

A fast incremental spectral clustering algorithm with cosine similarity

Guangliang Chen

Associate Professor of Statistics & Data Science
Hope College, Holland, Michigan, USA
cheng@hope.edu

ICDM 2023 IncrLearn Workshop
Shanghai, China

Introduction

Spectral clustering is a modern, powerful clustering approach. It uses the eigenvectors of a normalized graph Laplacian for embedding the data into a low-dimensional space for easy clustering.

However, it is well known to face two major challenges:

- scalability (speed and memory),
- out of sample extension.

We present a memory and speed efficient spectral clustering algorithm in the setting of *cosine similarity* that only uses the following efficient linear algebra operations:

- elementwise manipulation
- matrix-vector multiplication
- low-rank SVD

A spectral clustering algorithm

There are different formulations of spectral clustering; here we present the version by Ng, Jordan and Weiss (2001).

Input: Data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, # clusters k , scale parameter σ

Output: Clusters $\mathcal{C}_1, \dots, \mathcal{C}_k$

1: Construct a pairwise similarities matrix

$$\mathbf{W} = (w_{ij}), \quad w_{ij} = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma^2)}, \quad i \neq j$$

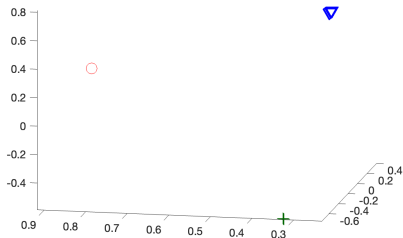
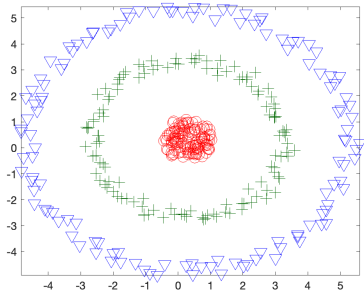
2: Find the row sums of \mathbf{W} and use them to define a diagonal (degree) matrix $\mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1})$. Let $\tilde{\mathbf{W}} = \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$.

3: Find the k largest eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ of $\tilde{\mathbf{W}}$ and use them to form an embedding matrix

$$\mathbf{X} \mapsto \mathbf{V} = [\mathbf{v}_1 \dots \mathbf{v}_k] \in \mathbb{R}^{n \times k}.$$

4: Apply k -means to group the rows of \mathbf{V} (after being normalized to have unit length) into k clusters.

A demonstration



Computational challenges

Spectral clustering is flexible, accurate and powerful, and has been successfully applied to tasks such as image segmentation, documents clustering and social community detection.

However, it is computationally intensive when being applied to large data sets.

- **Memory requirement:** $\mathcal{O}(n^2)$
- **Computational cost:**
 - (a) Construction of \mathbf{W} : $\mathcal{O}(n^2 d)$
 - (b) Decomposition of \mathbf{W} : $\mathcal{O}(n^3)$

Consequently, there has been a lot of work on making it scalable in **time** and/or **memory** to large data sets.

Speed scalability (ICPR18')

Given (low-dimensional or row-sparse) data $\mathbf{X} \in \mathbb{R}^{n \times d}$ with L_2 -normalized rows, the cosine similarity matrix is

$$\mathbf{W} = \mathbf{X}\mathbf{X}^T - \mathbf{I}.$$

First, we can compute the degree matrix \mathbf{D} directly from \mathbf{X} (without needing to construct \mathbf{W}):

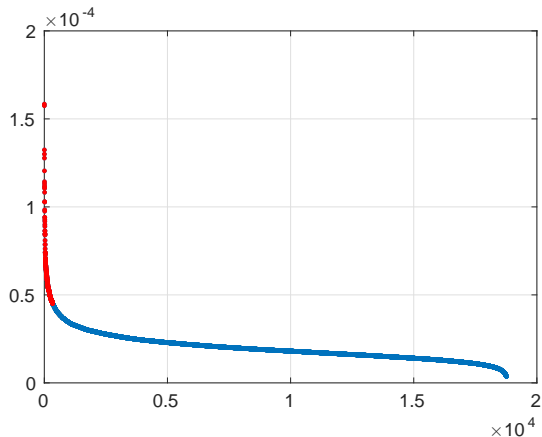
$$\mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1}) = \text{diag}((\mathbf{X}\mathbf{X}^T - \mathbf{I})\mathbf{1}) = \text{diag}(\mathbf{X}(\mathbf{X}^T\mathbf{1}) - \mathbf{1}_n).$$

Next, we write

$$\tilde{\mathbf{W}} = \mathbf{D}^{-1/2}(\mathbf{X}\mathbf{X}^T - \mathbf{I})\mathbf{D}^{-1/2} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T - \mathbf{D}^{-1}, \quad \tilde{\mathbf{X}} = \mathbf{D}^{-1/2}\mathbf{X}.$$

Finally, after removing a small fraction (α) of low-degree points (in order to make \mathbf{D}^{-1} nearly constant diagonal), we use the left singular vectors of $\tilde{\mathbf{X}}$ to approximate the eigenvectors $\tilde{\mathbf{U}}$ of $\tilde{\mathbf{W}}$.

Demonstration: Diagonal entries of \mathbf{D}^{-1} on a large data set (20 newsgroups)



The small fraction of points in red have relatively lowest degrees, typically being **outliers**. They are discarded so that **the reduced \mathbf{D}^{-1} matrix is nearly constant diagonal** (needed by the approximation).

Memory scalability

Suppose we encounter

- a **massive** data set, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$, that is too large to be fully loaded into computer memory, but we have access to small batches of the data through sampling, or
- an **online** data set that arrives sequentially, one point at a time.

In both settings, we would like to perform spectral clustering with the cosine similarity on the whole data set.

We address the memory challenge in the following steps:

- 1 We first develop a base algorithm using only a **single batch** of data to learn the nonlinear embedding and clustering rule.

Memory scalability

Suppose we encounter

- a **massive** data set, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$, that is too large to be fully loaded into computer memory, but we have access to small batches of the data through sampling, or
- an **online** data set that arrives sequentially, one point at a time.

In both settings, we would like to perform spectral clustering with the cosine similarity on the whole data set.

We address the memory challenge in the following steps:

- 1 We first develop a base algorithm using only a **single batch** of data to learn the nonlinear embedding and clustering rule.
- 2 We then present an **incremental learning** procedure to continuously refine them as more batches of data are drawn.

Memory scalability

Suppose we encounter

- a **massive** data set, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$, that is too large to be fully loaded into computer memory, but we have access to small batches of the data through sampling, or
- an **online** data set that arrives sequentially, one point at a time.

In both settings, we would like to perform spectral clustering with the cosine similarity on the whole data set.

We address the memory challenge in the following steps:

- ① We first develop a base algorithm using only a **single batch** of data to learn the nonlinear embedding and clustering rule.
- ② We then present an **incremental learning** procedure to continuously refine them as more batches of data are drawn.
- ③ Lastly, we explain how to **extend** both the nonlinear embedding and clustering rule learned on the training data to the rest of the data in \mathbf{X} , as they gradually become available.

Single batch learning

Assume a **small batch** of data of size $s \ll n$, denoted $\mathbf{X}_s \in \mathbb{R}^{s \times d}$, that has become available through sampling. We would like to take advantage of the computing procedure developed in previous work on the *speed scalability* of spectral clustering with cosine similarity.

The following equation establishes the bridge between the full data \mathbf{X} and the sample \mathbf{X}_s :

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{X}^T \mathbf{D}^{-1} \mathbf{X} = \sum_{i=1}^n \frac{1}{d_i} \mathbf{x}_i \mathbf{x}_i^T \approx \frac{n}{s} \sum_{i=1}^s \frac{1}{d_i} \mathbf{x}_i \mathbf{x}_i^T = \frac{n}{s} \tilde{\mathbf{X}}_s^T \tilde{\mathbf{X}}_s,$$

where $\tilde{\mathbf{X}}_s = \mathbf{D}_s^{-1/2} \mathbf{X}_s$ and \mathbf{D}_s represent the restrictions of $\tilde{\mathbf{X}}$ and \mathbf{D} to the sample \mathbf{X}_s , respectively:

$$\mathbf{D}_s = \text{diag}(\mathbf{d}_s), \quad \mathbf{d}_s = \mathbf{X}_s (\mathbf{X}^T \mathbf{1}) - \mathbf{1}_s \approx \frac{n}{s} \mathbf{X}_s (\mathbf{X}_s^T \mathbf{1}_s) - \mathbf{1}_s.$$

The previous equation, $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \frac{n}{s} \tilde{\mathbf{X}}_s^T \tilde{\mathbf{X}}_s$, allows us to use the rank- k SVD of the sample,

$$\tilde{\mathbf{X}}_s \approx \tilde{\mathbf{U}}_s \tilde{\Sigma}_s \tilde{\mathbf{V}}_s^T,$$

to estimate the nonlinear embedding used in the speed-scalable spectral clustering algorithm (ICPR18'), i.e.,

$$\mathbf{Y} := \tilde{\mathbf{U}} = \tilde{\mathbf{X}} \left(\tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right) \quad \leftarrow \quad \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T \approx \tilde{\mathbf{V}}_s \tilde{\mathbf{V}}_s^T \quad \text{and} \quad \tilde{\Sigma} \approx \sqrt{\frac{n}{s}} \tilde{\Sigma}_s$$

and apply it to the batch data $\mathbf{X}_s \in \mathbb{R}^{s \times d}$ as follows:

$$\mathbf{Y}_s := \tilde{\mathbf{X}}_s \left(\tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right) \approx \tilde{\mathbf{X}}_s \tilde{\mathbf{V}}_s \left(\sqrt{\frac{n}{s}} \tilde{\Sigma}_s \right)^{-1} = \sqrt{\frac{s}{n}} \tilde{\mathbf{X}}_s \tilde{\mathbf{V}}_s \tilde{\Sigma}_s^{-1} \in \mathbb{R}^{s \times k}.$$

We then perform k -means clustering in the \mathbf{Y}_s space.

Incremental learning

The embedding obtained on the initial batch of data can be refined by using more samples.

Now, consider a second batch of data \mathbf{X}_t that is independently sampled from the population, and let the combined sample be

$$\mathbf{X}_{s+t} = \begin{bmatrix} \mathbf{X}_s \\ \mathbf{X}_t \end{bmatrix} \in \mathbb{R}^{(s+t) \times d}.$$

By the same reasoning,

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \approx \frac{n}{s+t} \tilde{\mathbf{X}}_{s+t}^T \tilde{\mathbf{X}}_{s+t}.$$

However, we will not compute the SVD of $\tilde{\mathbf{X}}_{s+t}$ directly as there is a more efficient way of estimate the right singular vectors $\tilde{\mathbf{V}}$ of $\tilde{\mathbf{X}}$.

Write

$$\tilde{\mathbf{X}}_{s+t}^T \tilde{\mathbf{X}}_{s+t} = \begin{bmatrix} \tilde{\mathbf{X}}_s^T & \tilde{\mathbf{X}}_t^T \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_s \\ \tilde{\mathbf{X}}_t \end{bmatrix} = \tilde{\mathbf{X}}_s^T \tilde{\mathbf{X}}_s + \tilde{\mathbf{X}}_t^T \tilde{\mathbf{X}}_t.$$

Using the rank- k SVD of $\tilde{\mathbf{X}}_s$ (which has already been obtained earlier), we have the following approximation

$$\tilde{\mathbf{X}}_{s+t}^T \tilde{\mathbf{X}}_{s+t} \approx \tilde{\mathbf{V}}_s \tilde{\Sigma}_s^2 \tilde{\mathbf{V}}_s^T + \tilde{\mathbf{X}}_t^T \tilde{\mathbf{X}}_t = \underbrace{\begin{bmatrix} \tilde{\mathbf{V}}_s \tilde{\Sigma}_s & \tilde{\mathbf{X}}_t^T \end{bmatrix}}_{\tilde{\mathbf{X}}_{k,t}^T} \underbrace{\begin{bmatrix} \tilde{\Sigma}_s \tilde{\mathbf{V}}_s^T \\ \tilde{\mathbf{X}}_t \end{bmatrix}}_{\tilde{\mathbf{X}}_{k,t}} = \tilde{\mathbf{X}}_{k,t}^T \tilde{\mathbf{X}}_{k,t}$$

We can thus use the rank- k SVD of $\tilde{\mathbf{X}}_{k,t} \in \mathbb{R}^{(k+t) \times d}$ (which is smaller in size than $\tilde{\mathbf{X}}_{s+t} \in \mathbb{R}^{(s+t) \times d}$),

$$\tilde{\mathbf{X}}_{k,t} \approx \tilde{\mathbf{U}}_{k,t} \tilde{\Sigma}_{k,t} \tilde{\mathbf{V}}_{k,t}^T$$

to estimate $\tilde{\mathbf{V}}_{s+t}$, $\tilde{\Sigma}_{s+t}$ of $\tilde{\mathbf{X}}_{s+t}$ and embed \mathbf{X}_t accordingly:

$$\mathbf{Y}_t := \sqrt{\frac{s+t}{n}} \tilde{\mathbf{X}}_t \tilde{\mathbf{V}}_{s+t} \tilde{\Sigma}_{s+t}^{-1} \approx \sqrt{\frac{s+t}{n}} \tilde{\mathbf{X}}_t \tilde{\mathbf{V}}_{k,t} \tilde{\Sigma}_{k,t}^{-1}$$

The stopping criterion

Starting with an initial size s and fixing a step size t , we can iterate the above learning process:

$$\mathbf{s} \leftarrow \mathbf{s} + t, \quad \tilde{\mathbf{V}}_{\mathbf{s}} \leftarrow \tilde{\mathbf{V}}_{\mathbf{s}+t}$$

Convergence is determined based on Grassmannian metric:

$$g_{\mathbf{s}} = \left\| \tilde{\mathbf{V}}_{\mathbf{s}+t} \tilde{\mathbf{V}}_{\mathbf{s}+t}^T - \tilde{\mathbf{V}}_{\mathbf{s}} \tilde{\mathbf{V}}_{\mathbf{s}}^T \right\|_F = \sqrt{2k - 2 \left\| \tilde{\mathbf{V}}_{\mathbf{s}+t}^T \tilde{\mathbf{V}}_{\mathbf{s}} \right\|_F^2} = \sqrt{2 \sum_{j=1}^k \sin^2 \theta_j},$$

where $0 \leq \theta_1 \leq \dots \leq \theta_k \leq \frac{\pi}{2}$ are the principal angles between the column spaces of $\tilde{\mathbf{V}}_{\mathbf{s}+t}$ and $\tilde{\mathbf{V}}_{\mathbf{s}}$.

Empirically, we set s such that all $\theta_j \leq \theta_0$, i.e.,

$$g_{\mathbf{s}} < \sqrt{2 \cdot k \cdot \sin^2 \theta_0} = \sqrt{2k} \sin \theta_0.$$

Out of sample extension

Any new point, say $\mathbf{x}_0 \in \mathbb{R}^d$, is embedded as follows:

$$\mathbf{y}_0 = \sqrt{\frac{s}{n}} \left(d_0^{-1/2} \mathbf{x}_0^T \right) \tilde{\mathbf{V}}_s \tilde{\Sigma}_s^{-1} \in \mathbb{R}^k,$$

where

$$d_0 = \mathbf{x}_0^T \sum_{i=1}^n \mathbf{x}_i - 1 \approx \frac{n}{s} \mathbf{x}_0^T (\mathbf{X}_s^T \mathbf{1}_s) - 1.$$

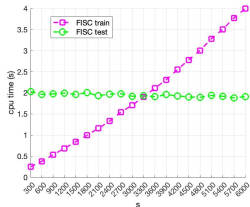
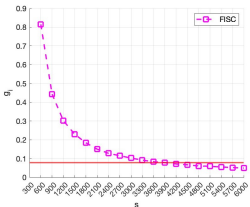
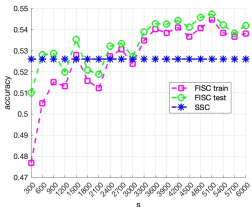
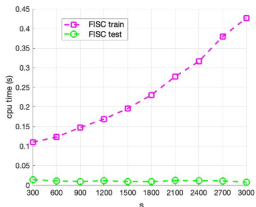
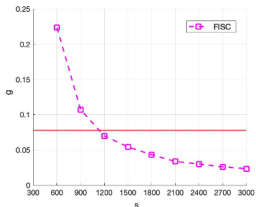
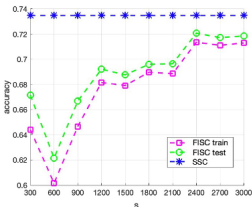
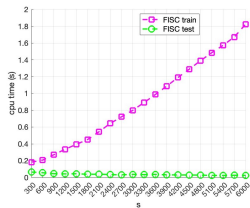
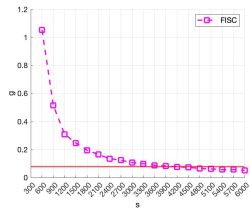
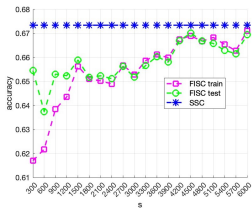
We then assign \mathbf{y}_0 to the closest centroids in embedding space.

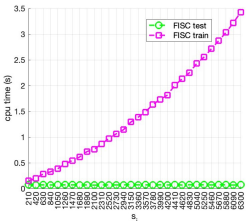
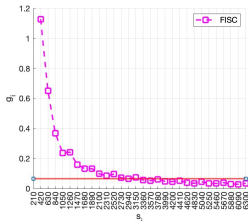
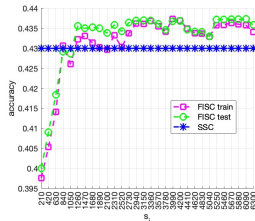
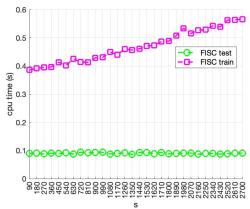
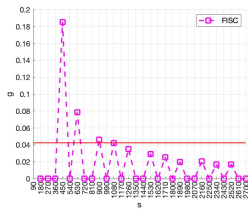
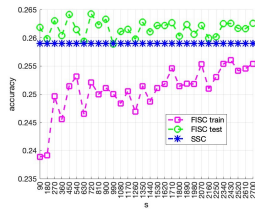
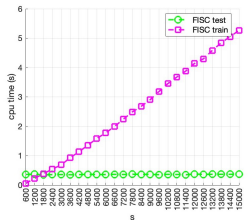
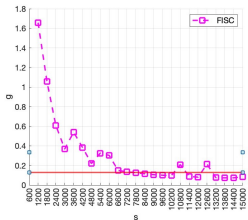
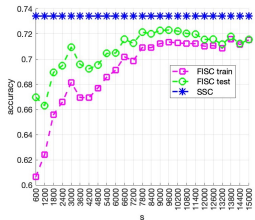
Experimental study

We conduct experiments on the following benchmark datasets to evaluate the performance of our algorithm - Fast, incremental spectral clustering (FISC), in terms of clustering accuracy and CPU time.

Data sets	n	p	k
usps	9,298	256	10
pendigit	10,992	16	10
mnist	70,000	184	10
20news	18,768	55,570	20
protein	24,387	357	3
covtype	581,012	54	7

All experiments were conducted with MATLAB R2021b on a desktop computer with 32 GB of RAM and a CPU with 4 cores.





Thank you for your attention!

Conclusions: We proposed a memory-efficient spectral clustering algorithm that uses small batches of data to effectively learn spectral embedding and clustering maps, as well as the out-of-sample extensions.

In particular, we introduced an incremental procedure that monitors the Grassmannian distances between consecutive iterations to check for convergence.

Future work:

- Statistical analysis of the sampling procedure
- Extension to other kinds of similarity such as Gaussian

Acknowledgement: This was joint work with Ran Li (my former graduate student at San José State University, who now works at Wells Fargo).