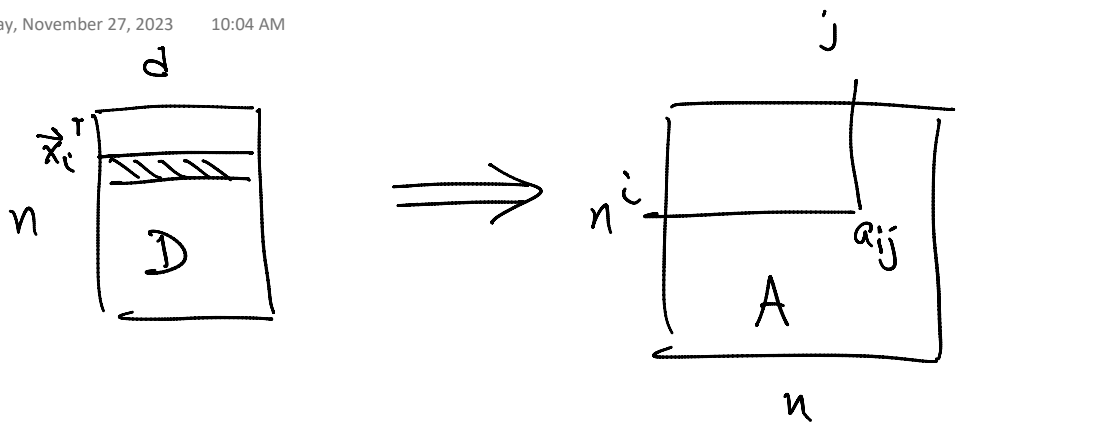


Lecture 23

Monday, November 27, 2023 10:04 AM



$n \times n$
symmetric
 $a_{ij} = a_{ji}$
 $a_{ij} \geq 0$

$A = \underbrace{K}_{\text{kernel matrix}}$

using Gaussian kernel

$$a_{ij} = e^{-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{2\sigma^2}}$$

↑ user chosen variance value

Markov matrix

$$M = \Delta^{-1} A$$



$$M = \begin{bmatrix} 1/d_1 & 0 & \dots \\ 0 & 1/d_2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$\Delta =$ degree matrix

$$\Delta = \begin{bmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & d_n \end{bmatrix}$$

$$d_i = \sum_{j=1}^n a_{ij}$$

sum of each row

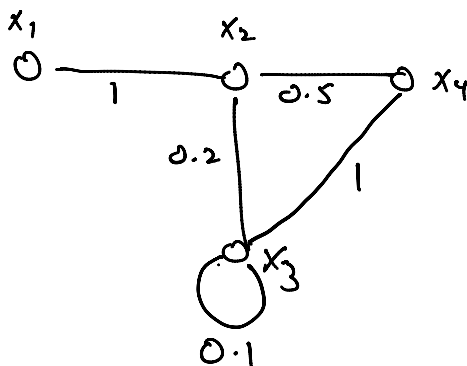
$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$$M = \Delta^{-1} \begin{bmatrix} \text{---} y_{d2} \text{---} \\ \text{---} 1/d_n \text{---} \end{bmatrix} \begin{bmatrix} \overline{a_{11}} \overline{a_{12}} \dots \overline{a_{1n}} \\ \overline{a_{21}} \overline{a_{22}} \dots \overline{a_{2n}} \\ \vdots \end{bmatrix} \quad \text{Sum of each row}$$

$$M = \begin{bmatrix} a_{11}/d_1 & a_{12}/d_1 & \dots & a_{1n}/d_1 \\ a_{21}/d_2 & a_{22}/d_2 & \dots & a_{2n}/d_2 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$$= \{ m_{ij} \}_{i,j=1 \dots n}$$

m_{ij} is the probability of 'jumping' from i to j



Each row is a probability vector

$$m_{ij} \geq 0$$

$$\sum_j m_{ij} = 1 \leftarrow \text{for each row}$$

	x_1	x_2	x_3	x_4
x_1	0	1	0	0
x_2	1	0	0.2	0.5
x_3	0	0.2	0.1	1
x_4	0	0.5	1	0

A

$$\Delta = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1.7 & 0 & 0 \\ 0 & 0 & 1.3 & 0 \\ 0 & 0 & 0 & 1.5 \end{bmatrix}$$

degree

Sim

$$M = \Delta^{-1} A =$$

0	1	0	0
0.59	0	0.12	0.29
0	0.12	0.08	0.75
0	0.33	0.67	0

Inflate $r=2$

0	1	0	0
0	1	0	0
0	1	0	0
0	1	0	0

$$M = \Delta A =$$

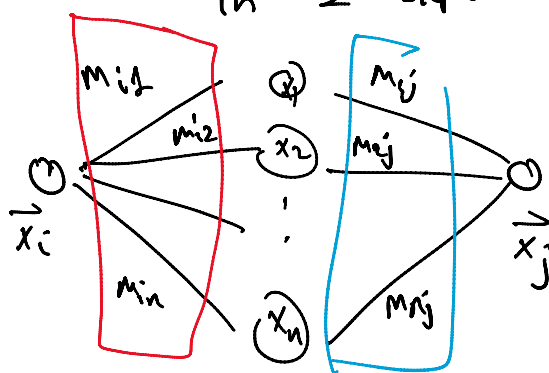
0.59	0	0.12	0.29
0	0.15	0.08	0.77
0	<u>0.33</u>	<u>0.67</u>	0

0	1/2	0	0
0	1/3	2/3	0
0	1/9	4/9	0

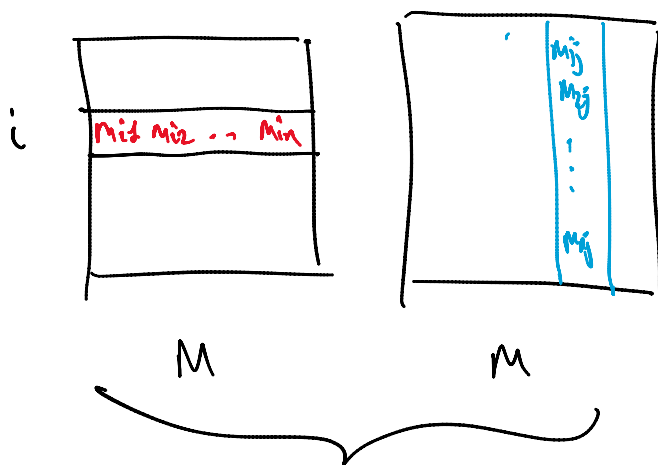
0.2 0.8
1/9 8/9
1/9 8/9

$M \equiv$ 1-step transition matrix

Q: What is the prob of ending up on \vec{x}_j starting from \vec{x}_i in 2 steps



M independent of step or time (homogeneous)



$\equiv M_{ij}^{(2)}$ prob from i to j in 2 steps

$M^2 \equiv$ 2-step transition matrix

$M^t \equiv$ t-step transition matrix

π is a initial prob vector

= initial prob of starting on node i (page i)

Col vector



what's the prob of ending up on j (buy j)

$$\vec{\pi} = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_n \end{bmatrix} = \begin{bmatrix} 0.01 \\ 0.1 \\ 0 \\ 0 \\ 0.5 \\ \vdots \end{bmatrix}$$

$$\Sigma = 1$$

$$\vec{x}_0 = \vec{\pi}$$

$$M \vec{x}_0 = \vec{x}_1$$

$$M \vec{x}_1 = \vec{x}_2$$

$$M(M \vec{x}_0) = \vec{x}_2$$

$$M^2 \vec{x}_0 = \vec{x}_2$$

\vdots

Converge to dominant Eigen vector
of M

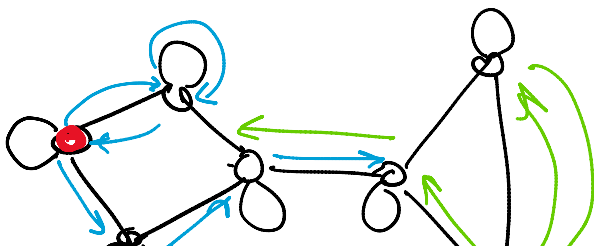
$$\pi, x_1, x_0 \in \mathbb{R}^n$$

final "destination"
probabilities

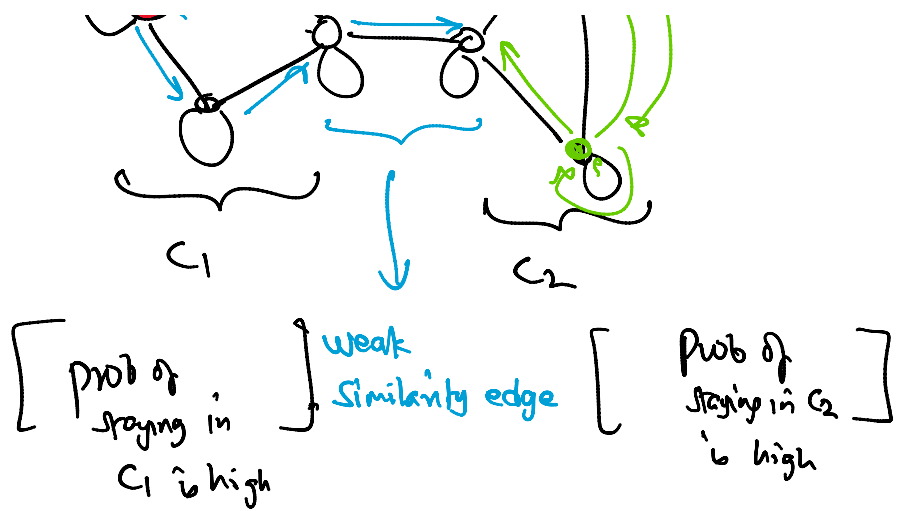
are independent of
the initial prob vector

Q: how to use M matrix for clustering

random walk



All nodes will
have self-loops
(for clustering)



(for clustering)

MCL: Markov chain clustering

starting with $M \leftarrow 1$ -step transition

a) walk one more step

$$\underline{M} = M \cdot M = M^2$$

b) Inflate M (prob that are high should be made even higher)
& prob that are low " " " even lower

$\gamma \equiv$ inflation parameter

- i) Compute γ -th power of each element of M
- ii) normalize

$$m_{ij} = \frac{m_{ij}^{(\gamma)}}{\sum_{j=1}^n m_{ij}^{(\gamma)}}$$

(each row should remain prob vector)

c) repeat a) & b) until convergence

c) repeat a) & b) until Convergence

$$\| M^{(t)} - M^{(t-1)} \|_F \leq \epsilon \quad \text{Stop}$$

(t) step number

sqr (sum of differences)

d) extract clusters

Using the final M we create a directed graph G
Sparse

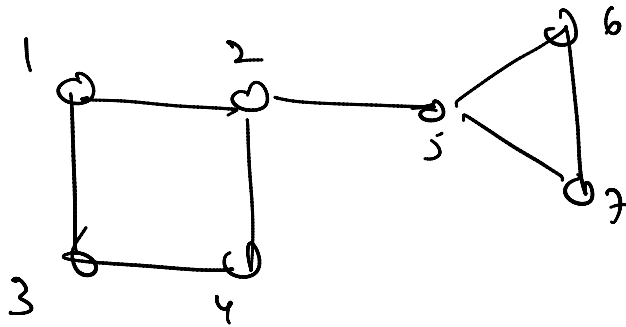
$$E = \{ (i, j) \mid \underline{M(i, j)} > 0 \}$$

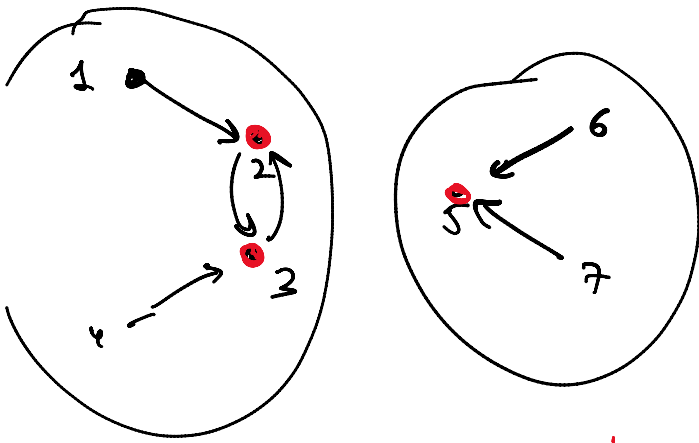
find weakly connected components in G

If $M(i, j) > 0 \Rightarrow i$ can transition to j

$M_{ij} > 0$

we also say that j is an
attractor for i





Weakly Connected Components

γ parameter

γ is large

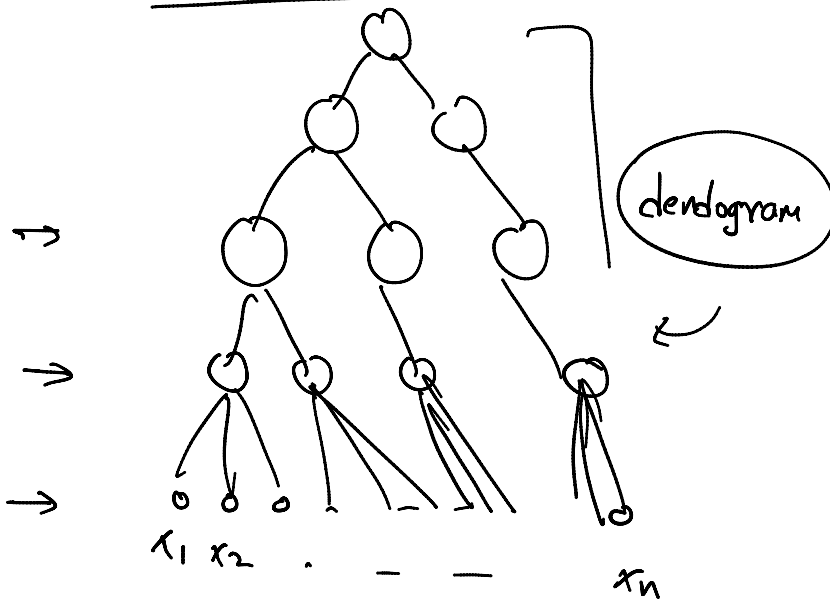
↳ many clusters, small

γ is small

↳ fewer, larger clusters

$\gamma \geq 1$

Hierarchical clustering



K-means \leftarrow hard

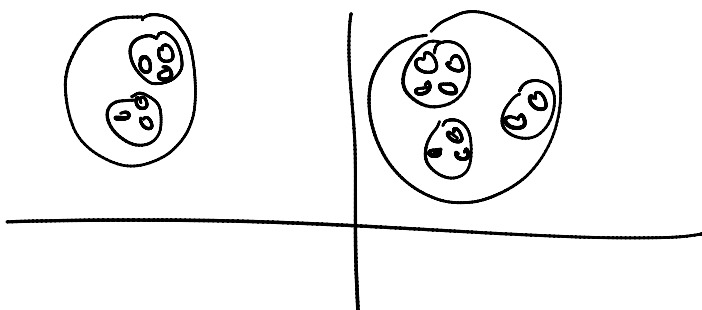
EM \leftarrow soft

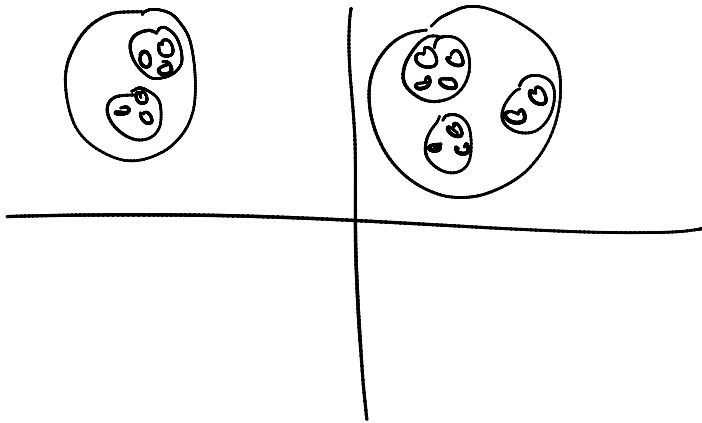
density \leftarrow non-convex overlaps

Spectral \rightarrow hard

ML \rightarrow overlaps

glar





Agglomerative hierarchical clustering ($K \leftarrow$ find K clusters)

$$1) \mathcal{C} = \{x_i \text{ in its own cluster}\}$$

$$= \{C_1, C_2, \dots, C_n\}$$

$$C_i = \{\vec{x}_i\}$$

repeat until $|\mathcal{C}| = K$ (or $|\mathcal{C}| = 1$)

full tree

find the closest pair of clusters
 C_i & C_j and merge

$$C_{ij} = C_i \cup C_j$$

$$\Omega(n^2)$$

to \mathcal{C} add C_{ij} and remove C_i and C_j

$$O(n^2 \log n)$$

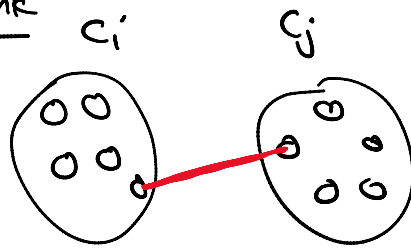
update distances from C_{ij} to all other C_k

$$s(C_{ij}, C_k) \quad \forall k$$

$$s(C_i, C_j) \quad \forall i, j$$

we need to define what is the distance between two clusters

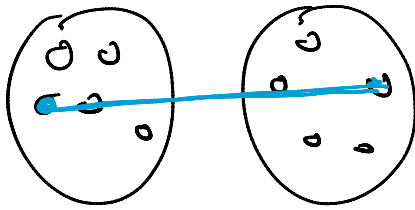
Single link



$$S(C_i, C_j) = \min \left\{ \|\vec{x} - \vec{y}\| \right\}$$

$\vec{x} \in C_i$
 $\vec{y} \in C_j$

Complete link



$$S(C_i, C_j) = \max \left\{ \|\vec{x} - \vec{y}\| \right\}$$

$\vec{x} \in C_i$
 $\vec{y} \in C_j$

Intuitively the farthest pair
includes all other pairs

Group average

$$S(C_i, C_j) = \text{avg} \left\{ \|\vec{x} - \vec{y}\| \right\} = \frac{\sum_{\vec{x} \in C_i} \sum_{\vec{y} \in C_j} \|\vec{x} - \vec{y}\|}{n_i \cdot n_j}$$

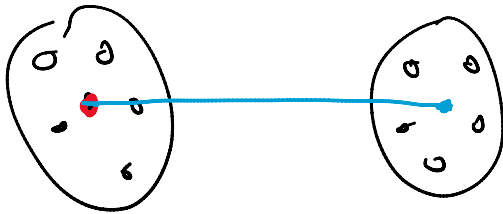
$\forall \vec{x} \in C_i$
 $\vec{y} \in C_j$

$$n_i = |C_i|$$

mean distance

$$S(C_i, C_j) = \|\vec{\mu}_i - \vec{\mu}_j\|$$

$$\vec{\mu}_i = \sum_{\vec{x} \in C_i} \vec{x} / n_i$$



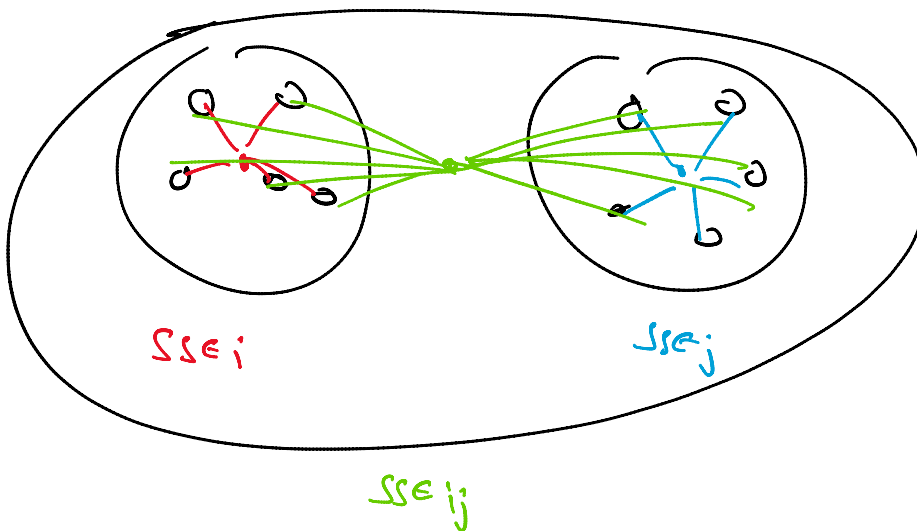
Minimum Variance (Wards Method)

$$\underline{S(C_i, C_j)} = \underbrace{\Delta SSE_{ij}}$$

Change in squared error

$$= SSE_{ij} - SSE_i - SSE_j = \left(\frac{n_i n_j}{n_i + n_j} \right) \|\vec{\mu}_i - \vec{\mu}_j\|^2$$

$$SSE_i = \sum_{\vec{x}_j \in C_i} \|\vec{x}_j - \vec{\mu}_i\|^2$$



Evaluate Cluster?

→ Unsupervised

external

when we
have some
ground truth
labels

$$T = \{T_1, T_2, \dots, T_K\}$$

ground truth clustering

$$\mathcal{C} = \{c_1, c_2, \dots, c_r\}$$

$$r = k$$

metrics: T vs \mathcal{C}

internal

no labels!

pairwise distances \equiv

proximity matrix $= P$

either in input space
or in feature space

relative (no labels)

\mathcal{C}_1 vs \mathcal{C}_2