**Case Study 2: Document Retrieval**

# Spectral Clustering

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

Emily Fox
February 12th, 2013

1

---

# Document Retrieval

- **Goal:** Retrieve documents of interest



ARTICLES

2

# Task 1: Find Similar Documents

■ **Setup**
  ☐ **Input:** Query article $X$
  ☐ **Output:** Set of k similar articles
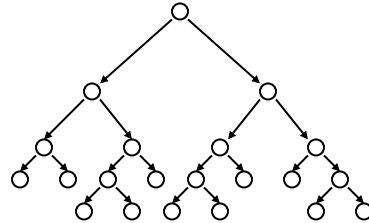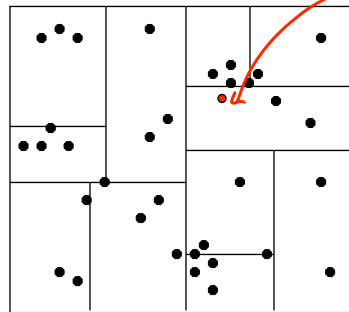


$X$

---

# k-Nearest Neighbor

■ Articles $\quad X = \{x^1, \ldots, x^N\}, \quad x^i \in \mathbb{R}^d$

■ Query: $\quad x \in \mathbb{R}^d$

■ k-NN
  ☐ Goal: Find k articles in $X$ closest $x$

  ☐ Formulation:

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

# Nearest Neighbor with KD Trees

query

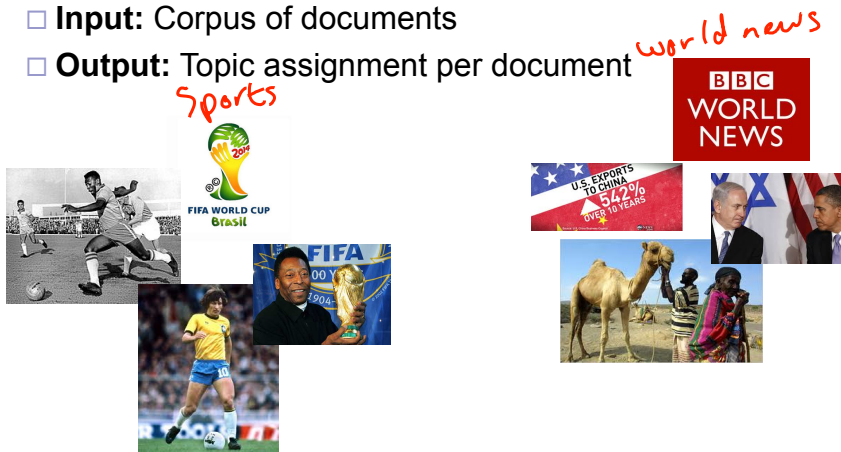- Traverse the tree looking for the nearest neighbor of the query point.

5

# Task 2: Cluster Documents

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

world news

Sports
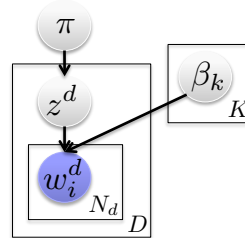
6

# A Generative Model

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



$$z^d \sim \pi \qquad d = 1, \ldots, D$$

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

↖ word prob. for cluster/topic $z^d$

Bayesian approach:

$$\pi \sim Dir(\alpha_1, \ldots, \alpha_K)$$

$$\beta_k \sim Dir(\lambda_1, \ldots, \lambda_V)$$

$\beta_k$ is a $V$-dim pmf ← size of vocab.

©Emily Fox 2013

7

---

# Inference

- **Two tasks**
  - ☐ **Point estimation:**
    - Expectation-Maximization (EM)
  - ☐ **Characterize posterior:**
    - Gibbs sampling
    - Variational methods
    - Stochastic variational inference

©Emily Fox 2013

8

4

# EM Algorithm

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

- **E-Step**

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

- **M-Step**

  Compute $\quad \hat{\theta}^{(t+1)} = \arg\max_{\theta} U(\theta, \hat{\theta}^{(t)})$
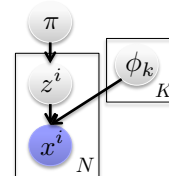
9

# Collapsed Gibbs Sampling

$$\pi \sim \text{Dir}(\alpha_1, \ldots, \alpha_K) \qquad z^i \sim \pi$$
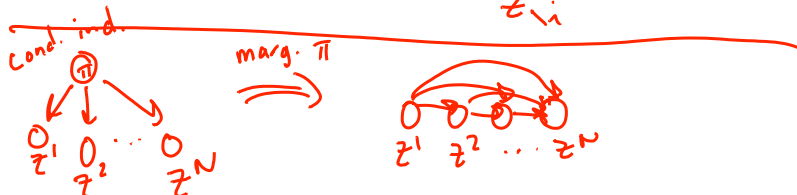$$\{\mu_k, \Sigma_k\} \sim F(\phi) \qquad x^i \mid z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Collapsed sampler

For $i = 1, \ldots, N$

$z^{i(t)} \sim p(z^i \mid z^{1(t)}, \ldots, z^{i-1(t)}, z^{i+1(t)}, \ldots, z^{N(t-1)}, X_{1:N}, \alpha, \lambda)$

$z_{\setminus i}^{(t)}$

cond. ind.

marg. $\pi$

$\pi$

$z^1 \quad z^2 \quad \cdots \quad z^N$

$z^1 \quad z^2 \quad \cdots \quad z^N$

10

5

# Task 3: Mixed Membership Model

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



EDUCATION

FINANCE

TECHNOLOGY

*mixed membership*

---

# Latent Dirichlet Allocation (LDA)



*each doc is a mixture of these corpus-wide topics*

*Topics*

*Documents*

*Topic proportions and assignments*

*each topic as a dist. over words* $\{\beta_k\}$

*every word is assigned to a topic*

*each doc has its own prevalence of topics in doc*

# Variational Methods

- Recall task: Characterize the posterior $p(\theta, z \mid x)$
  
  params ↗ ↑ latent vars → 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: $p(\theta, z \mid x)$
  
  Call the family $Q$ and want $q \in Q$ that is closest to $p(\theta, z \mid x)$

- Questions:
  - How do we measure "closeness"?
  - If the posterior is intractable, how can we approximate something we do not have to begin with?

14

---

# Variational Methods

- Similarity measure:

$$D(q(z,\theta) \| p(z,\theta \mid x)) = E_q[\log q(z,\theta)] - E_q[\log p(z,\theta \mid x)]$$
$$= E_q[\log q(z,\theta)] - E_q[\log p(z,\theta,x)]$$
$$-\mathcal{L}(q) \qquad \ast \log p(x)$$

- Evidence lower bound (ELBO)

$$\log p(x) = D(q(z,\theta) \| p(z,\theta \mid x)) + \mathcal{L}(q) \geq \mathcal{L}(q)$$

const.   add to a const

- Therefore, minimizing KL is equivalent to maximizing a lower bound on the marginal likelihood:
  - Max $\mathcal{L}(q)$ = min $D(q\|p)$ = max lower bound of $\log p(x)$   ← entropy of q

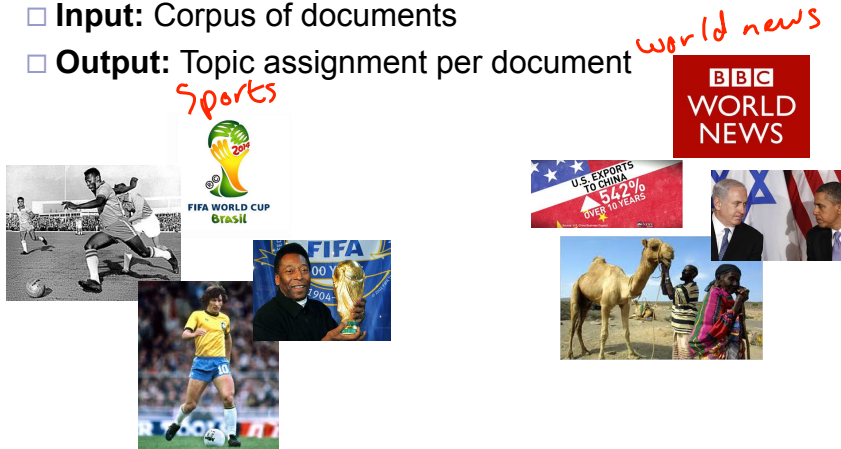$$\mathcal{L}(q) = E_q[\log p(\theta, z, x)] \ast E_q[\log q(\theta, z)]$$

15

7

# Task 2: Cluster Documents

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

*world news*

*Sports*

---
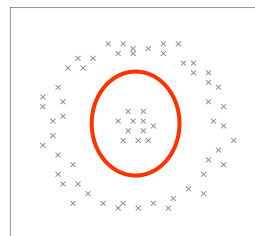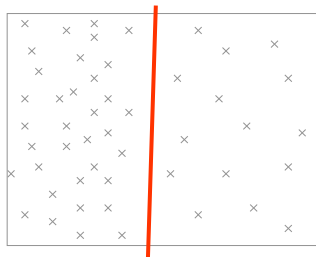
# New Approach: Spectral Clustering

- **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)
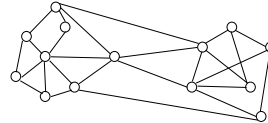
# Setup

- Data: $x^1, \ldots, x^N$
- Similarity metric:

- Similarity graph
  - Nodes
  - Edge weights

**G = {V, E}**

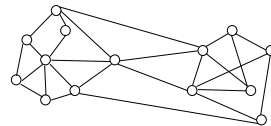- Problem: Want to partition graph such that edges between groups have low weights

# Types of Graphs

- **ε-neighborhood**:
  - Only include edges with distances < ε
  - Treat as unweighted

- **k-NN:**
  - Connect $v_i$ and $v_j$ if $v_j$ is a k-NN of $v_i$
  - Weighted by similarity $s_{ij}$
  - Directed → undirected

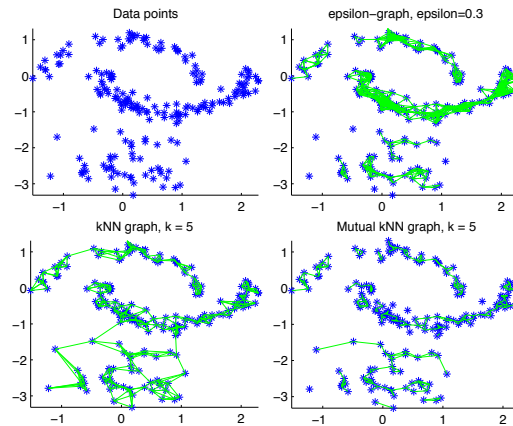- **Mutual k-NN:**
  - Same as k-NN, but only include mutual k-NN

# Issues with Choosing Graph

- Choosing graph construction techniques and parameters is non-trivial



From
von Luxburg
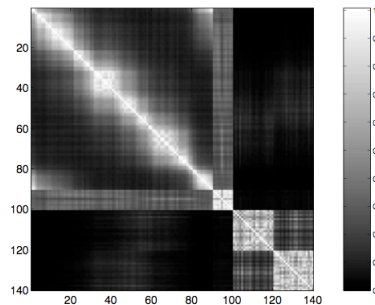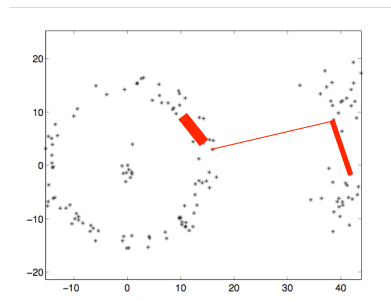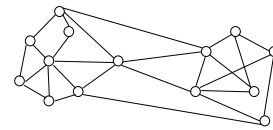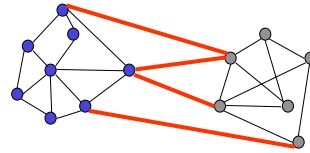2007

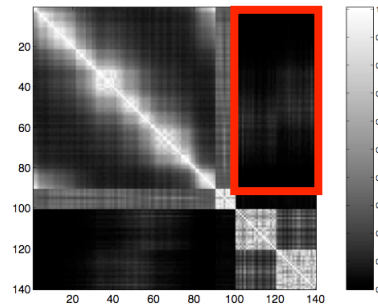# Graph Terminology I

- Weighted adjacency matrix

# Graph Cuts

- **Problem:** Partition graph such that edges between groups have low weights
- Define: $W(A, B) = \sum_{i \in A, j \in B} w_{ij}$



- MinCut problem:
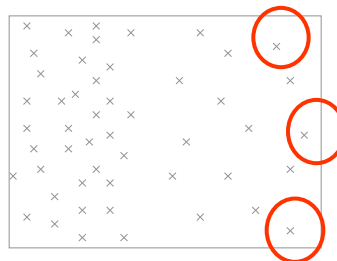


- Trivial to solve for *k*=2

23

# Issues with MinCut

- MinCut favors isolated clusters

24

# Cuts Accounting for Size

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



- First need more graph terminology…

# Graph Terminology II

- Two measures of size of a subset
  - □ Cardinality:

  $|A|$

  □ Volume:

  $\text{vol}(A)$

# Cuts Accounting for Size

- Ratio cuts (RatioCut)
  - $k$=2

  - General $k$


- Normalized cuts (Ncut)
  - $k$=2

  - General $k$


- Problem is NP-hard!  Look at relaxation.

---

# Graph Terminology III

- Degree



- Degree matrix

# Restating Cut Metric

# Restating Cut Metric

$$x^T \qquad W \qquad x$$

14

# Restating Cut Metric

$$x^T \qquad D \qquad x$$

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

---

# Restating Cut Metric

A                    A

$$x^T D x \qquad x^T W x \qquad x^T(D - W)x$$

# Graph Laplacian

- Definition:


- Facts:
  - □ Symmetric, positive semi-definite
  - □ Eigenvalues


  - □ Invariance to self-edges


  - □ Inner product in *L* space

33

# Relationship to Identifying Connected Components

- Proposition:
  - □ The multiplicity *k* of eigenvalue 0 of *L* is equal to the number of connected components


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

34

# 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

---

# Example – Mixture of Gaussians



From
von Luxburg
2007

# Graph Laplacians and Ratio Cuts

- Ratio cuts for *k*=2
- Define cluster indicator variables:



- Properties:



- RatioCut


- Reformulating RatioCut problem

# Relaxation to Formulation

- Let *f* be arbitrary continuous vector




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

# Mapping Back to Partition

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

- Cluster coordinates

- Return

# Ratio Cuts for General k

- Define cluster indicator variables:
$$F_{ij} = \begin{cases} 1/\sqrt{|A_j|} \\ 0 \end{cases} \qquad\qquad F'_{\mathcal{A}}F_{\mathcal{A}} = I$$

- RatioCut
$$\mathrm{RatioCut}(A_1, \ldots, A_k) = \sum_{i=1}^{k} f'_{\mathcal{A}i} L f_{\mathcal{A}i} = \mathrm{Tr}(F'_{\mathcal{A}} L F_{\mathcal{A}})$$

- Reformulating RatioCut problem
$$\min_{A_1, \ldots, A_k} \mathrm{Tr}(F'_{\mathcal{A}} L F_{\mathcal{A}})$$

- Relaxation
$$\min_{F \in R^{N \times k}} \mathrm{Tr}(F' L F)$$

# Ratio Cuts for General k

- Relaxation:

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

- Solution:

- To obtain partition:

# Graph Laplacians and Norm. Cuts

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

- Properties:

- Ncut

- Reformulating Ncut problem

# Relaxation to Formulation

- Let *f* be arbitrary continuous vector



- Rayleigh-Ritz Theorem

---

# Normalized Cuts for General k

- Define cluster indicator variables:

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

- Reformulating RatioCut problem

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

- Relaxation

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

- 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: $\mathrm{Ncut}(A, \bar{A}) = P(A \mid \bar{A}) + P(\bar{A} \mid A)$

- Proof:

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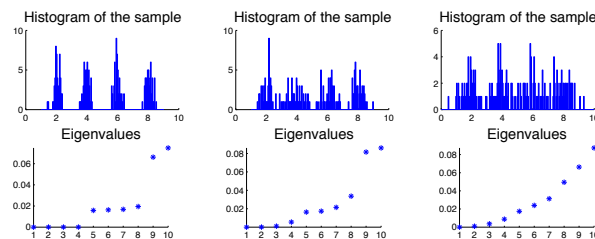
# Notes

- 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