\mathcal{A}} D F_{\mathcal{A}} &= I \end{aligned}$$

- Reformulating RatioCut problem

$$\min_{A_1, \dots, A_k} \text{Tr}(F'_{\mathcal{A}} L F_{\mathcal{A}}) \quad \text{s.t.} \quad F'_{\mathcal{A}} D F_{\mathcal{A}} = I$$

- Relaxation

$$\min_{H \in \mathbb{R}^{N \times k}} \text{Tr}(H' D^{-1/2} L D^{-1/2} H) \quad \text{s.t.} \quad H' H = I$$

- Solution:

- H is matrix of first k eigenvectors of L_{sym} , which is equivalent to the approximate F being the first k eigenvectors of L_{rw}

Random Walks on Graphs

- Stochastic process with random jumps from v_i to v_j wp:
- Transition matrix:
- Connection to graph Laplacian:
- Intuitively, want to partition graph s.t. random walk stays in cluster for a while and rarely jumps between clusters

Random Walks on Graphs

- Assume that stationary distribution exists and is unique. Then,
- Proposition: $N\text{cut}(A, \bar{A}) = P(A | \bar{A}) + P(\bar{A} | A)$
- Proof:

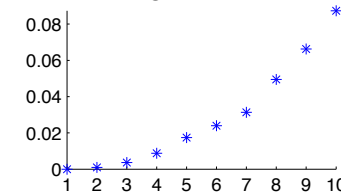
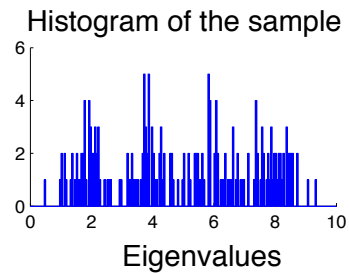
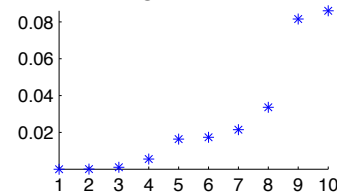
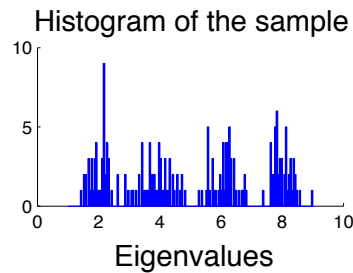
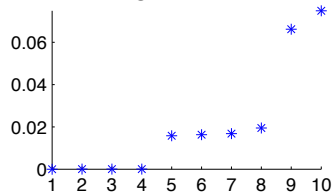
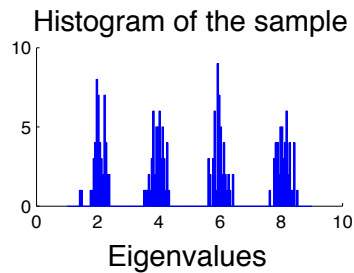
- Minimizing normalized cuts is equivalent to minimizing the probability of transitioning between clusters

Notes

- No guarantee to quality of approximation
- Sensitive to choice of similarity graph (see earlier)
- Which graph Laplacian to use?
 - If degrees in graph vary significantly, then Laplacians are quite different
 - In general, L_{rw} behaves the best
 - Volume gives better measure of within-cluster similarity than cardinality
 - Normalized cuts has consistency results, Ratio cuts does not

Notes

- Choosing the number of clusters k can be hard
 - Easy when clusters are well-separated



From
von Luxburg
2007

- k-means to return partition from solution to relaxation is *an* approach, but not the only