

Spectral clustering and stochastic block model

Homework should be submitted via Gradescope, by Monday afternoon: the code will be communicated by an announcement on Canvas. This homework requires some

For getting credit for the class, you are required to present solutions of some of these homeworks during the first 15 minutes of class starting on 1/20. Please, sign up for (at least) one slot, and be sure that your explanation lasts 15 minutes (or less). For these presentations, you are free to choose whatever format you prefer (slides, typed notes, handwriting, ...).

This week, the presentations will be:

- Monday 1/25: Questions (a), (b), (c)
- Wednesday 1/27: Problem 2, points (d), (e), (f).

Spectral clustering

Given a graph $G = (V = [n], E)$, let \mathbf{A}_G be its adjacency matrix, $\mathbf{d}_G = \mathbf{A}_G \mathbf{1}$ the vector of degrees and $\hat{d} = \langle \mathbf{d}_G, \mathbf{1} \rangle / n$ denote its average degree. Also, denote by $\mathbf{D}_G = \text{diag}(\mathbf{d}_G)$ be the diagonal matrix containing the degrees.

Spectral clustering is a general approach to generate a partition of V in k clusters. It has many variants and possible implementations, but at high level it proceeds as follows.

SPECTRAL CLUSTERING

Input : Data G , number of clusters k

Output : Vertex labels $\hat{\sigma} = (\hat{\sigma}_1, \dots, \hat{\sigma}_n)$, $\hat{\sigma}_i \in \{1, \dots, k\}$

- 1: Construct matrix $\mathbf{M}_G \in \mathbb{R}^{n \times n}$ from \mathbf{A}_G ;
- 2: Compute the top $(k-1)$ eigenvectors $\mathbf{v}_1(\mathbf{M}_G), \dots, \mathbf{v}_{k-1}(\mathbf{M}_G)$;
- 3: Let $\mathbf{x}_i \in \mathbb{R}^{k-1}$ the i -th row of the matrix $\mathbf{V} \in \mathbb{R}^{n \times (k-1)}$ whose i -th column is $\mathbf{v}_i(\mathbf{M}_G)$;
- 4: Cluster points $\{\mathbf{x}_i\}_{i \in [n]}$ into k groups, assigning labels $\hat{\sigma}_i \in \{1, \dots, k\}$;
- 5: Return vector $\hat{\sigma} = (\hat{\sigma}_1, \dots, \hat{\sigma}_n)$

Step 1 (construction of the matrix \mathbf{M}_G) and step 4 (clustering of points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^{k-1}$) admit several concrete implementation. We will consider a few of them.

- We will consider two specific choices for the matrix \mathbf{M}_G at step 1:

1. The centered adjacency matrix $\mathbf{M}_G = \mathbf{A}_G^{\text{cen}}$, where

$$\mathbf{A}_G^{\text{cen}} = \mathbf{A}_G - \frac{\hat{d}}{n} \mathbf{1} \mathbf{1}^\top. \quad (1)$$

2. The normalized (subtracted) Laplacian. Letting $\mathbf{z}_G = \mathbf{d}_G^{1/2} / \sqrt{n \hat{d}}$ (where $\mathbf{d}_G^{1/2}$ is the vector of square roots of degrees), we let $\mathbf{M}_G = \mathbf{L}_G$, where

$$\mathbf{L}_G = \mathbf{D}_G^{-1/2} \mathbf{A}_G \mathbf{D}_G^{-1/2} - \mathbf{z}_G \mathbf{z}_G^\top. \quad (2)$$

(We will study some interesting properties of the graph Laplacian later in the class.)

- We will consider two possibilities for the clustering step 4:

1. For $k = 2$, the points x_1, \dots, x_n generate at step 3 are real numbers. We partition them by thresholding at 0:

$$\hat{\sigma}_i = \begin{cases} 1 & \text{if } x_i \geq 0, \\ 2 & \text{if } x_i < 0, \end{cases} \quad (3)$$

2. For general k , we can cluster points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ using the k -means algorithm. (You can read about it and implement it, or read about it and use an available implementation.)

We will test spectral clustering in two different settings: (i) Synthetic data generated according to the stochastic block model; (ii) A real data set concerning political blogs.

In both cases, if $\boldsymbol{\sigma}$ denotes the true vertex labels, we will evaluate a proposed clustering $\hat{\boldsymbol{\sigma}}$ using the overlap

$$Q(\boldsymbol{\sigma}, \hat{\boldsymbol{\sigma}}) \equiv \frac{k}{k-1} \max_{\pi \in S_k} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\pi(\hat{\sigma}_i) = \sigma_i\}} - \frac{1}{k} \right\} \quad (4)$$

Here $\pi : [k] \rightarrow [k]$ is a permutation of the vertex labels $1, \dots, k$ and accounts for the fact that labels are arbitrary. S_k is the set for all such permutations. Note that $Q \in [0, 1]$ with $Q \approx 1$ corresponding to nearly perfect clustering, and $Q \approx 0$ to random clusters.

Stochastic block model

Recall that a graph from the two-groups symmetric block model $G = (V, E) \sim \mathbf{G}(n, a/n, b/n)$ has vertices $V = [n]$ and is generated as follows. Attribute to each vertex $i \in [n]$ a true label $\sigma_i \in \{+1, -1\}$ uniformly at random. Conditional on these labels, add edges independently with probabilities

$$\mathbb{P}\{(i, j) \in E | \boldsymbol{\sigma}\} = \begin{cases} a/n & \text{if } \sigma_i = \sigma_j, \\ b/n & \text{if } \sigma_i \neq \sigma_j. \end{cases} \quad (5)$$

- (a) Generate random graphs from the stochastic block model with $n = 10,000$, $b = 10$ and each of $a \in \{10, 15, 20, 25, 30\}$. In each case, plot the histogram of eigenvalues of the centered adjacency matrix $\mathbf{A}_G^{\text{cen}}$, along with the largest eigenvalue.
- (b) Generate 20 random graphs for each of the choices of parameters at the previous point (hence a total of 100 instances). Apply the spectral clustering algorithm described in the previous section, with $k = 2$, $\mathbf{M}_G = \mathbf{A}_G^{\text{cen}}$ the centered adjacency matrix, and using Eq. (3) to produce the vertex labels. Report the empirical overlap for each choice of the parameters (n, a, b) , averaged over the 20 samples.
- (c) Compare your results with the theory we studied in class.

A real dataset

We will study a political blogs dataset first compiled for the paper

Lada A. Adamic and Natalie Glance, *The political blogosphere and the 2004 US Election*, in Proceedings of the WWW-2005 Workshop on the Weblogging Ecosystem (2005)

The dataset `polblogs` contains a graph with $n = 1490$ vertices ('nodes') corresponding to political blogs. Each vertex has a 0 or -1 label (indicated as 'value') corresponding to the political orientation of that blog. We will consider this as the true label σ_i and try to reconstruct the vector $\boldsymbol{\sigma}$ from the graph using the spectral clustering approach introduced above.

For each of the implementations given below, report the resulting overlap $Q(\boldsymbol{\sigma}, \hat{\boldsymbol{\sigma}})$:

- (d) Use SPECTRAL CLUSTERING with $\mathbf{M}_G = \mathbf{A}_G^{\text{cen}}$ the centered adjacency matrix and using simple thresholding (as in Eq. (3)) to produce the clusters.
- (e) Repeat the experiment at the previous point, but replacing the centered adjacency matrix with the normalized Laplacian \mathbf{L}_G .
- (f) As in the previous point, we use the Laplacian matrix \mathbf{L}_G . However, here we assume that the number of clusters in G is estimated to be $k = 3$, and hence we use the top 2 principal components $\mathbf{v}_1(\mathbf{L}_G)$, $\mathbf{v}_2(\mathbf{L}_G)$. Produce a scatterplot of the datapoints $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^2$ produced at step 3 of the algorithm (it is useful to color-code them, using the true labels).
Use k -Means (with $k = 3$) to cluster the datapoints $\mathbf{x}_1, \dots, \mathbf{x}_n$. How would you evaluate the resulting clustering?