

Graph Learning from Smooth Signals

Gurusha Juneja
2019EE10480

April 2022

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Course Coordinator

Contents

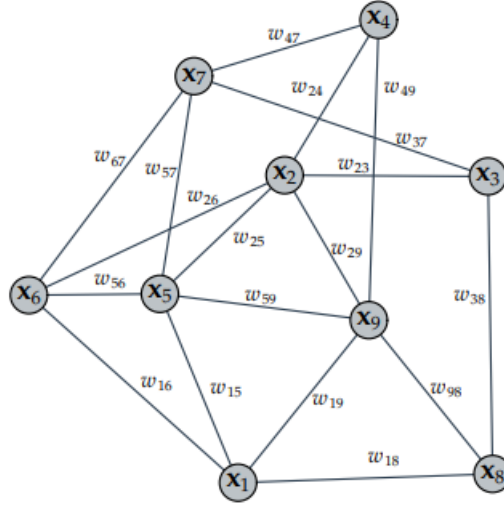
1	Introduction	2
2	Schematic of Graph Learning	2
3	Quantifying Smoothness	3
3.1	Learning graph under constraints	3
3.2	How to approach an optimisation problem?	4
3.3	How to enforce sparsity?	4
4	Graph Clustering	5
4.1	Constrained Laplacian Rank for Clustering	6
4.2	Kay-Fan Theorem	6
4.3	Solving for S and F	6
5	Spectral clustering Algorithm	7
6	Structure Learning	7
6.1	What is structured graph learning?	7
6.2	Spectral properties of structured graphs	8
6.3	Incorporating Spectral Properties into graph learning framework	9
6.4	Proposed Unified framework for Graph structure Learning	9
6.4.1	Formulation 1: Structure via Laplacian spectral constraints	9
6.4.2	Formulation 2: Structure via Adjacency spectral constraints	10
6.4.3	Formulation 3: Structure via joint spectral constraints . .	10
7	References	10

1 Introduction

Graph learning has become an important task with various applications. In these notes, we go through the algorithms to construct graph on a given dataset taking into consideration various constraints which help us make better connected graphs depicting real-life scenarios. One such way is learning graph from smooth signals, where we try to minimise the smoothness. Further in these notes, we explore how to approach problems that are highly combinatorial in nature, like learning structured graphs, and how using spectral properties of graphs can help us in doing so.

2 Schematic of Graph Learning

Given a data matrix $\mathbf{X} \in \mathbf{R}^{n \times p} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p]$, each column is assumed to represent one node on a graph and our goal is to obtain relationships between them as the weight of edges between these nodes as shown below,



Our goal is to obtain the weight matrix \mathbf{W} , where the weight of the edge between node i and node j is given by w_{ij} . The connectivity matrix \mathbf{C} , the weight matrix \mathbf{W} , the laplacian matrix \mathbf{L} are defined as follows,

$$[\mathbf{C}]_{ij} = \begin{cases} 1 & (i, j) \in E \\ 0 & otherwise \end{cases} \quad [\mathbf{W}]_{ij} = \begin{cases} w_{ij} & (i, j) \in E \\ 0 & otherwise \end{cases} \quad [\mathbf{L}]_{ij} = \begin{cases} -w_{ij} & (i, j) \in E \\ \sum_{j=1}^p w_{ij} & i = j \\ 0 & otherwise \end{cases} \quad (1)$$

Relation between \mathbf{W} and \mathbf{L} can be given as, $\mathbf{L} = \mathbf{D} - \mathbf{W}$, where $\mathbf{D} = \text{Diag}(\mathbf{W} \cdot \mathbf{1})$. Based on these matrices, various graph construction algorithms can be worked out. For example, sample correlation, where two nodes are connected if their pairwise correlation is greater than a given threshold, similarity based functions where $[\mathbf{W}]_{ij} = \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2})$ and $w_{ii} = 0$.

However, some solutions can be degenerate values like $w_{ij} = 0 \forall i, j$, hence some sort of regularisation is required, which we do by learning graphs from smooth signals.

3 Quantifying Smoothness

The energy of a system is defined as $\text{tr}(\mathbf{X}\mathbf{L}\mathbf{X}^T) = \sum_{i,j} w_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$. Smaller the value of the energy, smoother the signal \mathbf{X} over the graph \mathcal{G} . When \mathcal{G} is not available, we can learn it directly from \mathbf{X} and simultaneously minimise the smoothness by adding some regularisation term denoted by $h(\mathbf{L})$. Hence, the optimal laplacian is given by,

$$\hat{\mathbf{L}} := \underset{\mathbf{L}}{\text{argmin}} \text{tr}(\mathbf{X}\mathbf{L}\mathbf{X}^T) + \lambda h(\mathbf{L}) \quad (2)$$

$h(\mathbf{L})$ can be any function like $\|\mathbf{L}\|_1, \|\mathbf{L}\|_F^2, \log(\det(\mathbf{L}))$ etc.

Next we see how to formulate this problem as a convex optimisation problem and introduce some constraints.

3.1 Learning graph under constraints

Here we will introduce some constraints on the weight matrix so as to ensure that our solution becomes realistic, introducing the following constraints,

- Bounded weights: $0 \leq w_i \leq 1$
- sum of weights: $\sum_{i=1}^p w_{ij} = 1$
- Bounded energy: $\sum_{i=1}^p w_{ij}$

Hence the problem becomes,

$$\underset{w_{ij}}{\text{Minimise}} \sum_{i,j} w_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 + \frac{\lambda}{2} \sum_{j=1}^p w_{ij}^2 \quad (3)$$

$$\text{Subject to } \sum_{j=1}^p w_{ij} = 1 \quad w_{ij} \geq 0 \quad w_{ii} = 0 \quad \forall i, j \quad (4)$$

Since the objective function is convex and the constraints are all linear, this is a convex optimisation problem, we will get a closed form update rule.

When to use negative weights and when they don't make sense
They make sense when dealing with biology or chemistry of molecules, they don't make sense in real life scenarios.

3.2 How to approach an optimisation problem?

For the above optimisation problem, let us construct a matrix \mathbf{e} s.t. $e_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$ where \mathbf{e}_i as a vector with e_{ij} as the j^{th} element of the vector. This transforms the problem into:

$$\underset{\mathbf{w}_i}{\text{Minimise}} \frac{1}{2} \|\mathbf{w}_i + \frac{\mathbf{e}_i}{\lambda}\|_2^2 \quad (5)$$

$$\text{subject to } \mathbf{w}_i^T \mathbf{1} = 1, w_{ij} \geq 0, w_{ii} = 0 \quad (6)$$

Lagrangian function for this problem is given by the following expression,

$$\mathcal{L}(\mathbf{w}_i, \eta_i, \beta_i) = \frac{1}{2} \|\mathbf{w}_i + \frac{\mathbf{e}_i}{\lambda}\|_2^2 - \eta_i(\mathbf{w}_i^T \cdot \mathbf{1} - 1) - \beta_i^T \cdot \mathbf{w}_i \quad (7)$$

Where $\eta_i > 0, \beta_i \in \mathbf{R}^p$. Writing the KKT conditions,

1. Derivative of Lagrangian w.r.t. \mathbf{w}_i . so we have

$$\hat{\mathbf{w}}_i + \frac{\mathbf{e}_i}{\lambda} - \eta_i \mathbf{1} - \beta_i = 0 \quad (8)$$

2. Setting derivative of the j^{th} element of w_i , we have

$$\hat{w}_{ij} + \frac{e_{ij}}{\lambda} - \eta_i - \beta_{ij} = 0 \quad (9)$$

3. $\hat{w}_{ij}\beta_{ij} = 0$

We get,

$$\hat{w}_{ij} = (-\frac{e_{ij}}{\lambda} + \eta_i)^+ \quad (10)$$

where $a^+ = \max(0, a)$.

3.3 How to enforce sparsity?

Uptill now, we have got a graph that minimises the regularised errors subject to some constraints on the weights of the edges, but what if we want to enforce some sparsity conditions? This can be achieved in two ways,

1. Use l_1 norm.
2. choose λ s.t. each node has exactly $m (<< p)$ neighbours i.e. $\|\mathbf{w}_i\|_0 = m$

Exploring the second path, let us assume $e_{i1}, e_{i2} \dots e_{in}$ are arranged in increasing order.

Imposing the constraint $\|\mathbf{w}_i\|_0 = m$, we get $\hat{w}_{im} \geq 0$ and $\hat{w}_{im+1} = 0$, hence we have ,

$$-\frac{e_{im}}{\lambda_i} + 2\eta_i > 0 \text{ and } -\frac{e_{i(m+1)}}{\lambda_i} + \eta_i \leq 0 \quad (11)$$

combining the solution of \hat{w} with the constraint $\hat{\mathbf{w}} \cdot \mathbf{1} = 1$, we have ,

$$\sum_{j=1}^m \left(-\frac{e_{im}}{\lambda_i} + 2\eta_i \right) = 1 \implies \eta_i = \frac{1}{m} + \frac{1}{2m\lambda_i} \sum_{j=1}^m m e_{ij} \quad (12)$$

This leads to the following inequality on λ ,

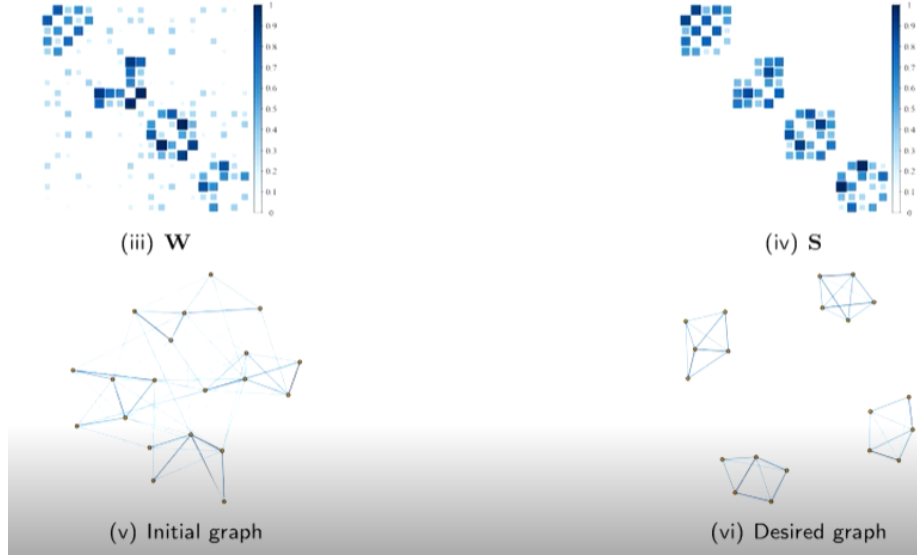
$$\frac{m}{2} e_{im} - \frac{1}{2} \sum_{j=1}^m e_{ij} < \lambda_i \leq \frac{m}{2} e_{i(m+1)} - \frac{1}{2} \sum_{j=1}^m e_{ij} \quad (13)$$

Therefore to obtain an optimal solution with exactly m non-zero values ($\|\hat{\mathbf{w}}_i\|_0 = m$), the maximal λ_i is,

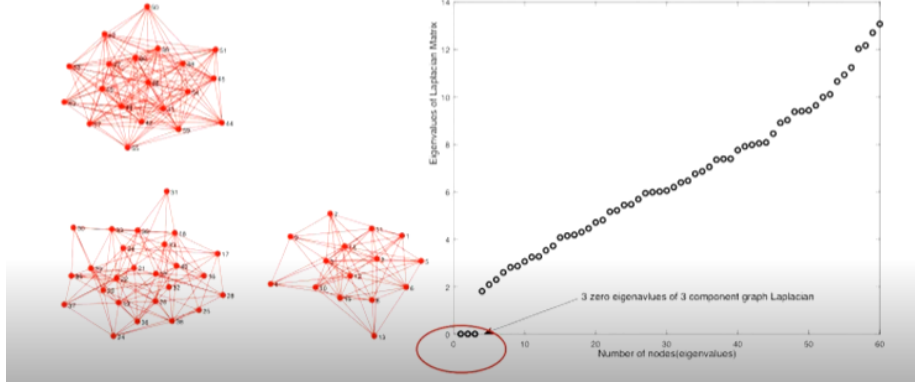
$$\lambda_i = \frac{m}{2} e_{i(m+1)} - \frac{1}{2} \sum_{j=1}^m e_{ij} \quad (14)$$

4 Graph Clustering

Our aim here is to learn a clustered graph from a raw graph. As shown in the figure below, on the left is the unclustered raw graph and on right is the clustered graph.



We can use the spectral properties of the graph. One property that can be used here is that the number of eigenvalues of the Laplacian is equal to the number of components present in the graph. As can be seen below, there are three components and there are three zero eigen values shown in the graph on the left.



4.1 Constrained Laplacian Rank for Clustering

In this algorithm, we first construct a weight matrix, where $w_{ij} = \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2})$. This might not be a clustered graph matrix, but we can learn a clustered graph matrix with the help of this weight matrix. Let us define a matrix L_s as $(\text{Diag}(\mathbf{S}\mathbf{1}) - \mathbf{S})$, hence the problem formulation is as follows:

$$\underset{\mathbf{S}=[\mathbf{s}_1, \dots, \mathbf{s}_p]}{\text{Minimise}} \|\mathbf{S} - \mathbf{W}\|_F^2 \quad (15)$$

$$\text{subject to } \mathbf{s}_i^T \mathbf{1} = 1, s_i \geq 0, s_{ii} = 0, \text{Rank}(\mathbf{L}_s) = p - k$$

4.2 Kay-Fan Theorem

To solve such optimisation problems, Kay Fan theorem says that the sum of top k eigen values is the minimum value of the trace of $\mathbf{F}^T \mathbf{L}_s \mathbf{F}$ subject to $\mathbf{F}^T \mathbf{F} = \mathbf{I}$. Formally stating,

$$\underset{\{\lambda_i(L_s)\}_{i=1}^k}{\text{minimise}} \sum_{i=1}^k \lambda_i(L_s) = \underset{\mathbf{F}}{\text{minimise}} \text{tr}(\mathbf{F}^T \mathbf{L}_s \mathbf{F}) \text{ s.t. } \mathbf{F}^T \mathbf{F} = \mathbf{I} \quad (16)$$

Using the Kay Fan theorem, we can force the first k eigen values to be zero using the followinf formulation,

$$\underset{\mathbf{S}, \mathbf{F}}{\text{minimise}} \|\mathbf{S} - \mathbf{W}\|_2^2 + \beta \text{tr}(\mathbf{F}^T \mathbf{L}_s \mathbf{F}) \quad (17)$$

$$\text{subject to } \mathbf{s}_i^T \mathbf{1} = 1, \mathbf{s}_{ij} \geq 0, \mathbf{s}_{ii} = 0, \mathbf{F}^T \mathbf{F} = \mathbf{I} \quad (18)$$

4.3 Solving for S and F

To solve for S and F, we alternatively solve two subproblems, which are as follows,

1. For S ,

$$\underset{s_{ij}}{\text{minimise}} \sum_j \|s_{ij} - w_{ij}\|_2^2 + \beta \sum_{i,j} \|f_i - f_j\|_2^2 s_{ij} \quad (19)$$

$$\text{subject to } \mathbf{s}_i^T \mathbf{1} = 1, \mathbf{s}_{ij} \geq 0, \mathbf{s}_{ii} = 0, \quad (20)$$

2. For F

$$\underset{F}{\text{minimise}} \text{tr}(\mathbf{F}^T \mathbf{L}_s \mathbf{F}) \text{ s.t. } \mathbf{F}^T \mathbf{F} = \mathbf{I} \quad (21)$$

since the two problems are independent, we can solve them independently.

5 Spectral clustering Algorithm

This algorithm was devised by Ng. et Al. Here they present a spectral clustering based algorithm, that is,

Given a set of points, $S = \{s_1, \dots, s_n\}$, we want to cluster them in k components,

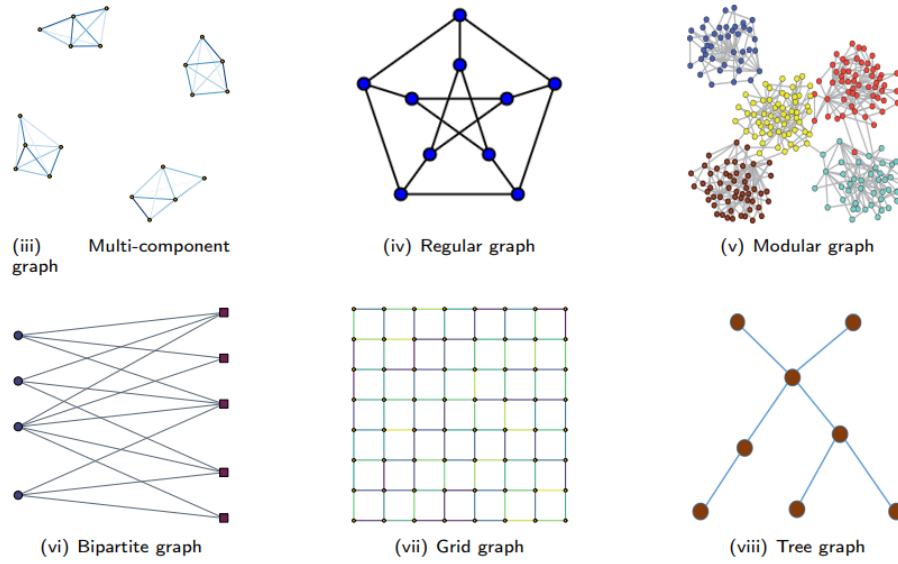
1. Form a matrix A s.t. $A_{ij} = \exp(-\|s_i - s_j\|_2^2 / 2\sigma^2)$ if $i \neq j$ and 0 if $i = j$.
2. Define a diagonal matrix D and construct the matrix $L = D^{-1/2} A D^{-1/2}$
3. Find x_1, \dots, x_k , the k largest Eigen values of L matrix and make a vector $X = [x_1, x_2, \dots, x_k]^T$ by stacking these k values.
4. Normalise X
5. Treating each row of X as a point in \mathbf{R}^k , cluster them into k components by K -means.
6. Assign each s_i to cluster j if and only if row i of the matrix X was assigned to cluster j .

6 Structure Learning

In many scenarios, knowing the structure provides us useful insights and helps us learn better graphs, in this section we explore how the structural properties can help us provide constraints for a combinatorial problem and solve it with help of spectral properties.

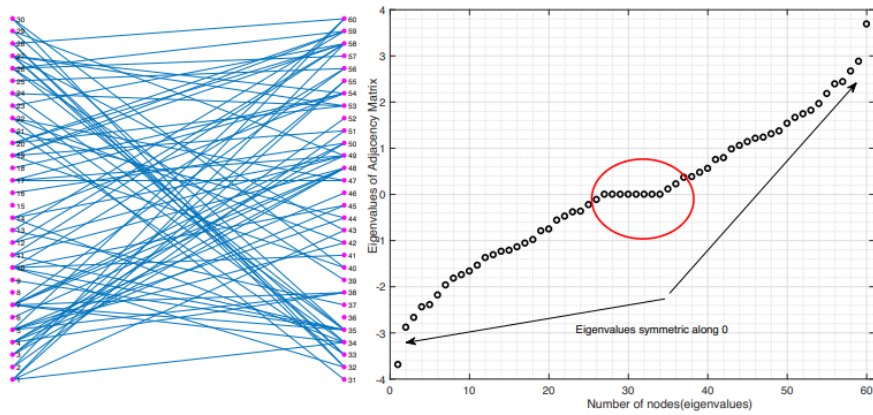
6.1 What is structured graph learning?

In many cases, we already know the structure of graph to be constructed for example, a multi component graph is required when we wish to do clustering. The figure below shows different types of graph structures.



6.2 Spectral properties of structured graphs

Learning from the structure of a graph can be highly combinatorial in nature and hence is NP-Hard, but many properties of graphs are stored in their eigen values. For example, in any bipartite graph the adjacency matrix would have symmetric zero centered Eigenvalues.



can be seen above, the eigenvalues are symmetric about zero. This is a very important combinatorial property encoded in the spectral properties of the graph.

6.3 Incorporating Spectral Properties into graph learning framework

These properties of structured graphs can be used as constraints in the graph learning framework. Hence the problem formulation becomes,

Given $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n$, we compute the sample co variance matrix $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}$. We want to estimate Θ which satisfies the laplacian constraints.

$$\mathbf{S}_\Theta = \{\Theta \in \mathbf{S}_+^p | \Theta_{ij} \leq 0 \ i \neq j \ \Theta \cdot \mathbf{1} = 0\} \quad (22)$$

And the eigenvalues $\lambda(\Theta)$ which satisfy the spectral properties.

6.4 Proposed Unified framework for Graph structure Learning

The below proposed formulation has converted the combinatorial constraints into analytical constraints.

$$\underset{\Theta}{\text{maximise}} \log gdet(\Theta) - tr(\Theta \mathbf{S}) - \alpha h(\Theta)$$

$$\text{subject to } \Theta \in S_\Theta, \lambda(T(\Theta)) \in S_T$$

Here, we are forcing Θ to belong to the set that satisfies the laplacian constraints and λ to belong to the set satisfying the spectral properties, but solving this problem is not easy, hence three formulations have been suggested below.

6.4.1 Formulation 1: Structure via Laplacian spectral constraints

$$\underset{\Theta}{\text{maximise}} \log(gdet(\Theta)) - tr(\Theta \mathbf{S}) - \alpha h(\Theta)$$

$$\text{subject to } \Theta \in S_\Theta, \Theta = U \text{Diag}(\lambda) U^T, \lambda(T(\Theta)) \in S_T, U^T U = I$$

where $\lambda = [\lambda_1 \dots \lambda_p]$ and U is the matrix of eigenvectors. This problem formulation can solve various graph structures depending upon the $h(\Theta)$ function. Some of those graphs are,

1. Connected Sparse Graphs with $S_\lambda = [\lambda_1 = 0, \lambda_2 \geq \lambda_1 \dots \lambda_n \geq \lambda_{n-1}]$
2. k-component graph with $S_\lambda = [\{\lambda\}_1^k = 0, \lambda_{k+1} \geq \lambda_k \dots \lambda_n \geq \lambda_{n-1}]$
3. d-regular graph $S_\lambda = [\{\lambda\}_1^k = 0, c_1 \leq \lambda_{k+1} \leq \lambda_k \dots \lambda_n \leq c_2]$ and $\text{Diag}(\Theta) = d.I$
4. Connected unweighted graph, cycle graph, star graph, minimum edge weighted tree, hypercube graph, maximum spanning tree, etc.

6.4.2 Formulation 2: Structure via Adjacency spectral constraints

We define an operator A that maps the laplacian matrix Θ to the corresponding adjacency matrix Θ_A .

$$[A(\Theta)]_{ij} = -[\Theta_{ij}] \text{ if } i \neq j, 0 \text{ if } i = j$$

$$\underset{\Theta, \psi, V}{\text{maximise}} \log gdet(\Theta) - tr(\Theta S) - \alpha h(\Theta)$$

$$\text{subject to } \Theta \in S_\Theta, A(\Theta) = V \text{Diag}(\psi) V^T, \psi \in S_\psi, V^T V = I$$

$S_\psi = \{\psi_i = -\psi_{p-i+1} \forall i = 1 \dots p\}$. The graphs that can be solved using the ψ formulation are,

1. Bipartite Graph
2. Connected Bipartite Graph

6.4.3 Formulation 3: Structure via joint spectral constraints

$$\underset{\Theta, \lambda, \psi, U, V}{\text{maximise}} \log gdet(\Theta) - tr(\Theta S) - \alpha h(\Theta)$$

$$\text{subject to } \Theta \in S_\Theta, A(\Theta) = V \text{Diag}(\psi) V^T, \psi \in S_\psi, V^T V = I,$$

$$\Theta = U \text{Diag}(\lambda) U^T, \lambda(T(\Theta)) \in S_T, U^T U = I$$

The structures that can be solved using this kind of approach are,

1. Regular Graph $\text{Diag}(\lambda) = \text{Diag}(d, \mathbf{1}) - \text{Diag}(\psi)$
2. k-component bipartite graph $S_\lambda = [\{\lambda\}_1^k = 0, \lambda_{k+1} \geq \lambda_k \dots \lambda_n \geq \lambda_{n-1}]$ and $S_\psi = \{\psi_i = -\psi_{p-i+1} \forall i = 1 \dots p\}$
3. Bipartite Regular Graph $S_\lambda = [\{\lambda\}_1^k = 0, c_1 \leq \lambda_{k+1} \leq \lambda_k \dots \lambda_n \leq c_2]$ and $S_\psi = \{\psi_i = -\psi_{p-i+1} \forall i = 1 \dots p\}$

7 References

1. On Spectral Clustering: Analysis and an algorithm by Andrew Ng, Michael Jordan, Yair Weiss, <https://ai.stanford.edu/~ang/papers/nips01-spectral.pdf>
2. Lecture 31 taught by Prof Sandeep Kumar
3. A Unified Framework for Structured Graph Learning via Spectral Constraints With: J. Ying, Ze. Vinicius, and D. Palomar, Journal of Machine Learning Research (JMLR), Vol. 21, no. 22, pp. 1-60, January 2020