# Graph (1)

- Consist of a set of nodes and a set of edges between those nodes

  - Most important model for applied mathematics

| continuous | discrete |
|---|---|
| function | vector |
| derivative | difference |
| integral | sum |
| calculus | linear algebra |

- Incidence matrix A (mxn)

$(m \text{ edges and } n \text{ nodes})$

edge $i = \text{row } i$, node $j \rightarrow k : -1$ in column $j$, $+1$ in column $k$

$N(\mathbf{A})$ contains all constant vectors: $\mathbf{x} = (c, c, \ldots, c)$

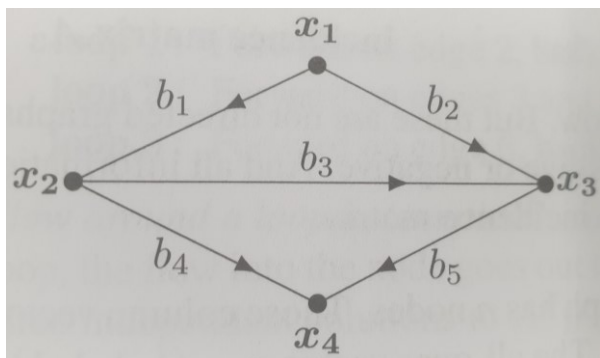$\dim N(\mathbf{A}) = 1$, $\dim C(\mathbf{A}) = \dim C(\mathbf{A}^T) = n - 1$, $\dim N(\mathbf{A}^T) = m - (n - 1)$

row space contains all constant vectors $\mathbf{x}$ with $x_1 + x_2 + \cdots + x_n = 0 (\mathbf{x} \perp \mathbf{1})$

bases $\begin{cases} N(\mathbf{A}): \text{ constant vector } \mathbf{1} \\ C(\mathbf{A}^T): (n-1) \text{ rows of } \mathbf{A} \text{ that produce a tree in the graph (a tree has no loop)} \\ C(\mathbf{A}): \text{ any } (n-1) \text{ columns of } \mathbf{A} \\ N(\mathbf{A}^T): \text{ flows around the } (m-n+1) \text{ small loops in the graph} \end{cases}$
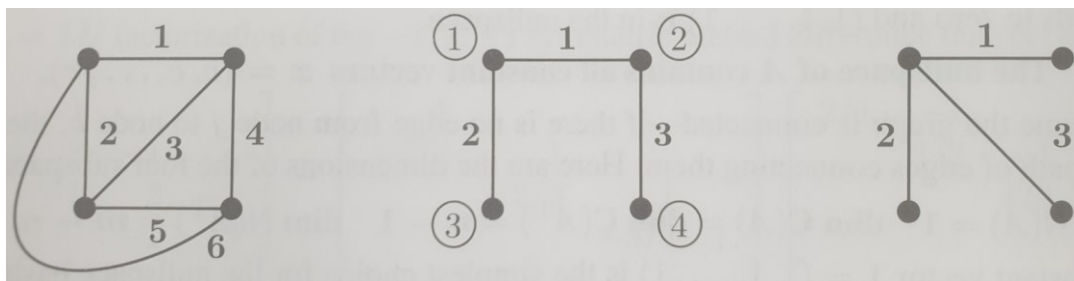
# Graph (2)

- ## Example: m=5, n=4



$$\mathbf{A} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 1 \end{bmatrix} \begin{matrix} \text{edges} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix}$$

$$\text{nodes} \quad 1 \quad 2 \quad 3 \quad 4$$

- ## Graph Laplacian matrix L=A$^T$A
  - Symmetric, positive semidefinite

$$\mathbf{L} = \mathbf{A}^T\mathbf{A} = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 3 & -1 \\ 0 & -1 & -1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & & & \\ & 3 & & \\ & & 3 & \\ & & & 2 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} = \mathbf{D} - \mathbf{B}$$

degree matrix
count edges into node

adjacency matrix
$b_{jk} = 1$: edge from $j$ to $k$

# Graph (2)

- Complete graph: every pair of nodes is connected by an edge, D=(n-1)I, B=all-ones minus I
- Tree: there are no loops in the connected graph



$$\mathbf{A}_1 = \begin{bmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 1 \\ -1 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_3 = \begin{bmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ -1 & 0 & 0 & 1 \end{bmatrix}$$

$$\underbrace{\phantom{xxxxxxxxxxxx}}_{\text{tree: } (n-1)}$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxx}}_{\text{complete graph: } m=\frac{1}{2}n(n-1)} \qquad \rightarrow \text{any graph: } (n-1) \leq m \leq \frac{1}{2}n(n-1)$$

# Kirchhoff's Current Law: $A^\mathsf{T}y=f$

- KCL = balance of currents (forces, money)
  - Flow into each node equals flow out from that node
  - Key to solving $A^\mathsf{T}y=0$ is to look at the small loops in the graph
  - (m-n+1) independent solutions
  - (number of nodes) – (number of edges) + (number of loops) = 1

$$\mathbf{A}_1 = \begin{bmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 1 \\ -1 & 0 & 0 & 1 \end{bmatrix} \rightarrow \underbrace{\mathbf{A}_1\mathbf{x} = \mathbf{0}}_{\substack{\mathbf{x}=(1,1,1,1) \\ \text{not interesting!}}}, \; \mathbf{A}_1{}^T\mathbf{y} = \begin{bmatrix} -1 & -1 & 0 & 0 & 0 & -1 \\ 1 & 0 & -1 & -1 & 0 & 0 \\ 0 & 1 & 1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \mathbf{0} \rightarrow \mathbf{y}_1 = \begin{bmatrix} -1 \\ 1 \\ -1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \; \mathbf{y}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ -1 \\ 1 \\ 0 \end{bmatrix}, \; \mathbf{y}_3 = \begin{bmatrix} 0 \\ -1 \\ 0 \\ 0 \\ -1 \\ 1 \end{bmatrix}$$

outer loop $= \mathbf{y}_1 + \mathbf{y}_2 + \mathbf{y}_3$

$\mathbf{A}_2{}^T\mathbf{y} = 0 \rightarrow \mathbf{y} = 0$

# AᵀCA Framework in Applied Mathematics

- Graphs are perfect examples for three equations in engineering, science, economics
  - Describe a system in steady state equilibrium
  - Balance laws: conservation of charge, balance of force, zero net income in economics, conservation of mass and energy, continuity of every kind

$$\begin{cases} \text{voltages } \mathbf{x} = (x_1, x_2, x_3, x_4) \text{ at the four nodes} \\ \text{currents } \mathbf{y} = (y_1, y_2, y_3, y_4, y_5, y_6) \text{ atalong the six edges} \end{cases}$$

$$\begin{cases} \text{Voltage differences across edges} \quad \mathbf{e} = \mathbf{Ax} \qquad e_1 = (\text{voltage at end node } 2) - (\text{voltage at end node } 1) \\ \text{Ohm's law on each edge} \quad \mathbf{y} = \mathbf{Ce} \qquad \text{current } y_1 = c_1 \text{ times } e_1 = (\text{conductance})(\text{voltage}) \\ \text{Kirchhoff's Law with current sources} \quad \mathbf{f} = \mathbf{A}^T\mathbf{y} \quad \text{current sources } \mathbf{f} \text{ into nodes balance the internal currents } \mathbf{y} \end{cases}$$

$$\rightarrow \mathbf{A}^T\mathbf{CAx} = \mathbf{f} \rightarrow \mathbf{Kx} = \mathbf{f}$$

$\mathbf{K}$ : symmetric, positive $\color{red}{\text{semi}}$definite $\xrightarrow[\substack{x_4=0 \\ n-1=3 \text{ unknown voltages}}]{\text{boundary condition}}$ $\underbrace{\text{reduced } \mathbf{K}}_{(3\times6)(6\times6)(6\times3)}$ : symmetric, invertible, positive definite

energy : $\mathbf{x}^T (\mathbf{A}^T\mathbf{CAx}) = (\mathbf{Ax})^T \mathbf{C}(\mathbf{Ax}) > 0$ if $\mathbf{x} \neq \mathbf{0}$

# A^TCA Framework in Applied Mathematics

- Liner regression: least squares applied to Ax=b

$$\begin{cases} \mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b} : \text{ Normal equation for the vector } \hat{\mathbf{x}} \text{ that best fits the data } \mathbf{b} \\ \mathbf{A}^T \mathbf{C} \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{C} \mathbf{b} : \text{ Least squares weighted by the inverse covariance matrix } \mathbf{C} = \mathbf{V}^{-1} \\ \min \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{\mathbf{C}}^2 : \text{ Minumum squared error } (\mathbf{b} - \mathbf{A}\mathbf{x})^T \mathbf{C} (\mathbf{b} - \mathbf{A}\mathbf{x}) \end{cases}$$

- Graph Laplacian Matrix

$$\mathbf{K} = \mathbf{A}^T \mathbf{C} \mathbf{A} : \text{ weighted graph Laplacian, } \mathbf{G} = \mathbf{A}^T \mathbf{A} : \text{ standard Laplacian} (\mathbf{C} = \mathbf{I})$$

$$\mathbf{A}^T \mathbf{A} = (\text{diagonal}) + (\text{off-diagonal}) = (\text{degree matrix}) - (\text{adjacency matrix}) = \mathbf{D} - \mathbf{B}$$
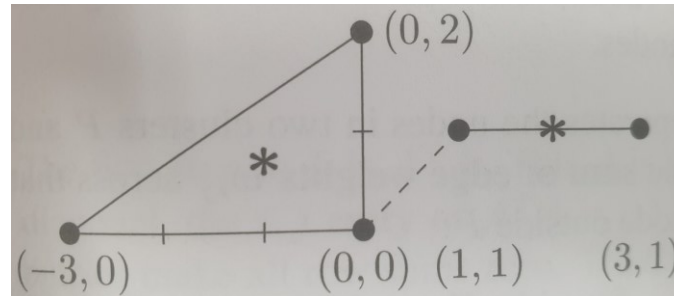
$$\begin{cases} \text{every row and column of } \mathbf{G} \text{ and } \mathbf{K} \text{ adds to zero because } \mathbf{x} = (1,\ldots,1) \text{ has } \mathbf{A}\mathbf{x} = \mathbf{0} \\ \mathbf{G} = \mathbf{A}^T \mathbf{A} \text{ is symmetric because edges go both ways (undirected graph)} \\ \text{The diagonal entry } (\mathbf{A}^T \mathbf{A})_{ii} \text{ counts the edges meeting at node } i: \text{ the degree} \\ \text{The off-diagonal entry is } (\mathbf{A}^T \mathbf{A})_{ij} = -1 \text{ when an edge connects node } i \text{ and } j \\ \mathbf{G} \text{ and } \mathbf{K} \text{ are positive semidefinite but not positive definite (because } \mathbf{A}\mathbf{x} = \mathbf{0} \text{ in } \mathbf{1}) \end{cases}$$

# Clustering

- ## How to understand a graph with many nodes?
  - Separate nodes into two or more clusters
  - Human Genome project: cluster genes that show highly correlated

- ## Break a graph in two pieces: clusters of nodes
  - Each cluster should contain roughly half of the nodes
  - The number of edges between clusters should be relatively small

- ## Examples
  - For load balancing in high computing, assign equal work to two processors
  - For social networks, identify two distinct groups
  - Segment an image
  - Reorder rows and columns of a matrix to make off-diagonal blocks sparse

# Example with Two Clusters



$n = 5$ nodes, $k = 2$ clusters

centroid $*$ : $\mathbf{c}_1 = (-1, 2/3)$, $\mathbf{c}_2 = (2,1) \leftarrow$ minimize the sume of squared distances $\left\| \mathbf{c} - \mathbf{a}_j \right\|^2$

Approximate an $m \times n$ matrix of $\mathbf{A}$ by $\underset{\text{low rank}}{\mathbf{CR}} = \underset{\substack{\text{only } k \text{ columns} \\ \text{centroids of clusters}}}{\underbrace{(m \times k)}} \; \underset{\substack{\text{single 1 and} \\ k-1 \text{ zeros}}}{\underbrace{(k \times n)}}$

$R_{ij} = 1 (\text{or } 0)$ if centroid $i$ is closest (or not) to the point $\mathbf{x}_j$

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 3 & 0 & -3 \\ 0 & 1 & 1 & 2 & 0 \end{bmatrix} \approx \begin{bmatrix} -1 & 2 & 2 & -1 & -1 \\ 2/3 & 1 & 1 & 2/3 & 2/3 \end{bmatrix}$$

$$\mathbf{A} \approx \mathbf{CR} = \begin{bmatrix} -1 & 2 \\ 2/3 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

# Four Methods for Clustering

- Spectral clustering (Fiedler vector)
  - using the graph Laplacian or the modularity matrix
- Minimum cut
- Weighted k-means

# Four Methods for Clustering

I. $\mathbf{A}^T \mathbf{C} \mathbf{A} \mathbf{z} = \lambda \mathbf{D} \mathbf{z} \rightarrow \mathbf{z}$ : Fiedler vector $\begin{cases} \lambda_1 = 0 \rightarrow \mathbf{z}_1 = (1,\ldots,1) \\ \lambda_2 \rightarrow \mathbf{z}_2 : \ +/- \text{ components indicate two clusters of nodes} \end{cases}$

II. $\mathbf{A}^T \mathbf{C} \mathbf{A} \rightarrow \underset{\substack{\text{modularity} \\ \text{matrix}}}{\mathbf{M}} = \mathbf{B} - \dfrac{1}{2m} \mathbf{d} \mathbf{d}^T$ where $\mathbf{d}$ : degrees of the $n$ nodes (number of edges adjacent to the nodes)

choose eigenvector that comes with the largest eigenvalue of $\mathbf{M}$

III. Find the minimum normalized cut that separates the nodes in two clusters P and Q

weight across cut: $links(P) = \sum w_{ij}$ for $i$ in $P$ and $j$ not in $P$

size of cluster: $size(P) = \sum w_{ij}$ for $i$ in $P$

normalized cut weight: $N_{cut}(P,Q) = \dfrac{links(P)}{size(P)} + \dfrac{links(Q)}{size(Q)} \xrightarrow{k-cut} N_{cut}(P_1,\ldots,P_k) = \sum_{i=1}^{k} \dfrac{links(P_i)}{size(P_i)}$

minimize $N_{cut}(P,Q)$ : good partition of the graph $\rightarrow$ application: segmentation of images

IV. nodes in the graph: $\mathbf{a}_1,\ldots,\mathbf{a}_n$, clusters $P$ and $Q$ have centers $\mathbf{c}_P \left( = \dfrac{\sum \mathbf{a}_i}{|P|} \right)$ and $\mathbf{c}_Q$

minimize the total squred distance from nodes to those centroids: $E = \sum_{i \text{ in } P} \|\mathbf{a}_i - \mathbf{c}_P\|^2 + \sum_{i \text{ in } Q} \|\mathbf{a}_i - \mathbf{c}_Q\|^2$

# Spectral Clustering (1)

$$\mathbf{A}^T\mathbf{CA} \xrightarrow[\text{the Laplacian}]{\mathbf{D} \text{ normalizes}} \mathbf{L} = \mathbf{D}^{-1/2}\mathbf{A}^T\mathbf{CAD}^{-1/2} = \mathbf{I} - \mathbf{N} \quad \text{where } n_{ij} = \frac{w_{ij}}{\sqrt{d_i d_j}} \left(\text{normalized weights}\right)$$

triangular graph: $n = 3$ nodes, $m = 3$ edges, $c_1, c_2, c_3 = w_{12}, w_{13}, w_{23}$

$$\mathbf{A}^T\mathbf{CA} = \mathbf{D} - \mathbf{W}$$

$$\left.\begin{bmatrix} w_{12} + w_{13} & -w_{12} & -w_{13} \\ -w_{21} & w_{21} + w_{23} & -w_{23} \\ -w_{31} & -w_{32} & w_{31} + w_{32} \end{bmatrix}\right\} \rightarrow \left\{\begin{array}{l} \mathbf{L} = \mathbf{D}^{-1/2}\mathbf{A}^T\mathbf{CAD}^{-1/2} \\ \begin{bmatrix} 1 & -n_{12} & -n_{13} \\ -n_{21} & 1 & -n_{23} \\ -n_{31} & -n_{32} & 1 \end{bmatrix} \end{array}\right.$$

$\mathbf{L} = \mathbf{I} - \mathbf{N}$ is like a correlation matrix in statistics

1. $\mathbf{L}$ is symmetric positive semidefinite: orthogonal eigenvectors, all eigenvalues $\lambda \geq 0$

2. The eigenvectors for $\lambda = 0$ is $\mathbf{u} = \left(\sqrt{d_1}, \ldots, \sqrt{d_n}\right)$. Then $\mathbf{Lu} = \mathbf{D}^{-1/2}\mathbf{A}^T\mathbf{CA1} = 0$.

3. The second eigenvector $\mathbf{v}$ of $\mathbf{L}$ minimizes the Rayleigh quotient on a subspace.

$$\left( \begin{array}{l} \lambda_2 = \text{smallest nonzero eigenvalue of } \mathbf{L} \rightarrow \min_{\substack{\text{subject to} \\ \mathbf{x}^T\mathbf{u}=0}} \frac{\mathbf{x}^T\mathbf{Lx}}{\mathbf{x}^T\mathbf{x}} = \frac{\mathbf{v}^T\mathbf{Lv}}{\mathbf{v}^T\mathbf{v}} = \lambda_2 \text{ at } \mathbf{x} = \mathbf{v} \\ \\ \text{upper bound for } \lambda_2, \text{ for any } \mathbf{x} \text{ orthogonal to the first eigenvector } \mathbf{u} = \mathbf{D}^{-1/2}\mathbf{1} \end{array} \right)$$

# Spectral Clustering
## normalized vs. unnormalized

$$\mathbf{Lv} = \mathbf{D}^{-1/2}\mathbf{A}^T\mathbf{CAD}^{-1/2}\mathbf{v} = \lambda\mathbf{v} \xrightarrow[\text{normalized Fiedler vector}]{\mathbf{z}=\mathbf{D}^{-1/2}\mathbf{v}}$$

$$\begin{cases} \mathbf{A}^T\mathbf{CAz} = \lambda\mathbf{Dz} \;\; \text{with} \;\; \mathbf{1}^T\mathbf{Dz} = 0 \\ \text{generalized eigenvalue problem} \\ \text{eigenvector for } \lambda = 0 \text{ is } \mathbf{1} \\ \text{the next eigenvector } \mathbf{z} \text{ is } \mathbf{D}\text{-orthogonal to } \mathbf{1} \\ \mathbf{A}^T\mathbf{CAz} = \lambda_2\mathbf{Dz} \end{cases}$$

$$\min_{\substack{\text{subject to} \\ \mathbf{x}^T\mathbf{u}=0}} \frac{\mathbf{x}^T\mathbf{Lx}}{\mathbf{x}^T\mathbf{x}} \xrightarrow{\mathbf{x}=\mathbf{D}^{1/2}\mathbf{y}} \min_{\substack{\text{subject to} \\ \mathbf{1}^T\mathbf{Dy}=0}} \frac{\mathbf{y}^T\mathbf{A}^T\mathbf{CAy}}{\mathbf{y}^T\mathbf{Dy}} = \frac{\sum\sum w_{ij}\left(y_i - y_j\right)^2}{\sum d_i y_i^2} = \lambda_2 \text{ at } \mathbf{y} = \mathbf{z}$$
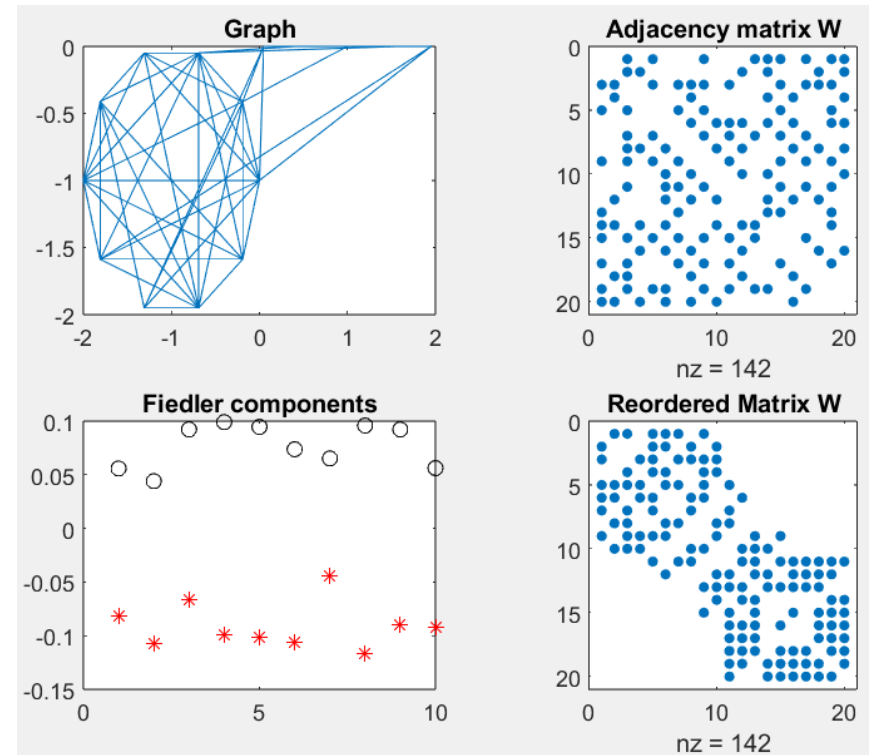
$\mathbf{Ay}$ : incidence matrix $\mathbf{A}$ gives the differences $\left(y_i - y_j\right)$

- Example: 20-node graph has two 10-node clusters P and Q (to find from z)
  - Create edges within P and Q with probability 0.7. Edges between nodes in P and Q have smaller probability 0.1. All edges have wrights $w_{ij}$=1. (C=I)

# Code: MATLAB

```matlab
N=10; W=zeros(2*N,2*N); % Generate 2N nodes in two clusters
rand('state',100) % rand repeats to give the same graph
for i=1:2*N-1
for j=i+1:2*N
p=0.7-0.6*mod(j-i,2); % p=0.1 when j-i is odd, 0.7 else
W(i,j)=rand<p; % Insert edges with probability p
end % The weights are wi,j=1 (or 0)
end % So far W is strictly upper triangular
W=W+W'; D=diag(sum(W)); % Adjacency matrix W, degress in D
G=D-W; [V,E]=eig(G,D); % Eigenvalues of Gx=(lambda)Dx in E
[a,b]=sort(diag(E)); z=V(:,b(2));% Fiedler eigenvector z for (lambda)2
plot(sort(z),'.-'); % Show +- groups of Fiedler components


theta=[1:N]*2*pi/N; x=zeros(2*N,1); y=x; % Angles to plot graph
x(1:2:2*N-1)=cos(theta)-1; x(2:2:2*N)=cos(theta)+1;
y(1:2:2*N-1)=sin(theta)-1; x(2:2:2*N)=sin(theta)+1;
print theta,x,y
subplot(2,2,1), gplot(W,[x,y]), title('Graph')
subplot(2,2,2), spy(W), title('Adjacency matrix W')
subplot(2,2,3), plot(z(1:2:2*N-1),'ko'), hold on
plot(z(2:2:2*N),'r*'), hold off, title('Fiedler components')
[c,d]=sort(z); subplot(2,2,4), spy(W(d,d)), title('Reordered Matrix W')
```

# Minimum Cut

$(\text{edge})$ weight across cut: $links(P) = \sum w_{ij}$ for $i$ in $P$ and $j$ not in $P$

size of cluster: $size(P) = \sum w_{ij}$ for $i$ in $P$

normalized cut weight: $Ncut(P,Q) = \dfrac{links(P)}{size(P)} + \dfrac{links(Q)}{size(Q)}$

normalized $K$-cut: $Ncut(P_1,\ldots,P_k) = \displaystyle\sum_{i=1}^{K} \dfrac{links(P_i)}{size(P_i)}$

$[\text{cuts connected to eigenvectors}]$

perfect indicator of a cut: vector $\mathbf{y}$ with all components equal to $p$ or $-q$ (two values only) $\rightarrow$ node $i$ goes $\begin{cases} \text{in } P \text{ if } y_i = p \\ \text{in } Q \text{ if } y_i = -q \end{cases}$

$\mathbf{1}^T \mathbf{Dy}$ will multiply one group of $d_i$ by $p$ and the other group by $-q$.

The first $d_i$ add to $size(P) = $ sum of $d_i$ $(i$ in $P)$.

The second group of $d_i$ add to $size(Q)$

$\left. \right\} \rightarrow \mathbf{1}^T \mathbf{Dy} = 0 \rightarrow psize(P) = qsize(Q)$

$$\dfrac{\mathbf{y}^T \mathbf{A}^T \mathbf{CAy}}{\mathbf{y}^T \mathbf{Dy}} = \dfrac{\sum\sum w_{ij}(y_i - y_j)^2}{\sum d_i y_i^2} = \dfrac{(p+q)^2 \, links(P,Q)}{p^2 size(P) + q^2 size(Q)} = \dfrac{(p+q)\, links(P,Q)}{psize(P)} = \dfrac{links(P,Q)}{size(P)} + \dfrac{links(P,Q)}{size(Q)} = Ncut(P,Q)$$

# Clustering by k-means

$n$ points $\mathbf{a}_1,\ldots,\mathbf{a}_n$ in d-dimensional space $\rightarrow$ partition those points into $k$ clusters

clusters $P_1,\ldots,P_k$ have centroids $\mathbf{c}_1,\ldots,\mathbf{c}_k$

$$\mathbf{c}_j = \frac{\text{sum of } \mathbf{a}\text{'s}}{\text{number of } \mathbf{a}\text{'s}} \quad \rightarrow \quad \text{minimize} \sum \|\mathbf{c} - \mathbf{a}\|^2 \text{ for all } \mathbf{a}\text{'s in cluster } P_j$$

clustering: to find the partition $P_1,\ldots,P_k$ with minimum total distance $D$ to centroids:

$$\text{minimize } D = \sum_{j=1}^{k} D_j = \sum_{j=1}^{k} \|\mathbf{c}_j - \mathbf{a}_i\|^2 \text{ for } \mathbf{a}_i \text{ in cluster } P_j$$

step 1: find the centroids $\mathbf{c}_j$ of the (old) clustering $P_1,\ldots,P_k$.

step 2: find the (new) clustering that puts $\mathbf{a}$ in $P_j$ if $\mathbf{c}_j$ is the closest centroid.

# Clustering by k-means: Weights and Kernel Method

weighted distance: $d(\mathbf{x}, \mathbf{a}_i) = w_i \|\mathbf{x} - \mathbf{a}_i\|^2$, $\mathbf{c}_j = \dfrac{\sum w_i \mathbf{a}_i}{\sum w_i} \left(\mathbf{a}_i \text{ in } P_j\right)$

distances to centroids only require dot product $\mathbf{a}_i \cdot \mathbf{a}_j$: $\left(\text{each } i \text{ in } P_j\right)$ $\|\mathbf{c}_j - \mathbf{a}_i\|^2 = \mathbf{c}_j \cdot \mathbf{c}_j - 2\mathbf{c}_j \cdot \mathbf{a}_i + \mathbf{a}_i \cdot \mathbf{a}_i$

Kernel method: weighted kernel matrix $\mathbf{K}$ has entries $\mathbf{a}_i \cdot \mathbf{a}_l$
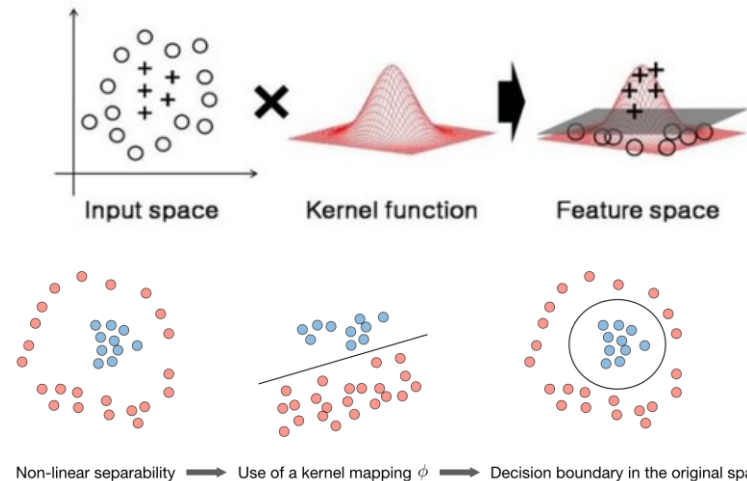
nodes are point $\mathbf{x}_i$ in input space $\rightarrow \mathbf{a}_i = \phi(\mathbf{x}_i)$ points in a high-dimensional <span style="color:red">feature space</span>

$\left(\text{sum over nodes in } P_j\right)$ $\sum \|\mathbf{c}_j - \mathbf{a}_i\|^2 = \dfrac{\sum w_i w_l \mathbf{K}_{il}}{\left(\sum w_i\right)^2} - 2\dfrac{\sum w_i \mathbf{K}_{il}}{\sum w_i} + \sum \mathbf{K}_{ii}$

$\left(\text{vision}\right)$ polynomial $\mathbf{K}_{il} = \left(\mathbf{x}_i \cdot \mathbf{x}_l + c\right)^d$

$\left(\text{statistics}\right)$ Gaussian $\mathbf{K}_{il} = \exp\left(-\dfrac{\|\mathbf{x}_i - \mathbf{x}_l\|^2}{2\sigma^2}\right)$

$\left(\text{neural networks}\right)$ Sigmoid $\mathbf{K}_{il} = \tanh\left(c\mathbf{x}_i \cdot \mathbf{x}_l + \theta\right)$



Input space     Kernel function     Feature space

Non-linear separability ➡ Use of a kernel mapping $\phi$ ➡ Decision boundary in the original space

for large data sets, $k$-means and $\text{eig}\left(\mathbf{A}^T \mathbf{C} \mathbf{A}, \mathbf{D}\right)$ will be expensive $\rightarrow \begin{cases} \text{random sampling} \\ \text{multilevel clustering} \end{cases}$

# Applications of Clustering

- Learning theory, training sets, neural networks, Hidden Markov Models
- Classification, regression, pattern recognition, Support Vector Machines
- Statistical learning, maximum likelihood, Bayesian statistics, spatial statistics, kriging, time series, ARMA models, stationary processes
- Social networks, organization theory
- Data mining, document indexing, image retrieval, kernel-based learning, Nystrom method, low rank approximation
- Bioinformatics, microarray data, systems biology
- Cheminformatics, drug design, decision trees
- Information theory, vector quantization, rate distortion theory, Bregman divergences
- Image segmentation, computer vision, texture, min cut
- Predictive control, feedback samples, robotics, adaptive control, Riccati equations