

FINAL PROJECT FOR THE COURSE ADVANCED FOUNDATIONS OF MACHINE LEARNING: ONLINE LEARNING OVER GRAPHS*

PAUL SOPHER LINTILHAC[†] AND THOMAS NANFENG LI[‡]

MAY 9TH, 2017

1 Introduction to Graph

1.1 Concepts and Definitions

We first summarise some key concepts of graph theory, for more detailed knowledge, we refer to Bapat (2014). A **simple graph**, that is, graph without loops and parallel edges, $G(V, E)$ consists of a finite set of **vertices** $V(G)$ and a set of **edges** $E(G)$ consisting of distinct, unordered pairs of vertices. $V(G) = \{v_1, v_2, \dots, v_n\}$ is called the vertex set with $n = |V(G)|$, $E(G) = \{e_{ij}\}$ is called the edge set with $m = |E(G)|$. An edge e_{ij} connects vertices v_i and v_j if they are **adjacent** or neighbours, which is denoted by $v_i \sim v_j$. The number of neighbours of a vertex v is called the **degree** of v and is denoted by $d(v)$, therefore, for each vertex, $d(v_i) = \sum_{v_i \sim v_j} 1$. If all the vertices of a graph have the same degree, the graph is **regular**, the vertices of an **Eulerian Graph** have even degree. A graph is **complete** if there is an edge between every pair of vertices.

$H(G)$ is a **sub-graph** of G if $V(H) \subseteq V(G)$ and $E(H) \subseteq E(G)$. A sub-graph $H(G)$ is an **induced sub-graph** of G if two vertices of $V(H)$ are adjacent if and only if they are adjacent in G . A **clique** is a complete sub-graph of a graph. A **path** of k vertices is a sequence of k distinct vertices such that consecutive vertices are adjacent. A **cycle** is a connected sub-graph where every vertex has exactly two neighbours. A graph containing no cycles is a **forest**. A connected forest is a **tree**.

We define incidence matrix of graph. Let $G(V, E)$ be a graph with $V(G) = \{v_1, v_2, \dots, v_n\}$ and $E(G) = \{e_1, e_2, \dots, e_m\}$. Suppose each edge of $G(V, E)$ is assigned an orientation, which is arbitrary but fixed. The vertex-edge **incidence matrix** of $G(V, E)$, denoted by $Q(G)$, is the $n \times m$ matrix defined as follows. The rows and the columns of $Q(G)$ are indexed by $V(G)$ and $E(G)$, respectively. The (i, j) entry of $Q(G)$ is 0 if vertex i and edge e_j are not incident, and otherwise it is -1 or 1 according as e_j originates or terminates at i , respectively. For instance, the incidence matrix $Q(G)$ of the graph that is shown in figure 1.1 is

$$Q(G) = \begin{bmatrix} 1 & -1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & -1 & 1 & -1 \end{bmatrix}. \quad (1.1)$$

We introduce adjacency matrix of graph. Let $G(V, E)$ be a graph with $V(G) = \{v_1, v_2, \dots, v_n\}$ and $E(G) = \{e_1, e_2, \dots, e_m\}$. The **adjacency matrix** of $G(V, E)$, denoted by $A(G)$, is the $n \times n$ matrix defined as follows. The rows and the columns of $A(G)$ are indexed by $V(G)$. If $i \neq j$ then the (i, j) entry of $A(G)$ is 0 for vertices i and j non-adjacent, and the (i, j) entry is 1 for i and j adjacent. The (i, j) entry of $A(G)$ is 0 for $i = j = 1, \dots, n$. For instance, the adjacency matrix $A(G)$ of the graph that is shown in

*New York University Courant Institute of Mathematical Sciences. Spring 2017. Professor Mehryar Mohri, Ph.D.

[†]New York University School of Engineering. psl274@nyu.edu

[‡]New York University School of Engineering. nl747@nyu.edu

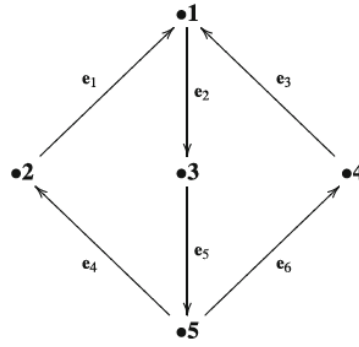


Figure 1.1: Example of Incidence Matrix of Graph

figure 1.2 is

$$A(G) = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}. \quad (1.2)$$

Clearly A is a symmetric matrix with zeros on the diagonal. The (i, j) entry of A^k is the number of walks of length k from i to j .

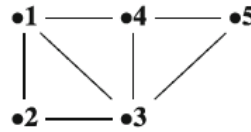


Figure 1.2: Example of Adjacency Matrix of Graph

We define degree matrix of graph. Let $G(V, E)$ be a graph with $V(G) = \{v_1, v_2, \dots, v_n\}$ and $E(G) = \{e_1, e_2, \dots, e_m\}$. The **degree matrix** $D(G)$ for $G(V, E)$ is a $n \times n$ diagonal matrix defined as

$$D(G)_{i,j} := \begin{cases} d(v_i) & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

According to this definition, the degree matrix of figure 1.2 is

$$A(G) = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}. \quad (1.3)$$

Weighted graph $G(V, E, W)$ is a graph with real edge weights given by $w : E \rightarrow \mathbb{R}$. Here, the weight $w(e)$ of an edge e indicates the similarity of the incident vertices, and a missing edge corresponds to zero similarity. The **weighted adjacency matrix** $W(G)$ of the graph $G(V, E, W)$ is defined by

$$W_{ij} := \begin{cases} w(e) & \text{if } e = (i, j) \in E \\ 0 & \text{otherwise.} \end{cases} \quad (1.4)$$

The weight matrix $\mathbf{W}(G)$ can be, for instance, the k -nearest neighbour matrix $\mathbf{W}(G)_{ij} = 1$ if and only if vertex v_i is among the k -nearest neighbours of v_j or vice versa, and is 0 otherwise. Another typical weight matrix is given by the Gaussian kernel of width σ

$$\mathbf{W}(G)_{ij} = e^{-\frac{\|v_i - v_j\|^2}{2\sigma^2}}. \quad (1.5)$$

Then the **degree matrix** $\mathbf{D}(G)$ for weighted graph $G(V, E, W)$ is defined by

$$\mathbf{D}(G)_{i,i} := \sum_j \mathbf{W}(G)_{ij} \quad (1.6)$$

1.2 Graph Laplacian

The graph Laplacian $\mathbf{L}(G)$ is defined in two different ways. The **normalized graph Laplacian** is

$$\mathbf{L}(G) := \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{W} \mathbf{D}^{-\frac{1}{2}}, \quad (1.7)$$

and the **unnormalized graph Laplacian** is

$$\mathbf{L}(G) := \mathbf{D} - \mathbf{W}. \quad (1.8)$$

Let us consider an example to understand the graph Laplacian of the graph that is shown in figure 1.3. Suppose $\mathbf{f} : V \rightarrow \mathbb{R}$ is a real-valued function on the set of the vertices of graph $G(V, E)$ such that it

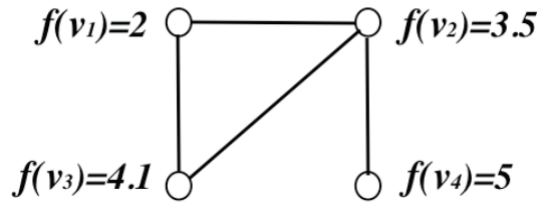


Figure 1.3: Real-Valued Functions on a Graph

assigns a real number to each graph vertex. Therefore, $\mathbf{f} = (f(v_1), f(v_2), \dots, f(v_n))^T \in \mathbb{R}^n$ is a vector indexed by the vertices of graph. Its adjacency matrix is

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}. \quad (1.9)$$

Hence, the eigenvectors of the adjacency matrix, $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$, can be viewed as eigenfunctions $\mathbf{A}\mathbf{f} = \lambda\mathbf{f}$. The adjacency matrix can be viewed as an operator

$$\mathbf{g} = \mathbf{A}\mathbf{f} \quad (1.10)$$

$$g(i) = \sum_{i \sim j} f(j),$$

and it can also be viewed as a quadratic form

$$\mathbf{f}^T \mathbf{A} \mathbf{f} = \sum_{e_{ij}} f(i) f(j). \quad (1.11)$$



Assume that each edge in the graph have an arbitrary but fixed orientation, which is shown in figure 1.4. Then the incidence matrix of the graph is

$$\mathbf{Q} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}. \quad (1.12)$$

Therefore the co-boundary mapping of the graph $\mathbf{f} \rightarrow \mathbf{Q}\mathbf{f}$ implies $(\mathbf{Q}\mathbf{f})(e_{ij}) = f(v_j) - f(v_i)$ is

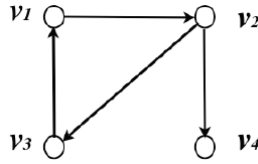


Figure 1.4: Orientation of the Graph

$$\begin{bmatrix} f(2) - f(1) \\ f(1) - f(3) \\ f(3) - f(2) \\ f(4) - f(2) \end{bmatrix} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} f(1) \\ f(2) \\ f(3) \\ f(4) \end{bmatrix}. \quad (1.13)$$

If we let

$$\mathbf{L} = \mathbf{Q}^T \mathbf{Q}, \quad (1.14)$$

then we have

$$(\mathbf{L}\mathbf{f})(v_i) = \sum_{v_i \sim v_j} [f(v_i) - f(v_j)]. \quad (1.15)$$

Hence, the connection between the Laplacian and the adjacency matrices is

$$\mathbf{L} = \mathbf{D} - \mathbf{A} = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}, \quad (1.16)$$

where the degree matrix \mathbf{D} is

$$\mathbf{D} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (1.17)$$

If we consider undirected weighted graphs, which is each edge e_{ij} is weighted by w_{ij} , then the Laplacian as an operator is

$$(\mathbf{L}\mathbf{f})(v_i) = \sum_{v_i \sim v_j} w_{ij} [f(v_i) - f(v_j)]. \quad (1.18)$$

Its quadratic form is

$$\mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{e_{ij}} w_{ij} [f(v_i) - f(v_j)]^2. \quad (1.19)$$

The intuition behind a Laplacian matrix is the following. If, for instance, we apply the Laplacian operator of formula 1.16 to the real-valued functions $\mathbf{f} = (f(v_1), f(v_2), f(v_3), f(v_4))^T$ of the set of the vertices of

graph $G(V, E)$, we have

$$(\mathbf{L}\mathbf{f})(v_i) = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} f(v_1) \\ f(v_2) \\ f(v_3) \\ f(v_4) \end{bmatrix}. \quad (1.20)$$

For simplicity, let us only look at the first element

$$\begin{aligned} (\mathbf{L}\mathbf{f})(v_i)_1 &= \begin{bmatrix} 2 & -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} f(v_1) \\ f(v_2) \\ f(v_3) \\ f(v_4) \end{bmatrix} \\ &= 2f(v_1) - f(v_2) - f(v_3) \\ &= -[f(v_2) - 2f(v_1) + f(v_3)] \\ &= -[f(v_2) - f(v_1) - f(v_1) + f(v_3)] \end{aligned} \quad (1.21)$$

If we label $f(v_1) = f_k$, $f(v_2) = f_{k+1}$, and $f(v_3) = f_{k-1}$, then we have

$$(\mathbf{L}\mathbf{f})(v_i)_1 = -[f_{k+1} - 2f_k + f_{k-1}]. \quad (1.22)$$

We recall that the second order derivative can be approximated by

$$\begin{aligned} f''(x) &= \frac{\frac{f(x+\Delta x) - f(x)}{\Delta x} - \frac{f(x) - f(x-\Delta x)}{\Delta x}}{\Delta x} \\ &= \frac{f(x+\Delta x) - 2f(x) + f(x-\Delta x)}{(\Delta x)^2}. \end{aligned} \quad (1.23)$$

Hence, we observe that the graph Laplacian is the negative numerator of the finite difference approximation of the second derivative. The Laplacian matrix \mathbf{L} is symmetric and positive semi-definite, and it has n non-negative, real-valued eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. The number of 0 eigenvalues of the Laplacian matrix \mathbf{L} is the number of connected components, because each connected component forms a block in the Laplacian matrix that only has edges within itself, and each block is the Laplacian for a small connected component and it has one zero eigenvalue, so the number of zeros is the number of blocks is the number of connected components.

2 Semi-Supervised Learning and Regularization over Graphs

Semi – supervised learning is halfway between supervised and unsupervised learning. In addition to unlabelled data, the algorithm is provided with some supervision information, but not necessarily for all examples. Often, this information will be the targets associated with some of the examples. In this case, the data set $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ can be divided into two parts: the points $\mathbf{X}_l = \{x_1, x_2, \dots, x_l\}$, for which labels $\mathbf{Y}_l = (y_1, y_2, \dots, y_l)$ are provided, and the points $\mathbf{X}_u = \{x_{l+1}, x_{l+2}, \dots, x_{l+u}\}$, the labels of which are not known.

The common denominator of semi-supervised learning algorithms over graphs is that the data are represented by the vertices of a graph, the edges of which are labelled with the pairwise distances of the incident vertices, and a missing edge corresponds to infinite distance. Most graph methods refer to the graph by utilizing the graph Laplacian.



2.1 Label Propagation

Given the graph $G(V, E, W)$, a simple idea for semi-supervised learning label propagation is to propagate labels on the graph. Starting with vertices $1, 2, \dots, l$ labelled l with their known label l or -1 and nodes $l+1, \dots, n$ labelled with 0 , each vertex starts to propagate its label to its neighbours, and the process is repeated until convergence. Zhou et al. (2004) proposed a label propagation algorithm, see algorithm 2.1, that uses graph Laplacian. Estimated labels on both labelled and unlabelled data are denoted by $\hat{\mathbf{Y}} = (\hat{\mathbf{Y}}_l, \hat{\mathbf{Y}}_u)$, where $\hat{\mathbf{Y}}_l$ may be allowed to differ from the given labels $\mathbf{Y}_l = (y_1, y_2, \dots, y_l)$. At each step a vertex i receives a contribution from its neighbours j and an additional small contribution given by its initial value.

Algorithm 2.1 Graph Laplacian Label Propagation Algorithm

Compute weight matrix \mathbf{W} from formula 1.5 for $i \neq j$ and $\mathbf{W}_{ii} = 0$

Compute the diagonal degree matrix \mathbf{D} by $\mathbf{D}_{ii} = \sum_j \mathbf{W}_{ij}$

Compute the normalized graph Laplacian $\mathbf{L} = \mathbf{D}^{-\frac{1}{2}} \mathbf{W} \mathbf{D}^{-\frac{1}{2}}$

Initialize $\hat{\mathbf{Y}}^{(0)} = (y_1, y_2, \dots, y_l, 0, 0, \dots, 0)$

Choose a parameter $\alpha \in [0, 1]$

Iterate

$$\hat{\mathbf{Y}}^{(t+1)} = \alpha \mathbf{L} \hat{\mathbf{Y}}^{(t)} + (1 - \alpha) \hat{\mathbf{Y}}^{(0)}$$

until convergence to $\hat{\mathbf{Y}}^{(\infty)}$

Label point v_i by the sign of $\hat{y}_i^{(\infty)}$

Bengio, Delalleau and Roux (2006) proposed a similar label propagation scheme based on the Jacobi iterative method for linear systems. It uses an additional regularization term ϵ for better numerical stability, which is shown in algorithm 2.2.

Algorithm 2.2 Jacobi Iterative Label Propagation Algorithm

Compute weight matrix \mathbf{W} from formula 1.5 such that $\mathbf{W}_{ii} = 0$

Compute the diagonal degree matrix \mathbf{D} by $\mathbf{D}_{ii} = \sum_j \mathbf{W}_{ij}$

Choose a parameter $\alpha \in (0, 1)$ and a small $\epsilon > 0$

$\mu = \frac{\alpha}{1-\alpha} \in (0, +\infty)$

Compute the diagonal matrix \mathbf{A} by $\mathbf{A}_{ii} = \mathbf{I}_l(i) + \mu \mathbf{D}_{ii} + \mu \epsilon$

Initialize $\hat{\mathbf{Y}}^{(0)} = (y_1, y_2, \dots, y_l, 0, 0, \dots, 0)$

Iterate

$$\hat{\mathbf{Y}}^{(t+1)} = \mathbf{A}^{-1} (\mu \mathbf{W} \hat{\mathbf{Y}}^{(t)} + \hat{\mathbf{Y}}^{(0)})$$

until convergence to $\hat{\mathbf{Y}}^{(\infty)}$

Label point v_i by the sign of $\hat{y}_i^{(\infty)}$

2.2 Minimization of Cost Function

A method that is equivalent to label propagation algorithms is minimising a cost function derived from the graph $G(V, E, W)$. Given a labelling $\hat{\mathbf{Y}} = (\hat{\mathbf{Y}}_l, \hat{\mathbf{Y}}_u)$, consistency with the initial labelling can be measured by

$$\sum_{i=1}^l (\hat{y}_i - y_i)^2 = \|\hat{\mathbf{Y}}_l - \mathbf{Y}_l\|^2. \quad (2.1)$$

By following the smoothness assumption, if two points x_1 and x_2 in a high-density region are close, then so should be the corresponding outputs y_1 and y_2 , we consider a penalty term of the form

$$\begin{aligned} \frac{1}{2} \sum_{i,j=1}^n \mathbf{W}_{ij} (\hat{y}_i - \hat{y}_j)^2 &= \frac{1}{2} \left(2 \sum_{i=1}^n \hat{y}_i^2 \sum_{j=1}^n \mathbf{W}_{ij} - 2 \sum_{i,j=1}^n \mathbf{W}_{ij} \hat{y}_i \hat{y}_j \right) \\ &= \hat{\mathbf{Y}}^T (\mathbf{D} - \mathbf{W}) \hat{\mathbf{Y}} \\ &= \hat{\mathbf{Y}}^T \mathbf{L} \hat{\mathbf{Y}}. \end{aligned} \quad (2.2)$$

This means we penalize rapid changes in $\hat{\mathbf{Y}}$ between points that are close, which is given by the similarity matrix \mathbf{W} . However, if there is noise in the available labels, it may be beneficial to allow the algorithm to relabel the labelled data. This could also help generalization in a noise-free setting where, for instance, a positive sample had been drawn from a region of space mainly filled with negative samples.

Based on this observation, Belkin, Matveeva and Niyogi (2004) proposed a more general cost criterion involving a trade-off between 2.1 and 2.2, which is shown in algorithm 2.3 and algorithm 2.4. The paper assumed that $G(V, E, W)$ is connected and that the vertices of the graph are numbered, and only partial information, $f(\mathbf{x}_i) = y_i$, $1 \leq i \leq k$, is given. The labels can potentially be noisy. They also allow data points to have multiplicities, i.e. each vertex of the graph may appear more than once with same or different y value. They precondition the data by mean subtracting first. That is we take

$$\tilde{\mathbf{y}} = (y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_k - \bar{y}), \quad (2.3)$$

where $\bar{y} = \frac{1}{k} \sum y_i$.

In the algorithm 2.3, \mathbf{S} is a smoothness matrix, e.g. $S = L$ or $S = L^p$, $p \in \mathbb{N}$. The condition $\sum f_i = 0$ is needed to make the algorithm stable. The solution to the quadratic problem above is not hard to obtain by standard linear algebra considerations. We have the objective function $\frac{1}{k} \sum_i (f_i - \tilde{y}_i)^2 + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f}$ and constraint $\sum f_i = 0$. Therefore, the Lagrangian is

$$\begin{aligned} \mathcal{L} &= \frac{1}{k} \sum_i (f_i - \tilde{y}_i)^2 + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f} - \mu \left(\sum f_i - 0 \right) \\ &= \frac{1}{k} \mathbf{f}^T \mathbf{f} - \frac{2}{k} \mathbf{f}^T \tilde{\mathbf{y}} + \frac{1}{k} \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f} - \mu \mathbf{f}^T \mathbf{1} \\ &= \mathbf{f}^T \left(\frac{1}{k} \mathbf{I} + \gamma \mathbf{S} \right) \mathbf{f} - \frac{2}{k} \mathbf{f}^T (\tilde{\mathbf{y}} + \mu \mathbf{1}) + \frac{1}{k} \tilde{\mathbf{y}}^T \tilde{\mathbf{y}}. \end{aligned} \quad (2.4)$$

Then by taking $\frac{\partial \mathcal{L}}{\partial \mathbf{f}} = 0$, we have

$$\tilde{\mathbf{f}} = (\mathbf{I} + k\gamma \mathbf{S})^{-1} (\tilde{\mathbf{y}} + \mu \mathbf{1}). \quad (2.5)$$

Here, μ is chosen so that the resulting vector \mathbf{f} is orthogonal to $\mathbf{1}$. Denote by $s(\mathbf{f})$ the functional

$$s : \mathbf{f} \rightarrow \sum_i f_i. \quad (2.6)$$

Since s is linear, we obtain

$$0 = s(\tilde{\mathbf{f}}) = s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \tilde{\mathbf{y}}) + s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \mathbf{1}). \quad (2.7)$$

Therefore we can write

$$\mu = - \frac{s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \tilde{\mathbf{y}})}{s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \mathbf{1})} \quad (2.8)$$

Algorithm 2.3 Tichonov Regularization with Parameter $\gamma \in \mathbb{R}$

The objective is to minimize the square loss function plus the smoothness penalty.

$$\tilde{\mathbf{f}} = \arg \min_{\substack{\mathbf{f}=(f_1, \dots, f_n) \\ \sum f_i=0}} \frac{1}{k} \sum_i (f_i - \tilde{y}_i)^2 + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f}$$

In the algorithm 2.4, \mathbf{S} is a smoothness matrix, e.g. $S = L$ or $S = L^p$, $p \in \mathbb{N}$. However, here we are not allowing multiple vertices in the sample. We partition \mathbf{S} as

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_1 & \mathbf{S}_2 \\ \mathbf{S}_2^T & \mathbf{S}_3 \end{bmatrix}, \quad (2.9)$$

where \mathbf{S}_1 is a $k \times k$ matrix, \mathbf{S}_2 is a $k \times (n - k)$ matrix, and \mathbf{S}_3 is a $(n - k) \times (n - k)$ matrix. Let $\tilde{\mathbf{f}}$ be the values of \mathbf{f} , where the function is unknown, $\tilde{\mathbf{f}} = (f_{k+1}, \dots, f_n)$. By applying the similar Lagrangian



multiplier method, we have

$$\tilde{\mathbf{f}} = \mathbf{S}_3^{-1} \mathbf{S}_2^T \left((\tilde{y}_1, \dots, \tilde{y}_k)^T + \mu \mathbf{1} \right), \quad (2.10)$$

and

$$\mu = -\frac{s(\mathbf{S}_3^{-1} \mathbf{S}_2^T \tilde{\mathbf{y}})}{s(\mathbf{S}_3^{-1} \mathbf{S}_2^T \mathbf{1})} \quad (2.11)$$

Algorithm 2.4 Interpolated Regularization without Parameter

Here we assume that the values y_1, y_2, \dots, y_k have no noise.

Thus the optimization problem is to find a function of maximum smoothness satisfying $f(\mathbf{x}_i) = \tilde{y}_i, 1 \leq i \leq k$
 $\tilde{\mathbf{f}} = \arg \min_{\mathbf{f}=(\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_k, f_{k+1}, \dots, f_n)} \mathbf{f}^T \mathbf{S} \mathbf{f}$
 $\sum_{f_i=0}$

We now investigate generalization bounds for this graph semi-supervised learning regularization through algorithmic stability. Belkin, Matveeva and Niyogi (2004) defined the empirical error $R_k(f)$, which is a measure of how well we do on the training set, and the generalization error $R(f)$, which is the expectation of how well we do on all labelled or unlabelled points, in the following way

$$R_k(f) = \frac{1}{k} \sum_{i=1}^k (f(\mathbf{x}_i) - y_i)^2 \quad (2.12)$$

$$R(f) = \mathbb{E}_\mu [(f(\mathbf{x}) - y(\mathbf{x}))^2],$$

where f_T maps a given set of examples T to \mathbb{R} , i.e. $f_T : V \rightarrow \mathbb{R}$, the expectation is taken over an underlying distribution μ . In theorem 5, they showed that for data samples of size $k \geq 4$ with multiplicity of at most t , γ -regularization using the smoothness functional S is a $\left(\frac{3M\sqrt{tk}}{(k\gamma\lambda_1 - t)^2} + \frac{4M}{k\gamma\lambda_1 - t} \right)$ -stable algorithm, assuming that the denominator $k\gamma\lambda_1 - t$ is positive. Bousquet and Elisseeff (2001) showed that for a β -stable algorithm $T \rightarrow f_T$ we have

$$\mathbb{P}(|R_k(f_T) - R(f_T)| > \epsilon + \beta) \leq 2 \exp \left(-\frac{k\epsilon^2}{2(k\beta + K + M)^2} \right), \quad \forall \epsilon > 0. \quad (2.13)$$

Therefore, by following the derivation in the theorem 11.1 in Mohri, Rostamizadeh and Talwalkar (2012), we have with probability $1 - \delta$

$$|R_k(f_T) - R(f_T)| < \epsilon + \beta \quad (2.14)$$

where $\delta = 2 \exp \left(-\frac{k\epsilon^2}{2(k\beta + K + M)^2} \right)$. We then solve for ϵ in this expression for δ and get

$$\epsilon = \sqrt{\frac{2 \ln \frac{2}{\delta}}{k}} (k\beta + K + M), \quad (2.15)$$

then plug into inequality 2.14 and rearrange terms, then, with probability $1 - \delta$, we have

$$|R_k(f_T) - R(f_T)| \leq \beta + \sqrt{\frac{2 \ln \frac{2}{\delta}}{k}} (k\beta + K + M), \quad (2.16)$$

where

$$\beta = \left(\frac{3M\sqrt{tk}}{(k\gamma\lambda_1 - t)^2} + \frac{4M}{k\gamma\lambda_1 - t} \right). \quad (2.17)$$

This is the theorem 1 of Belkin, Matveeva and Niyogi (2004).

3 Projection Algorithm for Graph Online Learning

A different approach was proposed by Herbster, Pontil and Wainer (2005) using projection and kernel methods. The method uses minimal norm interpolation, which I will later explain can be thought of as finding

the distribution over the vertices that minimizes the total flux induced over the graph, under the constraint that the labelled points are approximately correct.

3.1 Pseudo-Inverse of Graph Laplacian and Kernel of Hilbert

A key object in this approach is the pseudo-inverse of the graph Laplacian \mathbf{L}^+ can be thought of as the reproducing kernel of a particular Hilbert space H of real-valued functions over the vertices of the graph

$$\mathbf{f} : V \rightarrow \mathbb{R}^n \quad (3.1)$$

equipped with the inner product

$$\langle \mathbf{f}, \mathbf{g} \rangle = \mathbf{f}^T \mathbf{L} \mathbf{g}. \quad (3.2)$$

Over the entire Vector space of real-valued functions on the graph, $\langle \mathbf{f}, \mathbf{g} \rangle$ is only a semi-inner product, and therefore induces the semi-norm $\|\mathbf{g}\|$. If \mathbf{f}^* is an eigenvector of \mathbf{L} corresponding to an eigenvalue of 0, then $\langle \mathbf{f}^*, \mathbf{f}^* \rangle = 0$. These eigenvectors with eigenvalues of zero are precisely the eigenvectors that are piecewise constant. We can see this by noting that

$$\|\mathbf{g}\|^2 = \sum_{(i,j) \in E(G)} (\mathbf{g}_i - \mathbf{g}_j)^2. \quad (3.3)$$

This is because

$$\begin{aligned} \|\mathbf{g}\|^2 &= \langle \mathbf{g}, \mathbf{g} \rangle \\ &= \sum_{i,j=1}^n \mathbf{g}_i \mathbf{L}_{ij} \mathbf{g}_j \\ &= \sum_{i=j}^n \mathbf{g}_i \mathbf{L}_{ij} \mathbf{g}_j + \sum_{i \neq j}^n \mathbf{g}_i \mathbf{L}_{ij} \mathbf{g}_j \\ &= \sum_{i=1}^n \mathbf{g}_i^2 \mathbf{D}_{ii} + \sum_{i \neq j}^n \mathbf{g}_i \mathbf{A}_{ij} \mathbf{g}_j \\ &= \sum_{i,j=1}^n \mathbf{g}_i^2 \mathbf{A}_{ij} + \sum_{i \neq j}^n \mathbf{g}_i \mathbf{L}_{ij} \mathbf{g}_j \\ &= \sum_{i,j \in E(G)} \mathbf{g}_i^2 - 2 \mathbf{g}_i \mathbf{g}_j \\ &= \sum_{(i,j) \in E(G)} (\mathbf{g}_i - \mathbf{g}_j)^2. \end{aligned} \quad (3.4)$$

Please notice that the diagonal component of \mathbf{L} is \mathbf{D} , the off-diagonal component of \mathbf{L} is $-\mathbf{A}$, $\mathbf{D}_{ii} = \sum_{j=1}^n \mathbf{A}_{ij}$, and in the proof we pick up a factor of 2 from the second summation since E is the set of unordered pairs in the edge set, and therefore each edge gets counted twice.

As we can see, the sum vanishes if $\mathbf{g}_i = \mathbf{g}_j$ for all vertices connected by an edge, i.e. when \mathbf{g} is piecewise constant on each connected component. Therefore we must restrict our functions to be in the subspace of H spanned by the eigenvectors that have non-zero eigenvalues. Restricted to this subspace, $\langle \mathbf{f}, \mathbf{g} \rangle$ is indeed an inner product, and therefore induces a true norm $\|\mathbf{g}\|$. Since H is a Hilbert space, we can use the Riesz Representation theorem to conclude that the evaluation functional of \mathbf{f} at any vertex, $\mathbf{f}(v_i)$, being a linear functional $(\mathbf{f} + \mathbf{g})(v_i) = (\mathbf{f} + \mathbf{g})_i(v) = \mathbf{f}_i(v) + \mathbf{g}_i(v) = \mathbf{f}(v_i) + \mathbf{g}(v_i)$, can be represented as an inner product

$$\mathbf{f}(v_i) = \langle \mathbf{K}_i, \mathbf{f} \rangle = \mathbf{K}_i \mathbf{L} \mathbf{f}. \quad (3.5)$$

The pseudo-inverse \mathbf{L}^+ satisfies this property, which can be verified by noting that

$$\mathbf{f} = \mathbf{f}_i = \mathbf{L}_i^+ \mathbf{L} \mathbf{f}_i. \quad (3.6)$$



Thus \mathbf{L}^+ is the reproducing kernel of H .

3.2 Graph Laplacian and Heat Equation

We can get more intuition about the basic objects such as $\|\mathbf{g}\|$ and \mathbf{L}^+ by using the concept of conductance (it is no coincidence that the continuous version of the Laplacian appears in the heat equation). If we think of \mathbf{f} as some distribution over the nodes, we can think of $\mathbf{L}\mathbf{f}$ roughly as the flux induced at each of the nodes by that distribution over the graph. Then the form $\mathbf{g}^T \mathbf{L}\mathbf{f}$ represents some weighted measurement of this flux with the weights given by \mathbf{g} . For example, if \mathbf{f} is a binary vector representing the boundary of some open set on the graph $\langle \mathbf{f}, \mathbf{g} \rangle$ would be the flux produced by \mathbf{g} as measured on the boundary represented by \mathbf{f} . We can even use this formalism to quickly prove a version of Stokes' Theorem on graphs.

Bibliography

- Bapat, Ravindra B. (2014). *Graphs and Matrices*. London, UK: Springer-Verlag. 193 pp. ISBN: 9781447165682 (cit. on p. 1).
- Belkin, Mikhail, Irina Matveeva and Partha Niyogi (2004). *Regularization and Semi-supervised Learning on Large Graphs*. Ed. by John Shawe-Taylor and Yoram Singer. Learning Theory 17th Annual Conference on Learning Theory, COLT 2004, Banff, Canada, July 1-4, 2004. Proceedings. Berlin, Deutschland: Springer Science+Business Media. 654 pp. ISBN: 9783540222828 (cit. on pp. 7, 8).
- Bengio, Yoshua, Olivier Delalleau and Nicolas Le Roux (2006). *Semi-Supervised Learning*. Cambridge, Massachusetts, USA: The MIT Press. 528 pp. ISBN: 9780262033589 (cit. on p. 6).
- Bousquet, Olivier and André Elisseeff (2001). *Algorithmic Stability and Generalization Performance*. Ed. by Todd K. Leen, Thomas G. Dietterich and Volker Tresp. Advances in Neural Information Processing Systems 13 Proceedings of the 2000 Conference. Cambridge, MA, USA: The MIT Press. 1128 pp. ISBN: 9780262122412 (cit. on p. 8).
- Herbster, Mark, Massimiliano Pontil and Lisa Wainer (2005). *Online Learning over Graphs*. Ed. by Luc De Raedt and Stefan Wrobel. ICML '05 Proceedings of the 22nd International Conference on Machine Learning. New York, NY, USA: Association for Computing Machinery. 1113 pp. ISBN: 1595931805 (cit. on p. 8).
- Mohri, Mehryar, Afshin Rostamizadeh and Ameet Talwalkar (2012). *Foundations of Machine Learning*. Cambridge, Massachusetts, USA: The MIT Press. 432 pp. ISBN: 9780262018258 (cit. on p. 8).
- Zhou, Dengyong, Olivier Bousquet, Thomas Navin Lal, Jason Weston and Bernhard Schölkopf (2004). *Learning with Local and Global Consistency*. Ed. by Sebastian Thrun, Lawrence K. Saul and Bernhard Schölkopf. Advances in Neural Information Processing Systems 16 Proceedings of the 2003 Conference. Cambridge, Massachusetts, USA: The MIT Press. 1728 pp. ISBN: 9780262201520 (cit. on p. 6).
- Zhu, Xiaojin and Zoubin Ghahramani (2002). "Learning from Labelled and Unlabelled Data with Label Propagation". In: *Technical Report CMU-CALD-02-107, Carnegie Mellon University*.

