Lecture VIII: Spectral Clustering

Marina Meilă mmp@stat.washington.edu

Department of Statistics University of Washington

March 2022

Marina Meila (UW)



Graph clustering paradigms

MNCut and Random Walks

Embedding, the spectral mapping, lumpability



Similarity based clustering

• **Paradigm:** the features we observe are measures of similarity/dissimilarity between pairs of data points, e.g

	points	features		
Image segmentation	pixels	distance in color space or location, separated by a contour, belong to same texture		
Social network	people	friends, coworkers, phone calls, emails		
Text analysis	words	appear in same context		
The features are summarized by a single similarity measure S				

The features are summarized by a single similarity measure S_{ij}

• e.g
$$S_{ij} = e^{\sum_k \alpha_k \text{feature}_k(i,j)}$$
 for all points i, j

- We want to put points that are similar to each other in the same cluster, dissimilar points in different clusters
- Problem is often cast as a graph cut problem

• points = graph nodes, similarity
$$S_{ij}$$
 = weight of edge ij

٥

Paradigms for grouping

• Graph cuts

remove some edges \implies disconnected graph the groups are the connected components

- By "similar behavior" nodes i, j in the same group iff i, j "have the same pattern of connections" w.r.t other nodes
- By Embedding
- map nodes $V = \{1, 2, ..., n\} \longrightarrow \{x_1, x_2, ..., x_n\} \in \mathbb{R}^d$ then use standard classification and clustering methods

Definitions

• $V = \{1, 2, ..., n\}$

• node degree or volume

$$D_i = \sum_{j \in V} S_{ij}$$

• **volume** of cluster $C \subseteq V$

$$D_C = \sum_{i \in C} D_i$$

• cut between subsets $C, C' \subseteq V$

$$\sum_{i \in C} \sum_{j \in C'} S_{ij}$$

Multiway Normalized Cut of a partition Δ = {C_{1:K}} of V

$$MNCut(\Delta) = \sum_{k=1}^{K} \sum_{k' \neq k} \frac{Cut(C_k, C_{k'})}{D_{C_k}}$$

in particular, for K = 2,

$$MNCut(C, C') = Cut(C, C') \left(\frac{1}{D_C} + \frac{1}{D_{C'}}\right)$$

Motivation for MNCut



A random walks view



Define

$$P_{ij} = rac{S_{ij}}{D_i} \quad ext{for all } i, j \in V$$

- in matrix notation $P = D^{-1}S$ where $P = [P_{ij}], D = \text{diag}(D_1, \dots, D_n)$
- P defines a random walk over the graph nodes V

Grouping from the random walks point of view

• Idea: group nodes together if they transition in the same way to other clusters



i	P.,red	P _{i,yellow}
\bigcirc	1/5	4/5
\bigcirc	1/5	4/5
\bigcirc	1/5	4/5
0	1/5	4/5
\bigcirc	1/5	4/5
\bigcirc	1/5	4/5
\bigcirc	1/5	4/5
\circ	1/5	4/5
	2/3	1/3
	2/3	1/3
•	2/3	1/3
	2/3	1/3

... is the same as grouping by embedding

- embedding of V = mapping from V into \mathbb{R}^d
- Wanted: similar points embedded near each other ideally, points in the same cluster mapped to the same point in \mathbb{R}^d



Some questions

• Not all graphs embed perfectly



- How many dimensions do we need?
- Nice, but we need to know the clusters in advance...

Lumpability

- A vector v is piecewise constant w.r.t a clustering Δ iff $v_i = v_i$ whenever i, j in same $C \in \Delta$ Another look at Pic a piecewise $P_{\cdot,red} \equiv f^{red}$ constant function $P_{\cdot,yel} \equiv f^{yel}$ P_{i,yellow} red 4/5 fred 4/5 4/5 4/5 4/5 4/5 1/5 4/5 4/5 4/5 1/3 1/3 ۰ 1/3 • 1/3 1/3
 - **Theorem [Lumpability]**[Meila&Shi 2001] Let *S* be a similarity matrix and Δ a clustering with *K* clusters. Then *P* has *K* piecewise constant eigenvectors w.r.t Δ iff

$$\sum_{j \in C'} P_{ij} = R_{CC'} \text{ whenver } i \in C, \text{ for all } C, C' \in \Delta$$

The spectral mapping



Spectral clustering in a nutshell



The algorithm

Spectral clustering

An algorithm based on [?] and [?]. Spectral Clustering Algorithm

Input Similarity matrix S, number of clusters K

• Transform S: Set $D_i = \sum_{j=1}^n S_{ij}$, j = 1 : n the node degrees. Form the transition matrix $P = [P_{ij}]_{ij=1}^n$ with

$$P_{ij} \leftarrow S_{ij}/D_i$$
, for $i, j = 1 : n$

- **(a)** Compute the largest K eigenvalues $\lambda_1 = 1 \ge \lambda_2 \ge \ldots \ge \lambda_K$ and eigenvectors v_1, \ldots, v_K of P.
- Embed the data in principal subspace Let $V = [v_2 v_3 ... v_K] \in \mathbb{R}^{n \times K}$, $x_i \leftarrow i$ -th row of V. ● (orthogonal initialization) Find K initial centers by
 - () take μ_1 randomly from $x_1, \ldots x_n$
 - **2** for k = 2, ..., K set $\mu_k = \operatorname{argmin}_{x_i} \max_{k' < k} \mu_{k'}^T x_i$.
- **(a)** Run the K-means algorithm on the "data" $x_{1:n}$ starting from the centers $\mu_{1:K}$.

Properties of spectral clustering

- Arbitrary cluster shapes (main advantage)
- Elegant mathematically
- Practical up to medium sized problems
 - Running time (by Lanczos algorithm) $\mathcal{O}(nk)/iteration$.
- Works well when K known, not too large estimating K [?]
- Depend heavily on the similarity function (main problem) learning the similarities [?],[?],[?],[?]
- Outliers become separate clusters (user must adjust K accordingly!)
- Very popular, many variants which aim to improve on the above Diffusion maps [?]: normalize the eigenvectors λ^k_k ν^k
- Practical fix, when K large: only compute a fixed number of eigenvectors d < K. This avoids the effects of noise in lower ranked eigenvectors

Affinity propagation

- Idea Each item $i \in D$ finds an exemplar item $k \in D$ to "represent" it
- Affinity Propagation is to spectral clustering what Mean Shift is to K-means
- number of exemplars not fixed in advance
- quantities of interest
 - similarities s_{ij} , $i \neq j$ (given)
 - availability a_{ik} of k for i = how much support there is from other items for k to be an exemplar
 - responsibility r_{ik} that measures how fit is k to represent i, as compared to other possible candidates k'.
 - diagonal elements s_{ii} represent self-similarities
 - larger $s_{ii} \Rightarrow$ more likely *i* will become an exemplar \Rightarrow more clusters

Affinity Propagation

Affinity Propagation Algorithm [?]

Input Similarity matrix $S = [s_{ik}]_{ik=1}^n$, parameter $\lambda = 0.5$ Iterate the following steps until convergence

 $\begin{array}{l} \bullet \\ \bullet \\ a_{ik} \leftarrow 0 \text{ for } i, k = 1:n \\ \bullet \\ \bullet \\ \text{ for all } i \end{array}$

• Find the best exemplar for $i: s^* \leftarrow \max_k(s_{ik} + a_{ik}), A_i^* \leftarrow \max_k(s_{ik} + a_{ik})$ (can be a set of items)

g for all k update responsibilities

$$r_{ik} \leftarrow \begin{cases} s_{ik} - s^*, & \text{if } k \notin A_i^* \\ s_{ik} - \max_{k' \notin A_i^*} (s_{ik} + a_{ik}) & \text{otherwise} \end{cases}$$

I for all k update availabilities

$$\begin{array}{l} \bullet \quad a_{kk} \leftarrow \sum_{i \neq k} |r_{ik}|_{+} \text{ where } [r_{ik}]_{+} = r_{ik} \text{ if } r_{ik} > 0 \text{ and } 0 \text{ otherwise.} \\ \\ \bullet \quad \text{for all } i, a_{ik} \leftarrow \min\{0, r_{kk} + \sum_{i' \neq i, k} |r_{i'k}|_{+}\} \end{array}$$

• Assign an exemplar to *i* by $k(i) \leftarrow \underset{k'}{\operatorname{argmax}}(r_{ik'} + a_{ik'})$