

Lecture 16: Diffusion maps

Last time we talked about spectral embeddings,

$$F: X \rightarrow \mathbb{R}^J \quad F(\underline{x}_i) = ((\underline{v}_1)_i, (\underline{v}_2)_i, \dots, (\underline{v}_J)_i)^T$$

which formed the basis of the spectral clustering algorithm. At a high level, spectral clustering is nothing more than a two step procedure:

Step 1) Pre-process/embed data w/ F

Step 2) Apply standard clustering such as k-NN.

It is natural to ask whether our choice of F can be made more flexible or more mathematically meaningful. For starters, observe that the kernel

$$K^J(\underline{x}, \underline{x}') = F(\underline{x})^T F(\underline{x}'), \quad F: X \rightarrow \mathbb{R}^J$$

is inconsistent, in the sense that $\lim_{J \rightarrow \infty} K^J$ does not exist.

Throughout this section we will consider the random walk graph Laplacian

①

$$\underset{\text{RW}}{L} = I - A, \quad A = D^{-1}W$$

The material presented here is based on the seminal paper Coifman & Lafon, "Diffusion Maps" (2006).

Let $G = \{X, W\}$ be a weighted graph & for notational purposes we let $V = \{1, \dots, n\}$ be the set of vertices of G .

It is helpful to think of functions on the vertices of G : $u: V \rightarrow \mathbb{R}$ $u \equiv \underline{u} \in \mathbb{R}^n$

If \underline{u} is such that $u(i) \geq 0$ & $\sum_{i=1}^n u(i) = 1$ then

\underline{u} represents a probability distribution on V . & we write $\underline{u} \in P(V)$.

Observe that (assuming no lonely nodes)

$$A \underline{1} = \bar{D}^{-1} W \underline{1} = \bar{D}^{-1} \underline{d} = \underline{1}$$

Thus, A is right stochastic & hence a transition matrix for a discrete Markov chain on V .

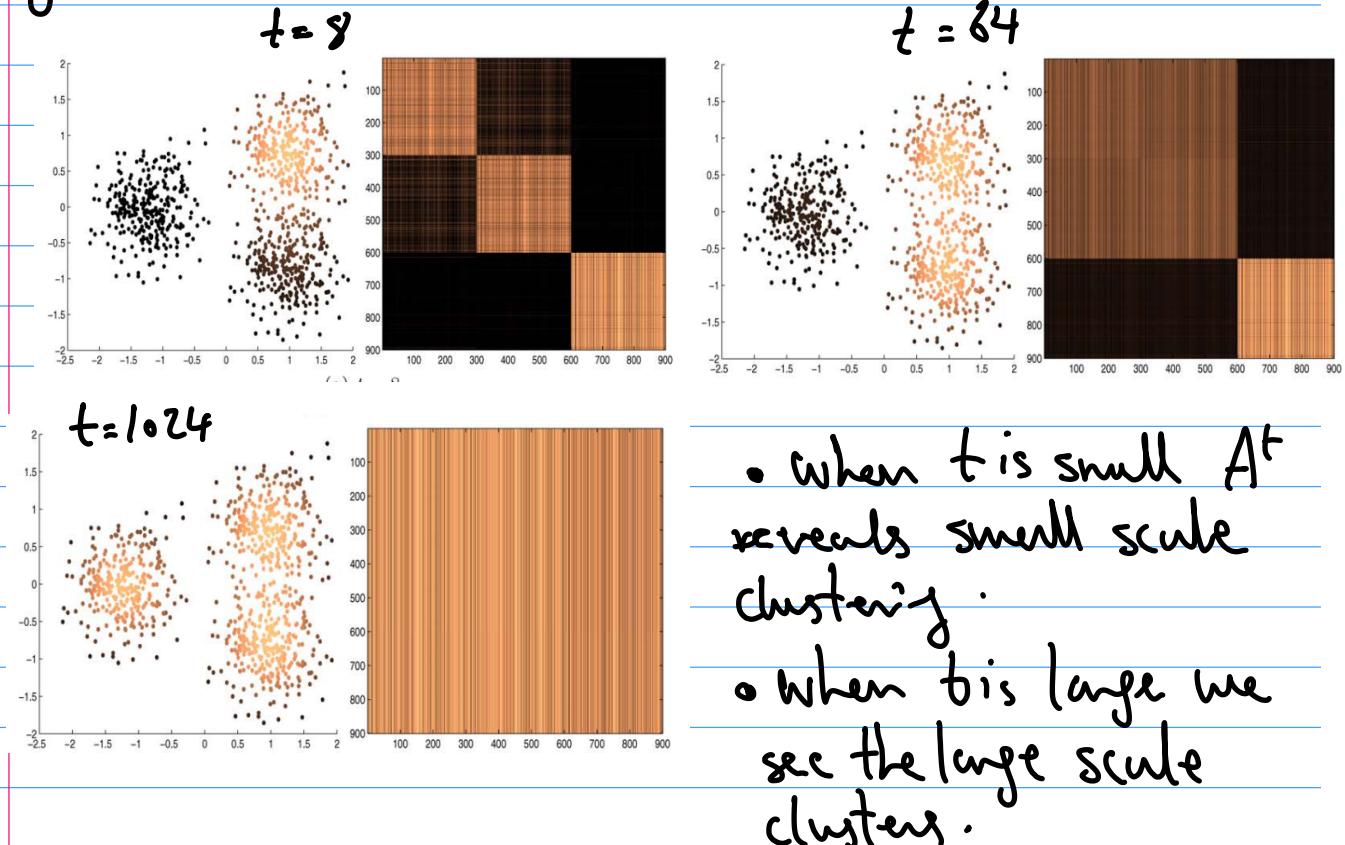
Indeed the entry A_{ij} denotes the probability of going from vertex i to vertex j .

Thus, if $\underline{u} \in P(V)$ denotes the initial state of the chain, the distribution at future time steps is given by

(2)

$$\underline{u}^t = \underline{u}^{t-1} A = \underline{u}^{t-2} A^2 = \underline{u}^0 A^t$$

The idea of Diffusion maps is to exploit the spectrum of A^t to reveal multi-scale geometric information in the data set.



- When t is small A^t reveals small scale clustering.
- When t is large we see the large scale clusters.

The idea of diffusion maps:

Let $P_t(i, \cdot)$ denote the distribution of the chain initiated at vertex $i \in V$ at time $t \geq 0$. Then the diffusion dist. between $\underline{x}_i, \underline{x}_j \in X$ is given by

$$D_t(\underline{x}_i, \underline{x}_j) := \left(\sum_{l=1}^n \frac{1}{d_l} |P_t(i, l) - P_t(j, l)|^2 \right)^{\frac{1}{2}}$$

$$= \|A_{i:}^t - A_{j:}^t\|_{\bar{D}}^2 \leftarrow \begin{matrix} i, j \text{ th} \\ \text{rows of } A^t \end{matrix}$$

Indeed, defining the maps

$$F: X \rightarrow \mathbb{R}^n \quad F(\underline{x}_i) = \begin{bmatrix} P_t(i, 1) \\ P_t(i, 2) \\ \vdots \\ P_t(i, n) \end{bmatrix}$$

we see that

$$D_t^2(\underline{x}_i, \underline{x}_j) = \|F(\underline{x}_i) - F(\underline{x}_j)\|_{\bar{D}}^2$$

The above definition is very intuitive but it is unclear whether it can be computed efficiently. Turns out a better formulation, related to the spectral embedding, is possible which we now derive.

Let $A = \bar{D}^{-\frac{1}{2}} W$ & define

$$B = \bar{D}^{\frac{1}{2}} A \bar{D}^{-\frac{1}{2}} = \bar{D}^{-\frac{1}{2}} W \bar{D}^{-\frac{1}{2}}$$

(4)

observe B is symm. in fact $I - B = L_{\text{normalized}}$

$$(\text{recall } L_{\text{Norm}} = D^{-1/2} (D - w) D^{1/2})$$

write $B = V \Lambda V^T$ (eigen decom)

& intum

$$A = D^{-1/2} B D^{1/2}$$

$$= D^{-1/2} (V \Lambda V^T) D^{1/2}$$

(since V is
orth)

$$= (D^{-1/2} V) \Lambda (V^T D^{1/2})$$

$$= Q \Lambda Q^{-1}$$

In otherwords, A has an eigen decom.

It has the same eig. values as B & columns
of $Q = D^{-1/2} V$ are the right eig. vectors of A

& rows of $Q^{-1} = V^T D^{1/2}$ are its left eigen
vectors. In other words,

$$A = \sum_{j=1}^n \lambda_j q_j r_j^T$$

$$\begin{aligned} q_j &= D^{-1/2} v_j \\ r_j &= D^{1/2} u_j \end{aligned}$$

Finally, we have also that

(5)

$$A^t = \sum_{j=1}^n \lambda_j^t q_j r_j^T$$

Now observe that

$$\langle \underline{r}_i, \underline{r}_j \rangle_{D^{-1}} = \underline{r}_i^T D^{-1} \underline{r}_j = \underline{v}_i^T D^{1/2} D^{-1} D^{1/2} \underline{v}_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

so that the \underline{r}_i form an orthonormal basis of $(\mathbb{R}^n, \langle \cdot, \cdot \rangle_{D^{-1}})$. Using this observation

we can see that (see Sec 8, la Porte et al

2008)

$$D_t(x_i, x_j) = \| A_{i:}^t - A_{j:}^t \|_{D^{-1}}^2$$

$$= \sum_{l=1}^n \lambda_l^{2t} |(\underline{q}_l)_i - (\underline{q}_l)_j|^2$$

$$= \| \mathcal{F}(x_i) - \mathcal{F}(x_j) \|_2^2$$

where

$$\mathcal{F}: X \rightarrow \mathbb{R}^n, \quad \mathcal{F}(x_i) = \begin{bmatrix} \lambda_1^t q_1(i) \\ \lambda_2^t q_2(i) \\ \vdots \\ \lambda_n^t q_n(i) \end{bmatrix}$$

(diffusion map)

(6)

Since $q_1 \propto \mathbb{1}$ & $\lambda_1 = 1$, some practitioners often choose to leave this eigen pair out of their algorithms.

Observe: The map \mathcal{F} is very similar to the spectral embedding

$$F(x_i) = \begin{bmatrix} (v_1)_i \\ (v_2)_i \\ \vdots \\ (v_j)_i \end{bmatrix} \quad \leftarrow v_i: \text{eig. vec. of D-W}$$

The construction of the diff. maps suggests families of embeddings of the form

$$\mathcal{F}_t(x_i) = \begin{bmatrix} \lambda_1^t q_1(i) \\ \lambda_2^t q_2(i) \\ \vdots \\ \lambda_n^t q_n(i) \end{bmatrix} \quad t \geq 0$$

where (λ_j, q_j) are eigen pairs of a Laplacian operator (normalized, unnormalized, etc.)

The map $\mathcal{F}(\underline{x}) \rightarrow \mathbb{R}^n$ defines a kernel $K(\underline{x}, \underline{x}') = \mathcal{F}(\underline{x})^T \mathcal{F}(\underline{x}')$ which we call the diffusion kernel.

Under appropriate assumptions one can show that this kernel coincides with the heat kernel on the manifold on which the data X lies.

(see Coifman & Lafon 2008 for discussion)

