



Jet Clustering with Spectral Clustering

Henry Day-Hall¹ Supervisors: Prof. Claire Shepherd-Themistocleous^{1,2}, Prof. Stefano Moretti¹, Prof. Srinandan Dasmahapatra¹, Dr. Emmanuel Olaiya²

> ¹University of Southampton, UK ²Rutherford Appleton Laboratory, UK

> > January 6, 2020

Table of Contents



Introduction

Results

Method

Jets

Southampton





A good jet cluttering algorithm will accurately match the kinematics of the partons chosen as tags.



- A good jet cluttering algorithm will accurately match the kinematics of the partons chosen as tags.
- ► This accuracy should vary smoothly with the cut-off parameter.



- A good jet cluttering algorithm will accurately match the kinematics of the partons chosen as tags.
- ► This accuracy should vary smoothly with the cut-off parameter.
- ► The jets formed should replicate higher level shape variables.

Results

Southampton



jet class = Spectral 0.0 < mean num jets < 20.0 Laplacien = unnormalised AffinityType = linear WithLaplacienScaling = False ExponentMultiplier = 1



Many attempts have been made to write a 'good' clustering algorithm. Most of them are not hierarchical, they are based on fitting a predefined model. This poses a challenge for jet clustering, we do not have a predefined number of clusters.

Clustering comparison

Southampton



Figure: Taken from

https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-nee

Aim of clustering



Let our points be nodes of a graph and the vertices carry a measure of the affinity, $a_{i,j}$.



Aim of clustering



We wish to split the points such that the severed affinities are minimised.



Often the optimum split by this metric will isolate one point. To avoid this small clusters are penalised.



These criteria result in RatioCut. If $W(A, B) = \sum_{i \in A, j \in B} a_{i,j}$ is the sum of the affinities that cross from A to B, and |A| is the number of nodes in A;

RatioCut
$$(A_1, A_2, \dots A_n) \equiv \frac{1}{2} \sum_{i=1}^n \frac{W(A_i, \bar{A}_i)}{|A_i|}$$

In the case of disconnected components (with zero affinity between clusters) this can be solved for with the eigenvalues of the matrix known as the graph Laplacien.



Let us imagine a graph, disconnected in n clusters.





Let us imagine a graph, disconnected in n clusters. Membership of cluster k is determined by the indicator vector h_k ;

$$h_{i,k} = \begin{cases} 1/\sqrt{|A_k|}, & \text{if point } i \in A_k \\ 0, & \text{otherwise} \end{cases}$$

The graph is represented by the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & -a_{1,3} & \dots \\ -a_{1,2} & \sum a_{2,i} & -a_{2,3} \\ -a_{1,3} & -a_{2,3} & \sum a_{3,i} \\ \vdots & & \ddots \end{bmatrix}$$

Then

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$



Let us imagine a graph, disconnected in n clusters. Membership of cluster k is determined by the indicator vector h_k ;

$$h_{i,k} = \begin{cases} 1/\sqrt{|A_k|}, & \text{if point } i \in A_k \\ 0, & \text{otherwise} \end{cases}$$

The graph is represented by the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & -a_{1,3} & \dots \\ -a_{1,2} & \sum a_{2,i} & -a_{2,3} \\ -a_{1,3} & -a_{2,3} & \sum a_{3,i} \\ \vdots & & \ddots \end{bmatrix}$$

Then

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$



Let us imagine a graph, disconnected in n clusters. Membership of cluster k is determined by the indicator vector h_k ;

$$h_{i,k} = \begin{cases} 1/\sqrt{|A_k|}, & \text{if point } i \in A_k \\ 0, & \text{otherwise} \end{cases}$$

The graph is represented by the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & -a_{1,3} & \dots \\ -a_{1,2} & \sum a_{2,i} & -a_{2,3} \\ -a_{1,3} & -a_{2,3} & \sum a_{3,i} \\ \vdots & & \ddots \end{bmatrix}$$

Then

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$

Then stack the of all clusters together

$$h'_k Lh_k = (H'LH)_{kk}$$

and the RatioCut aim discribed earlier is the trace;

$$\mathsf{RatioCut}(A_1, A_2, \dots A_n) \equiv \frac{1}{2} \sum_{i=1}^n \frac{W(A_i, \bar{A}_i)}{|A_i|} = \mathsf{Tr}(H'LH)$$

Where H'H = I. Trace minimisation in this form is done by finding the eigenvectors of L with smallest eigenvalues.

Generalising this to a graph that is not disconnected is just relaxing the requirements on the form of the indicator vectors; h_k .

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$

Then stack the of all clusters together

$$h'_k Lh_k = (H'LH)_{kk}$$

and the RatioCut aim discribed earlier is the trace;

$$\mathsf{RatioCut}(A_1, A_2, \dots A_n) \equiv \frac{1}{2} \sum_{i=1}^n \frac{W(A_i, \bar{A}_i)}{|A_i|} = \mathsf{Tr}(H'LH)$$

Where H'H = I. Trace minimisation in this form is done by finding the eigenvectors of L with smallest eigenvalues.

Generalising this to a graph that is not disconnected is just relaxing the requirements on the form of the indicator vectors; h_k .

$$h'_{k}Lh_{k} = \frac{1}{|A_{k}|} \sum_{i \in A_{k}, j \in A_{k}} \left(\delta_{i,j} \sum_{l} a_{l,i} - a_{i,j} \right) = \frac{W(A_{k}, \bar{A}_{k})}{|A_{k}|}$$

Then stack the of all clusters together

$$h'_k Lh_k = (H'LH)_{kk}$$

and the RatioCut aim discribed earlier is the trace;

$$\mathsf{RatioCut}(A_1, A_2, \dots A_n) \equiv \frac{1}{2} \sum_{i=1}^n \frac{W(A_i, \bar{A}_i)}{|A_i|} = \mathsf{Tr}(H'LH)$$

Where H'H = I. Trace minimisation in this form is done by finding the eigenvectors of L with smallest eigenvalues.

Generalising this to a graph that is not disconnected is just relaxing the requirements on the form of the indicator vectors; h_k .



To find n clusters from m points;

1. Identify affinities between all points; $a_{i,j}$.

Southampton

To find n clusters from m points;

- 1. Identify affinities between all points; $a_{i,j}$.
- 2. Construct the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & \dots \\ -a_{1,2} & \sum a_{2,i} & \\ \vdots & & \ddots \end{bmatrix}$$

Southampton

To find n clusters from m points;

- 1. Identify affinities between all points; $a_{i,j}$.
- 2. Construct the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & \dots \\ -a_{1,2} & \sum a_{2,i} & \\ \vdots & & \ddots \end{bmatrix}$$

3. Calculate the eigenvectors v of L corresponding to the n+1 smallest eigenvalues.

Southampton

To find n clusters from m points;

- 1. Identify affinities between all points; $a_{i,j}$.
- 2. Construct the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & \dots \\ -a_{1,2} & \sum a_{2,i} & \\ \vdots & & \ddots \end{bmatrix}$$

- 3. Calculate the eigenvectors v of L corresponding to the n+1 smallest eigenvalues.
- 4. Stack the eigenvectors (aside from the first) v into a matrix E that is n by m. Call E the eigenspace, each point in the original dataset is represented by one row.

Southampton

To find n clusters from m points;

- 1. Identify affinities between all points; $a_{i,j}$.
- 2. Construct the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & \dots \\ -a_{1,2} & \sum a_{2,i} & \\ \vdots & & \ddots \end{bmatrix}$$

- 3. Calculate the eigenvectors v of L corresponding to the n+1 smallest eigenvalues.
- 4. Stack the eigenvectors (aside from the first) v into a matrix E that is n by m. Call E the eigenspace, each point in the original dataset is represented by one row.
- 5. Cluster in the eigenspace, E, using knn.

Physics Process



To find ? clusters from m points;

- 1. Identify affinities between all points; $a_{i,j}$.
- 2. Construct the graph Laplacien;

$$L = \begin{bmatrix} \sum a_{1,i} & -a_{1,2} & \dots \\ -a_{1,2} & \sum a_{2,i} & \\ \vdots & & \ddots \end{bmatrix}$$

- 3. Calculate the eigenvectors v of L corresponding to the q + 1 smallest eigenvalues.
- 4. Stack the eigenvectors (aside from the first) v into a matrix E that is q by m. Call E the eigenspace, each point in the original dataset is represented by one row.
- 5. Cluster in the eigenspace, E, using with a hierarchical method.

Conclusions



This is a well motivated clustering method.

- ► The best hyperparameters need to be identified.
- ► It should be tested for IRC safety.
- It's replication of event shape variables should be tested.

These hurdles aside, the method shows potential when compared to traditional jet clustering algorithms.

Thank you for listening.