

# Diffusion Maps and Spectral Clustering

Uri Shaham

October 11, 2022

## 1 Spectral Clustering

### 1.1 Graph Laplacians

Spectral clustering is a clustering algorithm which is far more expressive than  $k$ -means, and is not limited to convex clusters.

**Definition 1.1** (Weighted graph). *A weighted graph is an undirected graph  $G = (V, E)$  where  $V = \{v_1, \dots, v_n\}$  and each two vertices  $v_i, v_j$  are connected with an edge with non-negative weight  $w_{ij} \geq 0$ . We denote the weight matrix as  $W$ .*

**Remark 1.2.** *When data lies in Euclidean space and the weights are not given, it is a common practice to use Gaussian kernel to compute the weights*

$$w_{ij} = \frac{\exp(-\|x_i - x_j\|^2)}{2\sigma^2},$$

for such bandwidth parameter  $\sigma$ .

**Remark 1.3.** *A common practice is to connect each point only to its  $k$  nearest neighbors, which is known as knn graph. The graph can be easily made symmetric, for example by using  $W := \frac{1}{2}(W + W^T)$ .*

**Definition 1.4** (Degree and Degree matrix). • *The degree of a vertex  $v_i$  is defined as the sum of weights on its edges  $d_i = \sum_j w_{ij}$ .*

• *The degree matrix  $D$  is a diagonal matrix with elements  $D_{ii} = d_i$ .*

**Definition 1.5** (Unnormalized graph Laplacian). *The unnormalized graph Laplacian is defined as  $L_{un} = D - W$ .*

Observe that  $L_{un}$  is symmetric.

**Proposition 1.6.** *For every vector  $f \in \mathbb{R}^n$ ,  $f^T L_{un} f = \sum_{i,j} w_{ij} (f_i - f_j)^2$ .*

*Proof.* Exercise. □

It follows that  $L_{un}$  is positive semi-definite

**Proposition 1.7.** *The smallest eigenvalue of  $L_{un}$  is 0, and its corresponding eigenvector is the constant vector  $\frac{1}{\sqrt{n}} \mathbb{1}$ .*

*Proof.* Let  $f = \frac{1}{\sqrt{n}}\mathbf{1}$ . Then

$$f^T L_{\text{un}} f = \frac{1}{n} D - W = 0,$$

hence  $f, 0$  are an eigenpair.  $\square$

**Proposition 1.8.** *The multiplicity of the zero eigenvalue equals the number  $k$  of connected components of  $G$ , and the corresponding eigenvectors are indicator vectors of the components.*

*Proof.* We know that 0 is an eigenvalue. Let  $f$  be a corresponding eigenvector. Then since  $f^T L_{\text{un}} f = 0$ , we have that  $f_i = f_j$  whenever  $w_{ij} > 0$ . For  $k > 1$  connected components,  $L_{\text{un}}$  is block diagonal, so its spectrum is the union of the spectra of all blocks.  $\square$

**Definition 1.9** (Normalized graph Laplacian). *We define two versions of normalized Laplacians:*

- *The random walk graph Laplacian is  $L_{\text{rw}} = D^{-1} L_{\text{un}} = I - D^{-1} W$ .*
- *The Symmetric graph Laplacian is  $L_{\text{sym}} = D^{-\frac{1}{2}} L_{\text{un}} D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$ .*

**Remark 1.10.** *Observe that  $L_{\text{rw}}$  and  $L_{\text{sym}}$  are similar matrices, and hence share the same spectrum. In particular, all eigenvalues of  $L_{\text{rw}}$  are real.*

**Proposition 1.11.**  *$L_{\text{sym}}$  is positive semi-definite (and hence also  $L_{\text{rw}}$ )*

*Proof.* Exercise.  $\square$

**Proposition 1.12.** *0 is an eigenvalue of  $L_{\text{rw}}$  (and hence also of  $L_{\text{sym}}$ ), associated with the constant eigenvector.*

*Proof.* Exercise.  $\square$

**Proposition 1.13.** *The multiplicity of the zero eigenvalue in  $L_{\text{rw}}$  and  $L_{\text{sym}}$  equals the number  $k$  of connected components of  $G$ , and the corresponding eigenvectors are indicator vectors of the components.*

*Proof.* Analogous to the proof of proposition 1.8.  $\square$

## 1.2 Spectral Clustering

We have seen that eigenvectors corresponding to the zero eigenvalues (to whom we refer from now the first eigenvectors) of the Laplacian indicate the connected components. It therefore makes sense to use these eigenvectors to represent the data and identify the cluster structure.

Spectral clustering works by representing the data using the first  $k$  eigenvectors, and running  $k$ -means on this representation. Unlike  $k$ -means it can handle non-convex data (see Figure 1. Spectral clustering can interestingly be motivated by a graph cut point of view (see section 5 in Von Luxburg's excellent tutorial).

**Remark 1.14.** *Since the first eigenvector is constant in  $L_{\text{un}}$  and  $L_{\text{rw}}$ , it is often omitted.*

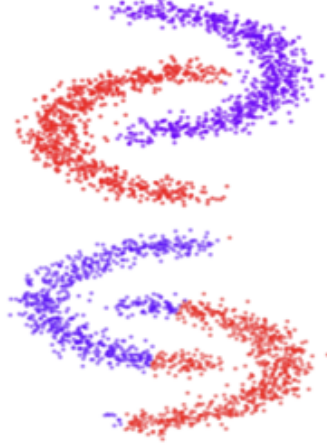


Figure 1: Spectral clustering (top) versus  $k$ -means on non-convex clusters.

## 2 Diffusion Maps

Diffusion map is a dimensionality reduction technique, that captures geometrical properties of data. In order to do so, it utilizes weighted graphs, which encode local similarity between pairs of points. This local interaction then allows to obtain global representations of the entire data.

Observe that  $P := D^{-1}W$  is a Markov matrix, and hence encodes transition probabilities between the vertices.  $P$  is a diffusion (averaging) operator, and defines directions of propagation. powers of  $P$  corresponds to multi-step walks.

**Proposition 2.1.**  $P$  has stationary distribution  $\pi$ , with  $\pi_i = \frac{d_i}{\sum_j d_j}$ , i.e.,  $\pi P = \pi$ .

*Proof.* Exercise. □

Let  $1 = \lambda_0 \geq \lambda_1 \geq \dots$  and  $\psi_0, \psi_1, \dots$  be the eigenvalues and eigenvectors of  $P$  (exercise: why  $\lambda_0 = 1$ ?)

**Definition 2.2** (Diffusion distance). *The diffusion distance between vertices  $v_i$  and  $v_j$  is defined as:*

$$D_t(v_i, v_j) = \|P_i^t - P_j^t\|_{\ell^2/d}^2 = \sum_{k=1}^n (P_{ik} - P_{jk})^2 / d_k,$$

where  $P_i^t$  is the  $i$ 'th row of  $P^t$ .

**Remark 2.3.** The notation  $\|x\|_{\ell^2/d}^2$  means that the length is  $\sqrt{\sum_i x_i^2 / d_i}$ .

Intuitively, if  $v_i$  and  $v_j$  are connected by a large number of paths, their diffusion probabilities will be similar, and hence their diffusion distance will be small.

The following theorem shows that the diffusion distance squared is in fact the Euclidean distance in the eigenspace of  $P$

**Theorem 2.4.**

$$D_t(v_i, v_j) = \left( \sum_{l \geq 1}^n \lambda_l^{2t} (\psi_{l,i} - \psi_{l,j})^2 \right)^{\frac{1}{2}}.$$

*Proof.* Consider  $A = D^{\frac{1}{2}} D^{-\frac{1}{2}} = D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$ . We already know that  $P$  has real eigenvalues (since  $P$  is similar to  $A$ , and  $A$  is symmetric). Let  $A = \sum \lambda_l \phi_l \phi_l^T$ . Then

$$\begin{aligned} P &= D^{-\frac{1}{2}} A D^{\frac{1}{2}} \\ &= \sum_l \lambda_l D^{-\frac{1}{2}} \phi_l \phi_l^T D^{\frac{1}{2}} \end{aligned} \tag{1}$$

The vectors  $D^{\frac{1}{2}} \phi_l$ ,  $l = 1, \dots, n$  are orthogonal with lengths  $\psi_l^T D \psi_l = \sum_i \psi_{li}^2 d_i$  (put another way, orthonormal in  $\|\cdot\|_{\frac{1}{d}}$ ). This means that we can view the  $i$ 'th row of  $P$  as an expansion in the that basis with coefficients  $\lambda_l D^{-\frac{1}{2}} \phi_{li} = \lambda_l \psi_{li}$ . Consequently,

$$\|P_i - P_j\|_{\ell^2/d}^2 = \sum_l \lambda^2 (\psi_{li} - \psi_{lj})^2$$

. Analogously, for  $P^t$  we have

$$\|P_i^t - P_j^t\|_{\ell^2/d}^2 = \sum_l \lambda^{2t} (\psi_{li} - \psi_{lj})^2$$

□

## 2.1 Representation

The above suggests that we can represent a node  $v_i$  using a feature vector

$$\Psi_t(v_i) = (\lambda_1^t \psi_{1,i}^t, \dots, \lambda_L^t \psi_{L,i}^t)^T,$$

for some  $L \leq n$ .

Recall that by proposition 1.11,  $L_{\text{sim}}$  is positive semi-definite. This implies that the entire spectrum of  $L_{\text{rw}}$  lies between 1 and 0, and consequently also the spectrum of  $P$ . Since the spectrum decays, for each  $t$  only eigenvalues  $\lambda^t$  have significant effect on the diffusion distances, which means that we can use dimensionality  $L$  such that  $\lambda_L^t$  is sufficiently small, and achieve dimensionality reduction.

## Homework

1. Prove proposition 1.6.
2. Prove proposition 1.11.
3. Prove proposition 1.12.
4. Generate data in two clusters in 2d, plot the data, Laplacian eigenvalues and eigenvectors.
5. Generate two-moon data in 2D, and separate the data using  $k$ -means and Spectral clustering.

6. Let  $P$  be a Markov matrix. Show that  $P^2$  encodes 2-steps transition probabilities.
7. Prove proposition 2.1
8. Prove that if  $\phi$  is an eigenvector of  $A$  with eigenvalue  $\lambda$ , then  $\psi = D^{-\frac{1}{2}}\phi$  is an eigenvector of  $P$  with the same eigenvalue.