

## Nonlinear dimensionality reduction (review)

How to map inputs  $\vec{x}_i \in \mathbb{R}^D$  to outputs  $\vec{y}_i \in \mathbb{R}^d$  with  $d \ll D$ ?

- Approach #1: probabilistic generative modeling
  - e.g. mixture of factor analyzers
  - strong parametric assumptions
  - latent (hidden) variables
  - parameter estimation by EM algorithm
  - potential local maxima in log-likelihood
- Approach #2: spectral methods
  - e.g. isomap, maximum variance unfolding
  - non-parametric, graph-based
  - low dimensional outputs from metric multidimensional scaling
  - tractable optimizations; no spurious local optima

## Clustering

How to map inputs  $\vec{x}_i \in \mathbb{R}^D$  to discrete outputs  $y \in \{1, 2, \dots, k\}$ ?

- Earlier in course: k-means (VQ), Gaussian mixture models
- Today: spectral clustering

Same distinctions arise as in NLDR.

For simplicity, focus on  $k=2$  binary partitioning.

(Finer clusters can be found by recursive subdivision.)

• Basic algorithm ( $k=2$ )

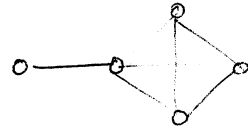
Step 1: graph construction

Define affinity matrix  $W_{ij} = e^{-\frac{1}{2\sigma^2} \|\vec{x}_i - \vec{x}_j\|^2}$

Define diagonal matrix  $D_{ii} = \sum_j W_{ij}$

Intuition:

- each input  $\vec{x}_i$  is node in graph
- nearby inputs connected by larger weights  $W_{ij}$
- row sums  $D_{ii}$  count "effective" # neighbors within radius  $\sigma$
- "sparsify" graphs by approximating  $W_{ij} = 0$  for  $\|\vec{x}_i - \vec{x}_j\|^2 \gg \sigma^2$



Step 2: constrained optimization over real-valued  $\vec{y} \in \mathbb{R}^N$

Minimize  $C(\vec{y}) = \frac{1}{2} \frac{\sum_{i,j} W_{ij} (y_i - y_j)^2}{\sum_i D_{ii} y_i^2}$  subject to  $\sum_i y_i D_{ii} = 0$

Numerator: favors  $y_i \approx y_j$  for large  $W_{ij}$

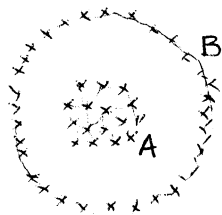
Constraint: prevents degenerate solution  $y_i = \text{constant}$  for all  $i$

Denominator: repels  $y_i$  from origin; more neighbors  $\rightarrow$  greater repulsion

Step 3: quantization

Threshold  $y_i$  on (say) mean or median to obtain  $k=2$  clusters.

Example:



Suppose  $W_{AB} \ll W_{AA}$

Then  $y_i \approx \begin{cases} +1 & \text{for } i \in A \\ -1 & \text{for } i \in B \end{cases}$   
will minimize  $C(\vec{y})$ .

But: how to minimize  $C(\vec{y})$  in step #2?

## Detour: generalized eigenvalue problem

Let  $A$  and  $B$  be square, symmetric matrices.

Also assume  $B$  is positive definite

- generalized Rayleigh quotient

$$Q(y) = \frac{y^T A y}{y^T B y}$$

To minimize:

$$\frac{\partial Q}{\partial y} = 0 \rightarrow 2 \left( \frac{A y}{y^T B y} \right) - \frac{y^T A y}{(y^T B y)^2} (2 B y) = 0$$

$$A y = \left( \frac{y^T A y}{y^T B y} \right) B y$$

$$(B^{-1} A) y = \underbrace{Q(y)}_{\text{scalar}} y \quad \text{this is an eigenvalue equation for } B^{-1} A$$

Solutions of  $\frac{\partial Q}{\partial y} = 0$  are eigenvectors of  $B^{-1} A$ .

$\min(Q(y)) =$  smallest eigenvalue of  $B^{-1} A$

$\operatorname{argmin} Q(y) =$  corresponding eigenvectors

- Symmetrized form: let  $z = B^{\frac{1}{2}} y$

$$\text{Then } Q(z) = \frac{(z^T B^{-\frac{1}{2}}) A (B^{-\frac{1}{2}} z)}{(z^T B^{-\frac{1}{2}}) B (B^{-\frac{1}{2}} z)} = \frac{z^T (B^{-\frac{1}{2}} A B^{-\frac{1}{2}}) z}{z^T z}$$

$$\text{To minimize: } \frac{\partial Q}{\partial z} = 0 \rightarrow (B^{-\frac{1}{2}} A B^{-\frac{1}{2}}) z = Q(z) z$$

Let  $z_\alpha$  denote  $\alpha$ th eigenvector of  $B^{-\frac{1}{2}} A B^{-\frac{1}{2}}$ , ...

Then  $z_\alpha^T z_\beta = 0$  for  $\alpha \neq \beta$

Substituting  $z = B^{\frac{1}{2}} y$  gives "generalized" orthogonality condition:

$$y_\alpha^T B y_\beta = 0 \text{ for } \alpha \neq \beta \text{ (different eigenvectors of } B^{-1} A)$$

↑ not symmetric

# Optimization for spectral clustering

- Cost function

$$C(y) = \frac{1}{2} \frac{\sum_{i,j} W_{ij} (y_i - y_j)^2}{\sum_i D_{ii} y_i^2} = \frac{\frac{1}{2} \left[ \sum_{i,j} W_{ij} (y_i^2 + y_j^2) - 2 \sum_{i,j} W_{ij} y_i y_j \right]}{\sum_i D_{ii} y_i^2}$$
$$= \frac{y^T (D - W) y}{y^T D y}$$

generalized eigenvalue problem  
 $A = D - W$   
 $B = D$

- Constraint

min  $C(y)$  subject to  $\sum y_i D_{ii} = 0$

$$\sum_i y_i D_{ii} = 0 \implies y^T D \mathbf{1} = 0 \text{ where } \mathbf{1} = \underbrace{(1, 1, 1, \dots, 1)}_{N \text{ ones}}^T$$

- Generalized eigenvalue equation

$$\frac{\partial C}{\partial y} = 0 \implies (D - W)y = \lambda D y$$

- Unconstrained minimum of  $C(y)$

Trivially:  $\min C(y) = 0$  by setting  $y_i = y_j = 1$

From eigenvalue equation:  $(D - W)\mathbf{1} = 0$  since  $D_{ii} = \sum_j W_{ij}$

Hence: smallest eigenvalue  $\lambda_0 = 0$  with eigenvector  $y_0 = \mathbf{1}$

However:  $y_0$  violates constraint b/c  $\mathbf{1}^T D \mathbf{1} = \sum_i D_{ii} \neq 0$ .

- Constrained minimum of  $C(y)$

Let  $y_1$  denote eigenvector of  $D^{-1}(D - W)$  with second smallest eigenvalue  $\lambda_1$ .

Also true that  $\frac{\partial C}{\partial y} \Big|_{y_1} = 0$ .

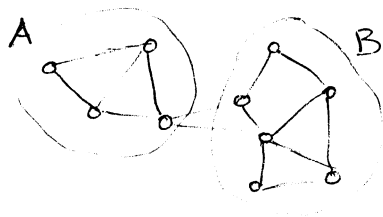
By orthogonality:  $y_1^T D y_0 = 0 \implies y_1^T D \mathbf{1} = 0$  satisfies constraint.

In sum:

$$\underset{y^T D \mathbf{1} = 0}{\operatorname{argmin}} [C(y)] = y_1 \text{ with } \min C(y) = \lambda_1.$$

This optimization has appeared in many different places...

## Graph partitioning



$$\text{cut}(A, B) = \sum_{i \in A} \sum_{j \in B} W_{ij}$$

"cut" of weighted graph  $W_{ij}$  into partitions A and B

• Minimizing  $\text{cut}(A, B)$  over partitions A and B can be done in poly-time. But results often yield highly unbalanced cuts.

• Normalized cut

$$N\text{cut}(A, B) = \frac{\text{cut}(A, B)}{|A|} + \frac{\text{cut}(A, B)}{|B|}$$

$$\text{where } |A| = \sum_{i \in A} D_{ii} \text{ and } |B| = \sum_{i \in B} D_{ii}$$

Penalizes unbalanced cuts.

Favors  $|A| \approx |B|$  with "equal mass" partitions.

However  $\pm$  minimizing  $N\text{cut}(A, B)$  is NP-complete.

• Eigenvalue relaxation

$$\text{Let } y_i = \begin{cases} 1 & \text{if } i \in A \\ 0 & \text{if } i \in B \end{cases}$$

$$N\text{cut}(A, B) = \frac{\sum_{ij} W_{ij} y_i (1 - y_j)}{\sum_i D_{ii} y_i} + \frac{\sum_{ij} W_{ij} y_i (1 - y_j)}{\sum_i D_{ii} (1 - y_i)}$$

$$\text{Let } b \triangleq \frac{\sum_i D_{ii} (y_i)}{\sum_i D_{ii} (1 - y_i)} \text{ ratio of } \frac{|A|}{|B|}$$

$$z_i \triangleq y_i - b(1 - y_i) \in \{1, -b\}$$

By construction:  $z^T D \mathbf{1} = \sum_i D_{ii} [y_i - b(1-y_i)] = 0$ .

It can be shown:

$$\min_{y \in \{0,1\}} \text{Ncut}(y) = \min_z \frac{z^T (D-W) z}{z^T D z} \text{ subject to } \begin{cases} z_i \in \{1, -b\} \\ z^T D \mathbf{1} = 0 \end{cases}$$

Spectral clustering relaxes  $z_i \in \mathbb{R}$  to be real-valued, then thresholds to obtain a discrete solution.